

---

# **datalad Documentation**

*Release 0.12.2*

**DataLad team**

**Feb 22, 2020**



---

## Contents

---

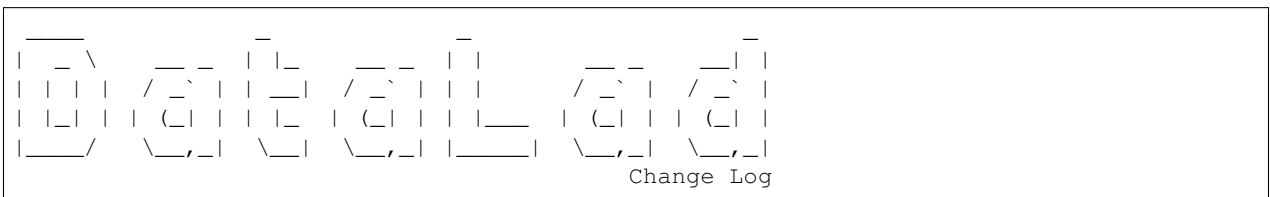
|                            |            |
|----------------------------|------------|
| <b>1 Content</b>           | <b>3</b>   |
| <b>Python Module Index</b> | <b>283</b> |
| <b>Index</b>               | <b>285</b> |



Welcome to DataLad's **technical documentation**. Information here is targeting software developers and is focused on the Python and command line APIs, as well as software design, employed technologies, and key features. Comprehensive **user documentation** with information on installation, basic operation, support, and (advanced) use case descriptions is available in the [DataLad handbook](#).



## 1.1 Change log



This is a high level and scarce summary of the changes between releases. We would recommend to consult log of the [DataLad git repository](#) for more details.

### 1.1.1 0.12.2 (Jan 28, 2020) – Smoothen the ride

Mostly a bugfix release with various robustifications, but also makes the first step towards versioned dataset installation requests.

#### Major refactoring and deprecations

- The minimum required version for GitPython is now 2.1.12. ([#4070](#))

#### Fixes

- The class for handling configuration values, `ConfigManager`, inappropriately considered the current working directory's dataset, if any, for both reading and writing when instantiated with `dataset=None`. This misbehavior is fairly inaccessible through typical use of DataLad. It affects `datalad.cfg`, the top-level configuration instance that should not consider repository-specific values. It also affects Python users that call `Dataset` with a path that does not yet exist and persists until that dataset is created. ([#4078](#))

- `update` saved the dataset when called with `--merge`, which is unnecessary and risks committing unrelated changes. (#3996)
- Confusing and irrelevant information about Python defaults have been dropped from the command-line help. (#4002)
- The logic for automatically propagating the ‘origin’ remote when cloning a local source didn’t properly account for relative paths. (#4045)
- Various fixes to file name handling and quoting on Windows. (#4049) (#4050)
- When cloning failed, error lines were not bubbled up to the user in some scenarios. (#4060)

### Enhancements and new features

- `clone` (and thus `install`)
  - now propagates the `reckless` mode from the superdataset when cloning a dataset into it. (#4037)
  - gained support for `ria+<protocol>://` URLs that point to RIA stores. (#4022)
  - learned to read “@version” from `ria+` URLs and install that version of a dataset (#4036) and to apply URL rewrites configured through Git’s `url.*.insteadOf` mechanism (#4064).
  - now copies `datalad.get.subdataset-source-candidate-<name>` options configured within the superdataset into the subdataset. This is particularly useful for RIA data stores. (#4073)
- Archives are now (optionally) handled with 7-Zip instead of `patool`. 7-Zip will be used by default, but `patool` will be used on non-Windows systems if the `datalad.runtime.use-patool` option is set or the 7z executable is not found. (#4041)

### 1.1.2 0.12.1 (Jan 15, 2020) – Small bump after big bang

Fix some fallout after major release.

#### Fixes

- Revert incorrect relative path adjustment to URLs in `clone`. (#3538)
- Various small fixes to internal helpers and test to run on Windows (#2566) (#2534)

### 1.1.3 0.12.0 (Jan 11, 2020) – Krakatoa

This release is the result of more than a year of development that includes fixes for a large number of issues, yielding more robust behavior across a wider range of use cases, and introduces major changes in API and behavior. It is the first release for which extensive user documentation is available in a dedicated [DataLad Handbook](#). Python 3 (3.5 and later) is now the only supported Python flavor.

#### Major changes 0.12 vs 0.11

- `save` fully replaces `add` (which is obsolete now, and will be removed in a future release).
- A new Git-annex aware `status` command enables detailed inspection of dataset hierarchies. The previously available `diff` command has been adjusted to match `status` in argument semantics and behavior.



- The ability to configure dataset procedures prior and after the execution of particular commands has been replaced by a flexible “hook” mechanism that is able to run arbitrary DataLad commands whenever command results are detected that match a specification.
- Support of the Windows platform has been improved substantially. While performance and feature coverage on Windows still falls behind Unix-like systems, typical data consumer use cases, and standard dataset operations, such as `create` and `save`, are now working. Basic support for data provenance capture via `run` is also functional.
- Support for Git-annex direct mode repositories has been removed, following the end of support in Git-annex itself.
- The semantics of relative paths in command line arguments have changed. Previously, a call `datalad save --dataset /tmp/myds some/relpath` would have been interpreted as saving a file at `/tmp/myds/some/relpath` into dataset `/tmp/myds`. This has changed to saving `$PWD/some/relpath` into dataset `/tmp/myds`. More generally, relative paths are now always treated as relative to the current working directory, except for path arguments of `Dataset` class instance methods of the Python API. The resulting partial duplication of path specifications between path and dataset arguments is mitigated by the introduction of two special symbols that can be given as dataset argument: `^` and `^.`, which identify the topmost superdataset and the closest dataset that contains the working directory, respectively.
- The concept of a “core API” has been introduced. Commands situated in the module `datalad.core` (such as `create`, `save`, `run`, `status`, `diff`) receive additional scrutiny regarding API and implementation, and are meant to provide longer-term stability. Application developers are encouraged to preferentially build on these commands.

### Major refactoring and deprecations since 0.12.0rc6

- `clone` has been incorporated into the growing core API. The public `--alternative-source` parameter has been removed, and a `clone_dataset` function with multi-source capabilities is provided instead. The `--reckless` parameter can now take literal mode labels instead of just being a binary flag, but backwards compatibility is maintained.
- The `get_file_content` method of `GitRepo` was no longer used internally or in any known DataLad extensions and has been removed. (#3812)
- The function `get_dataset_root` has been replaced by `rev_get_dataset_root`. `rev_get_dataset_root` remains as a compatibility alias and will be removed in a later release. (#3815)
- The `add_sibling` module, marked obsolete in v0.6.0, has been removed. (#3871)
- `mock` is no longer declared as an external dependency because we can rely on it being in the standard library now that our minimum required Python version is 3.5. (#3860)
- `download-url` now requires that directories be indicated with a trailing slash rather than interpreting a path as directory when it doesn't exist. This avoids confusion that can result from typos and makes it possible to support directory targets that do not exist. (#3854)
- The `dataset_only` argument of the `ConfigManager` class is deprecated. Use `source="dataset"` instead. (#3907)
- The `--proc-pre` and `--proc-post` options have been removed, and configuration values for `datalad.COMMAND.proc-pre` and `datalad.COMMAND.proc-post` are no longer honored. The new result hook mechanism provides an alternative for `proc-post` procedures. (#3963)

### Fixes since 0.12.0rc6

- `publish` crashed when called with a detached HEAD. It now aborts with an informative message. (#3804)
- Since 0.12.0rc6 the call to `update` in `siblings` resulted in a spurious warning. (#3877)

- `siblings` crashed if it encountered an annex repository that was marked as dead. (#3892)
- The update of `rerun` in v0.12.0rc3 for the rewritten `diff` command didn't account for a change in the output of `diff`, leading to `rerun --report` unintentionally including unchanged files in its diff values. (#3873)
- In 0.12.0rc5 `download-url` was updated to follow the new path handling logic, but its calls to `AnnexRepo` weren't properly adjusted, resulting in incorrect path handling when the called from a dataset subdirectory. (#3850)
- `download-url` called `git annex addurl` in a way that failed to register a URL when its header didn't report the content size. (#3911)
- With Git v2.24.0, saving new subdatasets failed due to a bug in that Git release. (#3904)
- With DataLad configured to stop on failure (e.g., specifying `--on-failure=stop` from the command line), a failing result record was not rendered. (#3863)
- Installing a subdataset yielded an “ok” status in cases where the repository was not yet in its final state, making it ineffective for a caller to operate on the repository in response to the result. (#3906)
- The internal helper for converting `git-annex`'s JSON output did not relay information from the “error-messages” field. (#3931)
- `run-procedure` reported relative paths that were confusingly not relative to the current directory in some cases. It now always reports absolute paths. (#3959)
- `diff` inappropriately reported files as deleted in some cases when `to` was a value other than `None`. (#3999)
- An assortment of fixes for Windows compatibility. (#3971) (#3974) (#3975) (#3976) (#3979)
- Subdatasets installed from a source given by relative path will now have this relative path used as ‘url’ in their `.gitmodules` record, instead of an absolute path generated by Git. (#3538)
- `clone` will now correctly interpret ‘~/...’ paths as absolute path specifications. (#3958)
- `run-procedure` mistakenly reported a directory as a procedure. (#3793)
- The cleanup for batched `git-annex` processes has been improved. (#3794) (#3851)
- The function for adding a version ID to an AWS S3 URL doesn't support URLs with an “s3://” scheme and raises a `NotImplementedError` exception when it encounters one. The function learned to return a URL untouched if an “s3://” URL comes in with a version ID. (#3842)
- A few spots needed to be adjusted for compatibility with `git-annex`'s new `--sameas` feature, which allows special remotes to share a data store. (#3856)
- The `swallow_logs` utility failed to capture some log messages due to an incompatibility with Python 3.7. (#3935)
- `siblings`
  - crashed if `--inherit` was passed but the parent dataset did not have a remote with a matching name. (#3954)
  - configured the wrong `pushurl` and `annexurl` values in some cases. (#3955)

### Enhancements and new features since 0.12.0rc6

- By default, datasets cloned from local source paths will now get a configured remote for any recursively discoverable ‘origin’ sibling that is also available from a local path in order to maximize automatic file availability across local annexes. (#3926)
- The new `result hooks mechanism` allows callers to specify, via local Git configuration values, DataLad command calls that will be triggered in response to matching result records (i.e., what you see when you call a command with `-f json_pp`). (#3903)

- The command interface classes learned to use a new `_examples_` attribute to render documentation examples for both the Python and command-line API. (#3821)
- Candidate URLs for cloning a submodule can now be generated based on configured templates that have access to various properties of the submodule, including its dataset ID. (#3828)
- DataLad’s check that the user’s Git identity is configured has been sped up and now considers the appropriate environment variables as well. (#3807)
- The `tag` method of `GitRepo` can now tag revisions other than `HEAD` and accepts a list of arbitrary `git tag` options. (#3787)
- When `get` clones a subdataset and the subdataset’s `HEAD` differs from the commit that is registered in the parent, the active branch of the subdataset is moved to the registered commit if the registered commit is an ancestor of the subdataset’s `HEAD` commit. This handling has been moved to a more central location within `GitRepo`, and now applies to any `update_submodule(..., init=True)` call. (#3831)
- The output of `datalad -h` has been reformatted to improve readability. (#3862)
- `unlock` has been sped up. (#3880)
- `run-procedure` learned to provide and render more information about discovered procedures, including whether the procedure is overridden by another procedure with the same base name. (#3960)
- `save now` (#3817)
  - records the active branch in the superdataset when registering a new subdataset.
  - calls `git annex sync` when saving a dataset on an adjusted branch so that the changes are brought into the mainline branch.
- `subdatasets` now aborts when its `dataset` argument points to a non-existent dataset. (#3940)
- `wtf now`
  - reports the dataset ID if the current working directory is visiting a dataset. (#3888)
  - outputs entries deterministically. (#3927)
- The `ConfigManager` class
  - learned to exclude `.datalad/config` as a source of configuration values, restricting the sources to standard Git configuration files, when called with `source="local"`. (#3907)
  - accepts a value of “override” for its `where` argument to allow Python callers to more convenient override configuration. (#3970)
- Commands now accept a `dataset` value of “^.” as shorthand for “the dataset to which the current directory belongs”. (#3242)

### 1.1.4 0.12.0rc6 (Oct 19, 2019) – some releases are better than the others

bet we will fix some bugs and make a world even a better place.

#### Major refactoring and deprecations

- DataLad no longer supports Python 2. The minimum supported version of Python is now 3.5. (#3629)
- Much of the user-focused content at <http://docs.datalad.org> has been removed in favor of more up to date and complete material available in the [DataLad Handbook](#). Going forward, the plan is to restrict <http://docs.datalad.org> to technical documentation geared at developers. (#3678)

- `update` used to allow the caller to specify which dataset(s) to update as a `PATH` argument or via the `--dataset` option; now only the latter is supported. Path arguments only serve to restrict which subdataset are updated when operating recursively. (#3700)
- Result records from a `get` call no longer have a “state” key. (#3746)
- `update` and `get` no longer support operating on independent hierarchies of datasets. (#3700) (#3746)
- The `run` update in 0.12.0rc4 for the new path resolution logic broke the handling of inputs and outputs for calls from a subdirectory. (#3747)
- The `is_submodule_modified` method of `GitRepo` as well as two helper functions in `gitrepo.py`, `kwargs_to_options` and `split_remote_branch`, were no longer used internally or in any known `DataLad` extensions and have been removed. (#3702) (#3704)
- The `only_remote` option of `GitRepo.is_with_annex` was not used internally or in any known extensions and has been dropped. (#3768)
- The `get_tags` method of `GitRepo` used to sort tags by committer date. It now sorts them by the tagger date for annotated tags and the committer date for lightweight tags. (#3715)
- The `rev_resolve_path` substituted `resolve_path` helper. (#3797)

### Fixes

- Correctly handle relative paths in `publish`. (#3799) (#3102)
- Do not erroneously discover directory as a procedure. (#3793)
- Correctly extract version from manpage to trigger use of manpages for `--help`. (#3798)
- The `cfg_yoda` procedure saved all modifications in the repository rather than saving only the files it modified. (#3680)
- Some spots in the documentation that were supposed appear as two hyphen’s were incorrectly rendered in the HTML output en-dash’s. (#3692)
- `create`, `install`, and `clone` treated paths as relative to the dataset even when the string form was given, violating the new path handling rules. (#3749) (#3777) (#3780)
- Providing the “^” shortcut to `--dataset` didn’t work properly when called from a subdirectory of a subdataset. (#3772)
- We failed to propagate some errors from `git-annex` when working with its JSON output. (#3751)
- With the Python API, callers are allowed to pass a string or list of strings as the `cfg_proc` argument to `create`, but the string form was mishandled. (#3761)
- Incorrect command quoting for SSH calls on Windows that rendered basic SSH-related functionality (e.g., `sshrun`) on Windows unusable. (#3688)
- Annex JSON result handling assumed platform-specific paths on Windows instead of the POSIX-style that is happening across all platforms. (#3719)
- `path_is_under()` was incapable of comparing Windows paths with different drive letters. (#3728)

### Enhancements and new features

- Provide a collection of “public” `call_git*` helpers within `GitRepo` and replace use of “private” and less specific `_git_custom_command` calls. (#3791)

- `status` gained a `--report-filetype`. Setting it to “raw” can give a performance boost for the price of no longer distinguishing symlinks that point to annexed content from other symlinks. (#3701)
- `save` disables file type reporting by `status` to improve performance. (#3712)
- `subdatasets` (#3743)
  - now extends its result records with a `contains` field that lists which `contains` arguments matched a given subdataset.
  - yields an ‘impossible’ result record when a `contains` argument wasn’t matched to any of the reported subdatasets.
- `install` now shows more readable output when cloning fails. (#3775)
- `SSHConnection` now displays a more informative error message when it cannot start the `ControlMaster` process. (#3776)
- If the new configuration option `datalad.log.result-level` is set to a single level, all result records will be logged at that level. If you’ve been bothered by DataLad’s double reporting of failures, consider setting this to “debug”. (#3754)
- Configuration values from `datalad -c OPTION=VALUE ...` are now validated to provide better errors. (#3695)
- `rerun` learned how to handle history with merges. As was already the case when cherry picking non-run commits, re-creating merges may result in conflicts, and `rerun` does not yet provide an interface to let the user handle these. (#2754)
- The `fsck` method of `AnnexRepo` has been enhanced to expose more features of the underlying `git fsck` command. (#3693)
- `GitRepo` now has a `for_each_ref_` method that wraps `git for-each-ref`, which is used in various spots that used to rely on `GitPython` functionality. (#3705)
- Do not pretend to be able to work in optimized (`python -O`) mode, crash early with an informative message. (#3803)

### 1.1.5 0.12.0rc5 (September 04, 2019) – .

Various fixes and enhancements that bring the 0.12.0 release closer.

#### Major refactoring and deprecations

- The two modules below have a new home. The old locations still exist as compatibility shims and will be removed in a future release.
  - `datalad.distribution.subdatasets` has been moved to `datalad.local.subdatasets` (#3429)
  - `datalad.interface.run` has been moved to `datalad.core.local.run` (#3444)
- The `lock` method of `AnnexRepo` and the `options` parameter of `AnnexRepo.unlock` were unused internally and have been removed. (#3459)
- The `get_submodules` method of `GitRepo` has been rewritten without `GitPython`. When the new `compat` flag is true (the current default), the method returns a value that is compatible with the old return value. This backwards-compatible return value and the `compat` flag will be removed in a future release. (#3508)

- The logic for resolving relative paths given to a command has changed (#3435). The new rule is that relative paths are taken as relative to the dataset only if a dataset *instance* is passed by the caller. In all other scenarios they're considered relative to the current directory.

The main user-visible difference from the command line is that using the `--dataset` argument does *not* result in relative paths being taken as relative to the specified dataset. (The undocumented distinction between “`rel/path`” and “`./rel/path`” no longer exists.)

All commands under `datalad.core` and `datalad.local`, as well as `unlock` and `addurls`, follow the new logic. The goal is for all commands to eventually do so.

### Fixes

- The function for loading JSON streams wasn't clever enough to handle content that included a Unicode line separator like U2028. (#3524)
- When `unlock` was called without an explicit target (i.e., a directory or no paths at all), the call failed if any of the files did not have content present. (#3459)
- `AnnexRepo.get_content_info` failed in the rare case of a key without size information. (#3534)
- `save` ignored `--on-failure` in its underlying call to `status`. (#3470)
- Calling `remove` with a subdirectory displayed spurious warnings about the subdirectory files not existing. (#3586)
- Our processing of `git-annex --json` output mishandled info messages from special remotes. (#3546)
- `create`
  - didn't bypass the “existing subdataset” check when called with `--force` as of 0.12.0rc3 (#3552)
  - failed to register the up-to-date revision of a subdataset when `--cfg-proc` was used with `--dataset` (#3591)
- The base downloader had some error handling that wasn't compatible with Python 3. (#3622)
- Fixed a number of Unicode py2-compatibility issues. (#3602)
- `AnnexRepo.get_content_annexinfo` did not properly chunk file arguments to avoid exceeding the command-line character limit. (#3587)

### Enhancements and new features

- New command `create-sibling-gitlab` provides an interface for creating a publication target on a GitLab instance. (#3447)
- `subdatasets` (#3429)
  - now supports path-constrained queries in the same manner as commands like `save` and `status`
  - gained a `--contains=PATH` option that can be used to restrict the output to datasets that include a specific path.
  - now narrows the listed subdatasets to those underneath the current directory when called with no arguments
- `status` learned to accept a plain `--annex` (no value) as shorthand for `--annex basic`. (#3534)
- The `.dirty` property of `GitRepo` and `AnnexRepo` has been sped up. (#3460)
- The `get_content_info` method of `GitRepo`, used by `status` and commands that depend on `status`, now restricts its git calls to a subset of files, if possible, for a performance gain in repositories with many files. (#3508)

- Extensions that do not provide a command, such as those that provide only metadata extractors, are now supported. (#3531)
- When calling `git-annex` with `--json`, we log standard error at the debug level rather than the warning level if a non-zero exit is expected behavior. (#3518)
- `create` no longer refuses to create a new dataset in the odd scenario of an empty `.git/` directory upstairs. (#3475)
- As of v2.22.0 Git treats a sub-repository on an unborn branch as a repository rather than as a directory. Our documentation and tests have been updated appropriately. (#3476)
- `addurls` learned to accept a `--cfg-proc` value and pass it to its `create` calls. (#3562)

### 1.1.6 0.12.0rc4 (May 15, 2019) – the revolution is over

With the replacement of the `save` command implementation with `rev-save` the revolution effort is now over, and the set of key commands for local dataset operations (`create`, `run`, `save`, `status`, `diff`) is now complete. This new core API is available from `datalad.core.local` (and also via `datalad.api`, as any other command).  
### Major refactoring and deprecations

- The `add` command is now deprecated. It will be removed in a future release.

#### Fixes

- Remove hard-coded dependencies on POSIX path conventions in SSH support code (#3400)
- Emit an `add` result when adding a new subdataset during `save` (#3398)
- SSH file transfer now actually opens a shared connection, if none exists yet (#3403)

#### Enhancements and new features

- `SSHConnection` now offers methods for file upload and download (`get()`, `put()`). The previous `copy()` method only supported upload and was discontinued (#3401)

### 1.1.7 0.12.0rc3 (May 07, 2019) – the revolution continues

Continues API consolidation and replaces the `create` and `diff` command with more performant implementations.

#### Major refactoring and deprecations

- The previous `diff` command has been replaced by the `diff` variant from the `datalad-revolution` extension. (#3366)
- `rev-create` has been renamed to `create`, and the previous `create` has been removed. (#3383)
- The procedure `setup_yoda_dataset` has been renamed to `cfg_yoda` (#3353).
- The `--nosave` of `addurls` now affects only added content, not newly created subdatasets (#3259).
- `Dataset.get_subdatasets` (deprecated since v0.9.0) has been removed. (#3336)
- The `.is_dirty` method of `GitRepo` and `AnnexRepo` has been replaced by `.status` or, for a subset of cases, the `.dirty` property. (#3330)
- `AnnexRepo.get_status` has been replaced by `AnnexRepo.status`. (#3330)

## Fixes

- `status`
  - reported on directories that contained only ignored files (#3238)
  - gave a confusing failure when called from a subdataset with an explicitly specified dataset argument and “.” as a path (#3325)
  - misleadingly claimed that the locally present content size was zero when `--annex basic` was specified (#3378)
- An informative error wasn’t given when a download provider was invalid. (#3258)
- Calling `rev-save PATH` saved unspecified untracked subdatasets. (#3288)
- The available choices for command-line options that take values are now displayed more consistently in the help output. (#3326)
- The new `pathlib`-based code had various encoding issues on Python 2. (#3332)

## Enhancements and new features

- `wtf` now includes information about the Python version. (#3255)
- When operating in an annex repository, checking whether `git-annex` is available is now delayed until a call to `git-annex` is actually needed, allowing systems without `git-annex` to operate on annex repositories in a restricted fashion. (#3274)
- The `load_stream` on helper now supports auto-detection of compressed files. (#3289)
- `create` (formerly `rev-create`)
  - learned to be speedier by passing a path to `status` (#3294)
  - gained a `--cfg-proc` (or `-c`) convenience option for running configuration procedures (or more accurately any procedure that begins with “`cfg_`”) in the newly created dataset (#3353)
- `AnnexRepo.set_metadata` now returns a list while `AnnexRepo.set_metadata_` returns a generator, a behavior which is consistent with the `add` and `add_` method pair. (#3298)
- `AnnexRepo.get_metadata` now supports batch querying of known annex files. Note, however, that callers should carefully validate the input paths because the batch call will silently hang if given non-annex files. (#3364)
- `status`
  - now reports a “`bytesize`” field for files tracked by Git (#3299)
  - gained a new option `eval_subdataset_state` that controls how the subdataset state is evaluated. Depending on the information you need, you can select a less expensive mode to make `status` faster. (#3324)
  - colors deleted files “red” (#3334)
- Querying repository content is faster due to batching of `git cat-file` calls. (#3301)
- The dataset ID of a subdataset is now recorded in the superdataset. (#3304)
- `GitRepo.diffstatus`
  - now avoids subdataset recursion when the comparison is not with the working tree, which substantially improves performance when diffing large dataset hierarchies (#3314)
  - got smarter and faster about labeling a subdataset as “modified” (#3343)



- `GitRepo.get_content_info` now supports disabling the file type evaluation, which gives a performance boost in cases where this information isn't needed. (#3362)
- The XMP metadata extractor now filters based on file name to improve its performance. (#3329)

### 1.1.8 0.12.0rc2 (Mar 18, 2019) – revolution!

#### Fixes

- `GitRepo.dirty` does not report on nested empty directories (#3196).
- `GitRepo.save()` reports results on deleted files.

#### Enhancements and new features

- Absorb a new set of core commands from the datalad-revolution extension:
  - `rev-status`: like `git status`, but simpler and working with dataset hierarchies
  - `rev-save`: a 2-in-1 replacement for `save` and `add`
  - `rev-create`: a ~30% faster `create`
- JSON support tools can now read and write compressed files.

### 1.1.9 0.12.0rc1 (Mar 03, 2019) – to boldly go ...

#### Major refactoring and deprecations

- Discontinued support for `git-annex direct-mode` (also no longer supported upstream).

#### Enhancements and new features

- `Dataset` and `Repo` object instances are now hashable, and can be created based on `pathlib Path` object instances
- Imported various additional methods for the `Repo` classes to query information and save changes.

### 1.1.10 0.11.8 (Oct 11, 2019) – annex-we-are-catching-up

#### Fixes

- Our internal command runner failed to capture output in some cases. (#3656)
- Workaround in the tests around `python` in `cPython >= 3.7.5` ‘;’ in the filename confusing `mimetypes` (#3769) (#3770)

#### Enhancements and new features

- Prepared for upstream changes in `git-annex`, including support for the latest `git-annex`
  - 7.20190912 auto-upgrades v5 repositories to v7. (#3648) (#3682)
  - 7.20191009 fixed treatment of `(larger/smaller)than` in `.gitattributes` (#3765)

- The `cfg_text2git` procedure, as well the `--text-no-annex` option of `create`, now configure `.gitattributes` so that empty files are stored in git rather than annex. (#3667)

### 1.1.11 0.11.7 (Sep 06, 2019) – python2-we-still-love-you-but-...

Primarily bugfixes with some optimizations and refactorings.

#### Fixes

- `addurls`
  - now provides better handling when the URL file isn't in the expected format. (#3579)
  - always considered a relative file for the URL file argument as relative to the current working directory, which goes against the convention used by other commands of taking relative paths as relative to the dataset argument. (#3582)
- `run-procedure`
  - hard coded “python” when formatting the command for non-executable procedures ending with “.py”. `sys.executable` is now used. (#3624)
  - failed if arguments needed more complicated quoting than simply surrounding the value with double quotes. This has been resolved for systems that support `shlex.quote`, but note that on Windows values are left unquoted. (#3626)
- `siblings` now displays an informative error message if a local path is given to `--url` but `--name` isn't specified. (#3555)
- `sshrun`, the command DataLad uses for `GIT_SSH_COMMAND`, didn't support all the parameters that Git expects it to. (#3616)
- Fixed a number of Unicode py2-compatibility issues. (#3597)
- `download-url` now will create leading directories of the output path if they do not exist (#3646)

#### Enhancements and new features

- The `annotate-paths` helper now caches subdatasets it has seen to avoid unnecessary calls. (#3570)
- A repeated configuration query has been dropped from the handling of `--proc-pre` and `--proc-post`. (#3576)
- Calls to `git annex find` now use `--in=.` instead of the alias `--in=here` to take advantage of an optimization that git-annex (as of the current release, 7.20190730) applies only to the former. (#3574)
- `addurls` now suggests close matches when the URL or file format contains an unknown field. (#3594)
- Shared logic used in the `setup.py` files of Datalad and its extensions has been moved to modules in the `_datalad_build_support/` directory. (#3600)
- Get ready for upcoming git-annex dropping support for direct mode (#3631)

### 1.1.12 0.11.6 (Jul 30, 2019) – am I the last of 0.11.x?

Primarily bug fixes to achieve more robust performance

## Fixes

- Our tests needed various adjustments to keep up with upstream changes in Travis and Git. (#3479) (#3492) (#3493)
- `AnnexRepo.is_special_annex_remote` was too selective in what it considered to be a special remote. (#3499)
- We now provide information about unexpected output when `git-annex` is called with `--json`. (#3516)
- Exception logging in the `__del__` method of `GitRepo` and `AnnexRepo` no longer fails if the names it needs are no longer bound. (#3527)
- `addurls` botched the construction of subdataset paths that were more than two levels deep and failed to create datasets in a reliable, breadth-first order. (#3561)
- Cloning a `type=git` special remote showed a spurious warning about the remote not being enabled. (#3547)

## Enhancements and new features

- For calls to `git` and `git-annex`, we disable automatic garbage collection due to past issues with `GitPython`'s state becoming stale, but doing so results in a larger `.git/objects/` directory that isn't cleaned up until garbage collection is triggered outside of `DataLad`. Tests with the latest `GitPython` didn't reveal any state issues, so we've re-enabled automatic garbage collection. (#3458)
- `rerun` learned an `--explicit` flag, which it relays to its calls to `[run][[]]`. This makes it possible to call `rerun` in a dirty working tree (#3498).
- The `metadata` command aborts earlier if a metadata extractor is unavailable. (#3525)

### 1.1.13 0.11.5 (May 23, 2019) – stability is not overrated

Should be faster and less buggy, with a few enhancements.

## Fixes

- `create-sibling` (#3318)
  - Siblings are no longer configured with a post-update hook unless a web interface is requested with `--ui`.
  - `git submodule update --init` is no longer called from the post-update hook.
  - If `--inherit` is given for a dataset without a superdataset, a warning is now given instead of raising an error.
- The internal command runner failed on Python 2 when its `env` argument had unicode values. (#3332)
- The safeguard that prevents creating a dataset in a subdirectory that already contains tracked files for another repository failed on Git versions before 2.14. For older Git versions, we now warn the caller that the safeguard is not active. (#3347)
- A regression introduced in v0.11.1 prevented `save` from committing changes under a subdirectory when the subdirectory was specified as a path argument. (#3106)
- A workaround introduced in v0.11.1 made it possible for `save` to do a partial commit with an annex file that has gone below the `annex.largefiles` threshold. The logic of this workaround was faulty, leading to files being displayed as typechanged in the index following the commit. (#3365)
- The `resolve_path()` helper confused paths that had a semicolon for SSH RIs. (#3425)

- The detection of SSH RIs has been improved. (#3425)

### Enhancements and new features

- The internal command runner was too aggressive in its decision to sleep. (#3322)
- The “INFO” label in log messages now retains the default text color for the terminal rather than using white, which only worked well for terminals with dark backgrounds. (#3334)
- A short flag `-R` is now available for the `--recursion-limit` flag, a flag shared by several subcommands. (#3340)
- The authentication logic for `create-sibling-github` has been revamped and now supports 2FA. (#3180)
- New configuration option `datalad.ui.progressbar` can be used to configure the default backend for progress reporting (“none”, for example, results in no progress bars being shown). (#3396)
- A new progress backend, available by setting `datalad.ui.progressbar` to “log”, replaces progress bars with a log message upon completion of an action. (#3396)
- DataLad learned to consult the `NO_COLOR` environment variable and the new `datalad.ui.color` configuration option when deciding to color output. The default value, “auto”, retains the current behavior of coloring output if attached to a TTY (#3407).
- `clean` now removes annex transfer directories, which is useful for cleaning up failed downloads. (#3374)
- `clone` no longer refuses to clone into a local path that looks like a URL, making its behavior consistent with `git clone`. (#3425)
- `wtf`
  - Learned to fall back to the `dist` package if `platform.dist`, which has been removed in the yet-to-be-release Python 3.8, does not exist. (#3439)
  - Gained a `--section` option for limiting the output to specific sections and a `--decor` option, which currently knows how to format the output as GitHub’s `<details>` section. (#3440)

### 1.1.14 0.11.4 (Mar 18, 2019) – get-ready

Largely a bug fix release with a few enhancements

#### Important

- 0.11.x series will be the last one with support for direct mode of `git-annex` which is used on crippled (no symlinks and no locking) filesystems. v7 repositories should be used instead.

#### Fixes

- Extraction of `.gz` files is broken without `p7zip` installed. We now abort with an informative error in this situation. (#3176)
- Committing failed in some cases because we didn’t ensure that the path passed to `git read-tree --index-output=...` resided on the same filesystem as the repository. (#3181)
- Some pointless warnings during metadata aggregation have been eliminated. (#3186)
- With Python 3 the LORIS token authenticator did not properly decode a response (#3205).

- With Python 3 downloaders unnecessarily decoded the response when getting the status, leading to an encoding error. (#3210)
- In some cases, our internal command Runner did not adjust the environment's PWD to match the current working directory specified with the `cwd` parameter. (#3215)
- The specification of the `pyliblzma` dependency was broken. (#3220)
- `search` displayed an uninformative blank log message in some cases. (#3222)
- The logic for finding the location of the aggregate metadata DB anchored the search path incorrectly, leading to a spurious warning. (#3241)
- Some progress bars were still displayed when `stdout` and `stderr` were not attached to a `tty`. (#3281)
- Check for `stdin/out/err` to not be closed before checking for `.isatty`. (#3268)

### Enhancements and new features

- Creating a new repository now aborts if any of the files in the directory are tracked by a repository in a parent directory. (#3211)
- `run` learned to replace the `{tmpdir}` placeholder in commands with a temporary directory. (#3223)
- `ducredit` support has been added for citing DataLad itself as well as datasets that an analysis uses. (#3184)
- The `eval_results` interface helper unintentionally modified one of its arguments. (#3249)
- A few DataLad constants have been added, changed, or renamed (#3250):
  - `HANDLE_META_DIR` is now `DATALAD_DOTDIR`. The old name should be considered deprecated.
  - `METADATA_DIR` now refers to `DATALAD_DOTDIR/metadata` rather than `DATALAD_DOTDIR/meta` (which is still available as `OLDMETADATA_DIR`).
  - The new `DATASET_METADATA_FILE` refers to `METADATA_DIR/dataset.json`.
  - The new `DATASET_CONFIG_FILE` refers to `DATALAD_DOTDIR/config`.
  - `METADATA_FILENAME` has been renamed to `OLDMETADATA_FILENAME`.

### 1.1.15 0.11.3 (Feb 19, 2019) – read-me-gently

Just a few of important fixes and minor enhancements.

#### Fixes

- The logic for setting the maximum command line length now works around Python 3.4 returning an unreasonably high value for `SC_ARG_MAX` on Debian systems. (#3165)
- DataLad commands that are conceptually “read-only”, such as `datalad ls -L`, can fail when the caller lacks write permissions because `git-annex` tries merging remote `git-annex` branches to update information about availability. DataLad now disables `annex.merge-annex-branches` in some common “read-only” scenarios to avoid these failures. (#3164)

## Enhancements and new features

- Accessing an “unbound” dataset method now automatically imports the necessary module rather than requiring an explicit import from the Python caller. For example, calling `Dataset.add` no longer needs to be preceded by `from datalad.distribution.add import Add` or an `import` of `datalad.api`. (#3156)
- Configuring the new variable `datalad.ssh.identityfile` instructs DataLad to pass a value to the `-i` option of `ssh`. (#3149) (#3168)

### 1.1.16 0.11.2 (Feb 07, 2019) – live-long-and-prosper

A variety of bugfixes and enhancements

#### Major refactoring and deprecations

- All extracted metadata is now placed under `git-annex` by default. Previously files smaller than 20 kb were stored in `git`. (#3109)
- The function `datalad.cmd.get_runner` has been removed. (#3104)

#### Fixes

- Improved handling of long commands:
  - The code that inspected `SC_ARG_MAX` didn’t check that the reported value was a sensible, positive number. (#3025)
  - More commands that invoke `git` and `git-annex` with file arguments learned to split up the command calls when it is likely that the command would fail due to exceeding the maximum supported length. (#3138)
- The `setup_yoda_dataset` procedure created a malformed `.gitattributes` line. (#3057)
- `download-url` unnecessarily tried to infer the dataset when `--no-save` was given. (#3029)
- `rerun` aborted too late and with a confusing message when a ref specified via `--onto` didn’t exist. (#3019)
- `run`:
  - `run` didn’t preserve the current directory prefix (“.”) on inputs and outputs, which is problematic if the caller relies on this representation when formatting the command. (#3037)
  - Fixed a number of unicode py2-compatibility issues. (#3035) (#3046)
  - To proceed with a failed command, the user was confusingly instructed to use `save` instead of `add` even though `run` uses `add` underneath. (#3080)
- Fixed a case where the helper class for checking external modules incorrectly reported a module as unknown. (#3051)
- `add-archive-content` mishandled the archive path when the leading path contained a symlink. (#3058)
- Following denied access, the credential code failed to consider a scenario, leading to a type error rather than an appropriate error message. (#3091)
- Some tests failed when executed from a `git worktree` checkout of the source repository. (#3129)
- During metadata extraction, batched annex processes weren’t properly terminated, leading to issues on Windows. (#3137)

- `add` incorrectly handled an “invalid repository” exception when trying to add a submodule. (#3141)
- Pass `GIT_SSH_VARIANT=ssh` to `git` processes to be able to specify alternative ports in SSH urls

## Enhancements and new features

- `search` learned to suggest closely matching keys if there are no hits. (#3089)
- `create-sibling`
  - gained a `--group` option so that the caller can specify the file system group for the repository. (#3098)
  - now understands SSH URLs that have a port in them (i.e. the “`ssh://[{}user@{}]host.xz[{}:port{}]/path/to/repo.git/`” syntax mentioned in `man git-fetch`). (#3146)
- Interface classes can now override the default renderer for summarizing results. (#3061)
- `run`:
  - `--input` and `--output` can now be shortened to `-i` and `-o`. (#3066)
  - Placeholders such as “`{inputs}`” are now expanded in the command that is shown in the commit message subject. (#3065)
  - `interface.run.run_command` gained an `extra_inputs` argument so that wrappers like `datalad-container` can specify additional inputs that aren’t considered when formatting the command string. (#3038)
  - “`-`” can now be used to separate options for `run` and those for the command in ambiguous cases. (#3119)
- The utilities `create_tree` and `ok_file_has_content` now support “.gz” files. (#3049)
- The Singularity container for 0.11.1 now uses `nd_freeze` to make its builds reproducible.
- A `publications` page has been added to the documentation. (#3099)
- `GitRepo.set_gitattributes` now accepts a `mode` argument that controls whether the `.gitattributes` file is appended to (default) or overwritten. (#3115)
- `datalad --help` now avoids using `man` so that the list of subcommands is shown. (#3124)

### 1.1.17 0.11.1 (Nov 26, 2018) – v7-better-than-v6

Rushed out bugfix release to stay fully compatible with recent `git-annex` which introduced v7 to replace v6.

## Fixes

- `install`: be able to install recursively into a dataset (#2982)
- `save`: be able to commit/save changes whenever files potentially could have swapped their storage between `git` and `annex` (#1651) (#2752) (#3009)
- `[aggregate-metadata][[]]`:
  - `dataset`’s itself is now not “aggregated” if specific paths are provided for aggregation (#3002). That resolves the issue of `-r` invocation aggregating all subdatasets of the specified dataset as well
  - also compare/verify the actual content checksum of aggregated metadata while considering subdataset metadata for re-aggregation (#3007)
- `annex` commands are now chunked assuming 50% “safety margin” on the maximal command line length. Should resolve crashes while operating on too many files at ones (#3001)

- `run` sidecar config processing (#2991)
- no double trailing period in docs (#2984)
- correct identification of the repository with symlinks in the paths in the tests (#2972)
- re-evaluation of dataset properties in case of dataset changes (#2946)
- `[text2git][]` procedure to use `ds.repo.set_gitattributes` (#2974) (#2954)
- Switch to use plain `os.getcwd()` if inconsistency with env var `$PWD` is detected (#2914)
- Make sure that credential defined in env var takes precedence (#2960) (#2950)

### Enhancements and new features

- `shub://datalad/datalad:git-annex-dev` provides a Debian buster Singularity image with build environment for `git-annex`. `tools/bisect-git-annex` provides a helper for running `git bisect` on `git-annex` using that Singularity container (#2995)
- Added `.zenodo.json` for better integration with Zenodo for citation
- `run-procedure` now provides names and help messages with a custom renderer for (#2993)
- Documentation: point to `datalad-revolution` extension (prototype of the greater DataLad future)
- `run`
  - support injecting of a detached command (#2937)
- `annex` metadata extractor now extracts `annex.key` metadata record. Should allow now to identify uses of specific files etc (#2952)
- Test that we can install from `http://datasets.datalad.org`
- Proper rendering of `CommandError` (e.g. in case of “out of space” error) (#2958)

### 1.1.18 0.11.0 (Oct 23, 2018) – Soon-to-be-perfect

`git-annex` 6.20180913 (or later) is now required - provides a number of fixes for v6 mode operations etc.

### Major refactoring and deprecations

- `datalad.consts.LOCAL_CENTRAL_PATH` constant was deprecated in favor of `datalad.locations.default-dataset configuration` variable (#2835)

### Minor refactoring

- "notneeded" messages are no longer reported by default results renderer
- `run` no longer shows commit instructions upon command failure when `explicit` is true and no outputs are specified (#2922)
- `get_git_dir` moved into `GitRepo` (#2886)
- `_gitpy_custom_call` removed from `GitRepo` (#2894)
- `GitRepo.get_merge_base` argument is now called `commitishes` instead of `treeishes` (#2903)



## Fixes

- `update` should not leave the dataset in non-clean state (#2858) and some other enhancements (#2859)
- Fixed chunking of the long command lines to account for decorators and other arguments (#2864)
- Progress bar should not crash the process on some missing progress information (#2891)
- Default value for `jobs` set to be "auto" (not None) to take advantage of possible parallel get if in `-g` mode (#2861)
- `wtf` must not crash if `git-annex` is not installed etc (#2865), (#2865), (#2918), (#2917)
- Fixed paths (with spaces etc) handling while reporting annex error output (#2892), (#2893)
- `__del__` should not access `.repo` but `._repo` to avoid attempts for instantiation etc (#2901)
- Fix up submodule `.git` right in `GitRepo.add_submodule` to avoid added submodules being non git-annex friendly (#2909), (#2904)
- `run-procedure` (#2905)
  - now will provide dataset into the procedure if called within dataset
  - will not crash if procedure is an executable without `.py` or `.sh` suffixes
- Use centralized `.gitattributes` handling while setting annex backend (#2912)
- `GlobbedPaths.expand(..., full=True)` incorrectly returned relative paths when called more than once (#2921)

## Enhancements and new features

- Report progress on `clone` when installing from “smart” git servers (#2876)
- Stale/unused `sth_like_file_has_content` was removed (#2860)
- Enhancements to `search` to operate on “improved” metadata layouts (#2878)
- Output of `git annex init` operation is now logged (#2881)
- New
  - `GitRepo.cherry_pick` (#2900)
  - `GitRepo.format_commit` (#2902)
- `run-procedure` (#2905)
  - procedures can now recursively be discovered in subdatasets as well. The uppermost has highest priority
  - Procedures in user and system locations now take precedence over those in datasets.

### 1.1.19 0.10.3.1 (Sep 13, 2018) – Nothing-is-perfect

Emergency bugfix to address forgotten boost of version in `datalad/version.py`.

### 1.1.20 0.10.3 (Sep 13, 2018) – Almost-perfect

This is largely a bugfix release which addressed many (but not yet all) issues of working with `git-annex` direct and version 6 modes, and operation on Windows in general. Among enhancements you will see the support of public S3 buckets (even with periods in their names), ability to configure new providers interactively, and improved `egrep` search backend.

Although we do not require with this release, it is recommended to make sure that you are using a recent `git-annex` since it also had a variety of fixes and enhancements in the past months.

#### Fixes

- Parsing of combined short options has been broken since DataLad v0.10.0. (#2710)
- The `datalad save` instructions shown by `datalad run` for a command with a non-zero exit were incorrectly formatted. (#2692)
- Decompression of zip files (e.g., through `datalad add-archive-content`) failed on Python 3. (#2702)
- Windows:
  - colored log output was not being processed by `colorama`. (#2707)
  - more codepaths now try multiple times when removing a file to deal with latency and locking issues on Windows. (#2795)
- Internal `git` fetch calls have been updated to work around a `GitPython` `BadName` issue. (#2712), (#2794)
- The progress bar for annex file transferring was unable to handle an empty file. (#2717)
- `datalad add-readme` halted when no aggregated metadata was found rather than displaying a warning. (#2731)
- `datalad rerun` failed if `--onto` was specified and the history contained no run commits. (#2761)
- Processing of a command's results failed on a result record with a missing value (e.g., absent field or subfield in metadata). Now the missing value is rendered as "N/A". (#2725).
- A couple of documentation links in the "Delineation from related solutions" were misformatted. (#2773)
- With the latest `git-annex`, several known V6 failures are no longer an issue. (#2777)
- In direct mode, commit changes would often commit annexed content as regular `Git` files. A new approach fixes this and resolves a good number of known failures. (#2770)
- The reporting of command results failed if the current working directory was removed (e.g., after an unsuccessful `install`). (#2788)
- When installing into an existing empty directory, `datalad install` removed the directory after a failed clone. (#2788)
- `datalad run` incorrectly handled inputs and outputs for paths with spaces and other characters that require shell escaping. (#2798)
- Globbing inputs and outputs for `datalad run` didn't work correctly if a subdataset wasn't installed. (#2796)
- Minor (in)compatibility with `git 2.19` - (no) trailing period in an error message now. (#2815)

#### Enhancements and new features

- Anonymous access is now supported for S3 and other downloaders. (#2708)

- A new interface is available to ease setting up new providers. (#2708)
- Metadata: changes to egrep mode search (#2735)
  - Queries in egrep mode are now case-sensitive when the query contains any uppercase letters and are case-insensitive otherwise. The new mode `egrepcs` can be used to perform a case-sensitive query with all lower-case letters.
  - Search can now be limited to a specific key.
  - Multiple queries (list of expressions) are evaluated using AND to determine whether something is a hit.
  - A single multi-field query (e.g., `pa* : findme`) is a hit, when any matching field matches the query.
  - All matching key/value combinations across all (multi-field) queries are reported in the `query_matched` result field.
  - egrep mode now shows all hits rather than limiting the results to the top 20 hits.
- The documentation on how to format commands for `datalad run` has been improved. (#2703)
- The method for determining the current working directory on Windows has been improved. (#2707)
- `datalad --version` now simply shows the version without the license. (#2733)
- `datalad export-archive` learned to export under an existing directory via its `--filename` option. (#2723)
- `datalad export-to-figshare` now generates the zip archive in the root of the dataset unless `--filename` is specified. (#2723)
- After importing `datalad.api`, `help(datalad.api)` (or `datalad.api?` in IPython) now shows a summary of the available DataLad commands. (#2728)
- Support for using `datalad` from IPython has been improved. (#2722)
- `datalad wtf` now returns structured data and reports the version of each extension. (#2741)
- The internal handling of gitattributes information has been improved. A user-visible consequence is that `datalad create --force` no longer duplicates existing attributes. (#2744)
- The “annex” metadata extractor can now be used even when no content is present. (#2724)
- The `add_url_to_file` method (called by commands like `datalad download-url` and `datalad add-archive-content`) learned how to display a progress bar. (#2738)

### 1.1.21 0.10.2 (Jul 09, 2018) – Thesecuriestever

Primarily a bugfix release to accommodate recent git-annex release forbidding `file://` and `http://localhost/` URLs which might lead to revealing private files if annex is publicly shared.

#### Fixes

- fixed testing to be compatible with recent git-annex (6.20180626)
- `download-url` will now download to current directory instead of the top of the dataset

## Enhancements and new features

- do not quote ~ in URLs to be consistent with quote implementation in Python 3.7 which now follows RFC 3986
- `run` support for user-configured placeholder values
- documentation on native git-annex metadata support
- handle 401 errors from LORIS tokens
- `yoda` procedure will instantiate `README.md`
- `--discover` option added to `run-procedure` to list available procedures

### 1.1.22 0.10.1 (Jun 17, 2018) – OHBM polish

This is a minor bugfix release.

#### Fixes

- Be able to use `backports.lzma` as a drop-in replacement for `pyliblzma`.
- Give help when not specifying a procedure name in `run-procedure`.
- Abort early when a downloader received no filename.
- Avoid `rerun` error when trying to unlock non-available files.

### 1.1.23 0.10.0 (Jun 09, 2018) – The Release

This release is a major leap forward in metadata support.

#### Major refactoring and deprecations

- Metadata
  - Prior metadata provided by datasets under `.datalad/meta` is no longer used or supported. Metadata must be reaggregated using 0.10 version
  - Metadata extractor types are no longer auto-guessed and must be explicitly specified in `datalad.metadata.nativetype` config (could contain multiple values)
  - Metadata aggregation of a dataset hierarchy no longer updates all datasets in the tree with new metadata. Instead, only the target dataset is updated. This behavior can be changed via the `-update-mode` switch. The new default prevents needless modification of (3rd-party) subdatasets.
  - Neuroimaging metadata support has been moved into a dedicated extension: <https://github.com/datalad/datalad-neuroimaging>
- Crawler
  - moved into a dedicated extension: <https://github.com/datalad/datalad-crawler>
- `export_tarball` plugin has been generalized to `export_archive` and can now also generate ZIP archives.
- By default a dataset X is now only considered to be a super-dataset of another dataset Y, if Y is also a registered subdataset of X.

## Fixes

A number of fixes did not make it into the 0.9.x series:

- Dynamic configuration overrides via the `-c` option were not in effect.
- `save` is now more robust with respect to invocation in subdirectories of a dataset.
- `unlock` now reports correct paths when running in a dataset subdirectory.
- `get` is more robust to path that contain symbolic links.
- symlinks to subdatasets of a dataset are now correctly treated as a symlink, and not as a subdataset
- `add` now correctly saves staged subdataset additions.
- Running `datalad save` in a dataset no longer adds untracked content to the dataset. In order to add content a path has to be given, e.g. `datalad save .`
- `wtf` now works reliably with a DataLad that wasn't installed from Git (but, e.g., via pip)
- More robust URL handling in `simple_with_archives` crawler pipeline.

## Enhancements and new features

- Support for DataLad extension that can contribute API components from 3rd-party sources, incl. commands, metadata extractors, and test case implementations. See <https://github.com/datalad/datalad-extension-template> for a demo extension.
- Metadata (everything has changed!)
  - Metadata extraction and aggregation is now supported for datasets and individual files.
  - Metadata query via `search` can now discover individual files.
  - Extracted metadata can now be stored in XZ compressed files, is optionally annexed (when exceeding a configurable size threshold), and obtained on demand (new configuration option `datalad.metadata.create-aggregate-annex-limit`).
  - Status and availability of aggregated metadata can now be reported via `metadata --get-aggregates`
  - New configuration option `datalad.metadata.maxfieldsize` to exclude too large metadata fields from aggregation.
  - The type of metadata is no longer guessed during metadata extraction. A new configuration option `datalad.metadata.nativetype` was introduced to enable one or more particular metadata extractors for a dataset.
  - New configuration option `datalad.metadata.store-aggregate-content` to enable the storage of aggregated metadata for dataset content (i.e. file-based metadata) in contrast to just metadata describing a dataset as a whole.
- `search` was completely reimplemented. It offers three different modes now:
  - ‘egrep’ (default): expression matching in a plain string version of metadata
  - ‘textblob’: search a text version of all metadata using a fully featured query language (fast indexing, good for keyword search)
  - ‘autofield’: search an auto-generated index that preserves individual fields of metadata that can be represented in a tabular structure (substantial indexing cost, enables the most detailed queries of all modes)
- New extensions:

- `addurls`, an extension for creating a dataset (and possibly subdatasets) from a list of URLs.
- `export_to_figshare`
- `extract_metadata`
- `add_readme` makes use of available metadata
- By default the `wtf` extension now hides sensitive information, which can be included in the output by passing `--sensitive=some` or `--sensitive=all`.
- Reduced startup latency by only importing commands necessary for a particular command line call.
- `create`:
  - `-d <parent> --nosave` now registers subdatasets, when possible.
  - `--fake-dates` configures dataset to use fake-dates
- `run` now provides a way for the caller to save the result when a command has a non-zero exit status.
- `datalad rerun` now has a `--script` option that can be used to extract previous commands into a file.
- A DataLad Singularity container is now available on [Singularity Hub](#).
- More casts have been embedded in the [use case section of the documentation](#).
- `datalad --report-status` has a new value 'all' that can be used to temporarily re-enable reporting that was disabled by configuration settings.

### 1.1.24 0.9.3 (Mar 16, 2018) – pi+0.02 release

Some important bug fixes which should improve usability

#### Fixes

- `datalad-archives` special remote now will lock on acquiring or extracting an archive - this allows for it to be used with `-J` flag for parallel operation
- `relax` introduced in 0.9.2 demand on git being configured for datalad operation - now we will just issue a warning
- `datalad ls` should now list “authored date” and work also for datasets in detached HEAD mode
- `datalad save` will now save original file as well, if file was “git mv”ed, so you can now `datalad run git mv old new` and have changes recorded

#### Enhancements and new features

- `--jobs` argument now could take `auto` value which would decide on # of jobs depending on the # of available CPUs. `git-annex > 6.20180314` is recommended to avoid regression with `-J`.
- memoize calls to RI meta-constructor – should speed up operation a bit
- `DATALAD_SEED` environment variable could be used to seed Python RNG and provide reproducible UUIDs etc (useful for testing and demos)

### 1.1.25 0.9.2 (Mar 04, 2018) – it is (again) better than ever

Largely a bugfix release with a few enhancements.

## Fixes

- Execution of external commands (`git`) should not get stuck when lots of both stdout and stderr output, and should not lose remaining output in some cases
- Config overrides provided in the command line (`-c`) should now be handled correctly
- Consider more remotes (not just tracking one, which might be none) while installing subdatasets
- Compatibility with `git` 2.16 with some changed behaviors/annotations for submodules
- Fail `remove` if `annex drop` failed
- Do not fail operating on files which start with dash (`-`)
- URL unquote paths within S3, URLs and DataLad RIs (`///`)
- In non-interactive mode fail if authentication/access fails
- Web UI:
  - refactored a little to fix incorrect listing of submodules in subdirectories
  - now auto-focuses on search edit box upon entering the page
- Assure that extracted from tarballs directories have executable bit set

## Enhancements and new features

- A log message and progress bar will now inform if a tarball to be downloaded while getting specific files (requires `git-annex > 6.20180206`)
- A dedicated `datalad rerun` command capable of rerunning entire sequences of previously `run` commands. **Reproducibility through VCS. Use “run“ even if not interested in “rerun“**
- Alert the user if `git` is not yet configured but `git` operations are requested
- Delay collection of previous `ssh` connections until it is actually needed. Also do not require `:` while specifying `ssh` host
- AutomagicIO: Added proxying of `isfile`, `lzma.LZMAFile` and `io.open`
- Testing:
  - added `DATALAD_DATASETS_TOPURL=http://datasets-tests.datalad.org` to run tests against another website to not obscure access stats
  - tests run against temporary `HOME` to avoid side-effects
  - better unit-testing of interactions with special remotes
- `CONTRIBUTING.md` describes how to setup and use `git-hub` tool to “attach” commits to an issue making it into a PR
- `DATALAD_USE_DEFAULT_GIT` env variable could be used to cause DataLad to use default (not the one possibly bundled with `git-annex`) `git`
- Be more robust while handling not supported requests by `annex` in special remotes
- Use of `swallow_logs` in the code was refactored away – less mysteries now, just increase logging level
- `wtf` plugin will report more information about environment, externals and the system

### 1.1.26 0.9.1 (Oct 01, 2017) – “DATALAD!”(JBTM)

Minor bugfix release

#### Fixes

- Should work correctly with subdatasets named as numbers of bool values (requires also GitPython >= 2.1.6)
- Custom special remotes should work without crashing with git-annex >= 6.20170924

### 1.1.27 0.9.0 (Sep 19, 2017) – isn't it a lucky day even though not a Friday?

#### Major refactoring and deprecations

- the `files` argument of `save` has been renamed to `path` to be uniform with any other command
- all major commands now implement more uniform API semantics and result reporting. Functionality for modification detection of dataset content has been completely replaced with a more efficient implementation
- `publish` now features a `--transfer-data` switch that allows for a disambiguous specification of whether to publish data – independent of the selection which datasets to publish (which is done via their paths). Moreover, `publish` now transfers data before repository content is pushed.

#### Fixes

- `drop` no longer errors when some subdatasets are not installed
- `install` will no longer report nothing when a Dataset instance was given as a source argument, but rather perform as expected
- `remove` doesn't remove when some files of a dataset could not be dropped
- `publish`
  - no longer hides error during a repository push
  - `publish` behaves “correctly” for `--since=` in considering only the differences the last “pushed” state
  - data transfer handling while publishing with dependencies, to github
- improved robustness with broken Git configuration
- `search` should search for unicode strings correctly and not crash
- robustify git-annex special remotes protocol handling to allow for spaces in the last argument
- UI credentials interface should now allow to Ctrl-C the entry
- should not fail while operating on submodules named with numerics only or by bool (true/false) names
- crawl templates should not now override settings for `largefiles` if specified in `.gitattributes`

#### Enhancements and new features

- **Exciting new feature** `run` command to protocol execution of an external command and rerun computation if desired. See [screencast](#)
- `save` now uses Git for detecting with subdatasets need to be inspected for potential changes, instead of performing a complete traversal of a dataset tree



- `add` looks for changes relative to the last committed state of a dataset to discover files to add more efficiently
- `diff` can now report untracked files in addition to modified files
- `[uninstall]()` will check itself whether a subdataset is properly registered in a superdataset, even when no superdataset is given in a call
- `subdatasets` can now configure subdatasets for exclusion from recursive installation (`datalad-recursiveinstall` submodule configuration property)
- precrafted pipelines of `[crawl]()` now will not override `annex.largefiles` setting if any was set within `.gitattributes` (e.g. by `datalad create --text-no-annex`)
- framework for screencasts: `tools/cast*` tools and sample cast scripts under `doc/casts` which are published at [datalad.org/features.html](http://datalad.org/features.html)
- new project YouTube channel
- tests failing in direct and/or v6 modes marked explicitly

### 1.1.28 0.8.1 (Aug 13, 2017) – the best birthday gift

Bugfixes

#### Fixes

- Do not attempt to `update` a not installed sub-dataset
- In case of too many files to be specified for `get` or `copy_to`, we will make multiple invocations of underlying `git-annex` command to not overflow command line
- More robust handling of unicode output in terminals which might not support it

#### Enhancements and new features

- Ship a copy of `numpy.testing` to facilitate `[test]()` without requiring `numpy` as dependency. Also allow to pass to command which `test(s)` to run
- In `get` and `copy_to` provide actual original requested paths, not the ones we deduced need to be transferred, solely for knowing the total

### 1.1.29 0.8.0 (Jul 31, 2017) – it is better than ever

A variety of fixes and enhancements

#### Fixes

- `publish` would now push merged `git-annex` branch even if no other changes were done
- `publish` should be able to publish using relative path within SSH URI (`git hook` would use relative paths)
- `publish` should better tolerate publishing to pure `git` and `git-annex` special remotes

## Enhancements and new features

- `plugin` mechanism came to replace `export`. See `export_tarball` for the replacement of `export`. Now it should be easy to extend datalad's interface with custom functionality to be invoked along with other commands.
- Minimalistic coloring of the results rendering
- `publish/copy_to` got progress bar report now and support of `--jobs`
- minor fixes and enhancements to crawler (e.g. support of recursive removes)

### 1.1.30 0.7.0 (Jun 25, 2017) – when it works - it is quite awesome!

New features, refactorings, and bug fixes.

#### Major refactoring and deprecations

- `add-sibling` has been fully replaced by the `siblings` command
- `create-sibling`, and `unlock` have been re-written to support the same common API as most other commands

#### Enhancements and new features

- `siblings` can now be used to query and configure a local repository by using the sibling name `here`
- `siblings` can now query and set annex preferred content configuration. This includes `wanted` (as previously supported in other commands), and now also `required`
- New `metadata` command to interface with datasets/files `meta-data`
- Documentation for all commands is now built in a uniform fashion
- Significant parts of the documentation of been updated
- Instantiate GitPython's Repo instances lazily

#### Fixes

- API documentation is now rendered properly as HTML, and is easier to browse by having more compact pages
- Closed files left open on various occasions (Popen PIPEs, etc)
- Restored basic (consumer mode of operation) compatibility with Windows OS

### 1.1.31 0.6.0 (Jun 14, 2017) – German perfectionism

This release includes a **huge** refactoring to make code base and functionality more robust and flexible

- outputs from API commands could now be highly customized. See `--output-format`, `--report-status`, `--report-type`, and `--report-type` options for `datalad` command.
- effort was made to refactor code base so that underlying functions behave as generators where possible
- input paths/arguments analysis was redone for majority of the commands to provide unified behavior

## Major refactoring and deprecations

- `add-sibling` and `rewrite-urls` were refactored in favor of new `siblings` command which should be used for siblings manipulations
- `'datalad.api.alwaysrender'` config setting/support is removed in favor of new outputs processing

## Fixes

- Do not flush manually git index in pre-commit to avoid “Death by the Lock” issue
- Deployed by `publish post-update` hook script now should be more robust (tolerate directory names with spaces, etc.)
- A variety of fixes, see [list of pull requests and issues closed](#) for more information

## Enhancements and new features

- new `annotate-paths` plumbing command to inspect and annotate provided paths. Use `--modified` to summarize changes between different points in the history
- new `clone` plumbing command to provide a subset (install a single dataset from a URL) functionality of `install`
- new `diff` plumbing command
- new `siblings` command to list or manipulate siblings
- new `subdatasets` command to list subdatasets and their properties
- `drop` and `remove` commands were refactored
- `benchmarks/` collection of [Airspeed velocity](#) benchmarks initiated. See reports at <http://datalad.github.io/datalad/>
- crawler would try to download a new url multiple times increasing delay between attempts. Helps to resolve problems with extended crawls of Amazon S3
- `CRCNS` crawler pipeline now also fetches and aggregates meta-data for the datasets from datacite
- overall optimisations to benefit from the aforementioned refactoring and improve user-experience
- a few stub and not (yet) implemented commands (e.g. `move`) were removed from the interface
- Web frontend got proper coloring for the breadcrumbs and some additional caching to speed up interactions. See <http://datasets.datalad.org>
- Small improvements to the online documentation. See e.g. [summary of differences between git/git-annex/datalad](#)

### 1.1.32 0.5.1 (Mar 25, 2017) – cannot stop the progress

A bugfix release

## Fixes

- `add` was forcing addition of files to annex regardless of settings in `.gitattributes`. Now that decision is left to annex by default

- `tools/testing/run_doc_examples` used to run doc examples as tests, fixed up to provide status per each example and not fail at once
- `doc/examples`
  - `3rdparty_analysis_workflow.sh` was fixed up to reflect changes in the API of 0.5.0.
- progress bars
  - should no longer crash **datalad** and report correct sizes and speeds
  - should provide progress reports while using Python 3.x

### Enhancements and new features

- `doc/examples`
  - `nipype_workshop_dataset.sh` new example to demonstrate how new super- and sub- datasets were established as a part of our datasets collection

### 1.1.33 0.5.0 (Mar 20, 2017) – it’s huge

This release includes an avalanche of bug fixes, enhancements, and additions which at large should stay consistent with previous behavior but provide better functioning. Lots of code was refactored to provide more consistent code-base, and some API breakage has happened. Further work is ongoing to standardize output and results reporting ([#1350](#))

### Most notable changes

- requires `git-annex`  $\geq 6.20161210$  (or better even  $\geq 6.20161210$  for improved functionality)
- commands should now operate on paths specified (if any), without causing side-effects on other dirty/staged files
- `save`
  - `-a` is deprecated in favor of `-u` or `--all-updates` so only changes known components get saved, and no new files automatically added
  - `-S` does no longer store the originating dataset in its commit message
- `add`
  - can specify commit/save message with `-m`
- `add-sibling` and `create-sibling`
  - now take the name of the sibling (remote) as a `-s` (`--name`) option, not a positional argument
  - `--publish-depends` to setup publishing data and code to multiple repositories (e.g. github + web-serve) should now be functional see [this comment](#)
  - got `--publish-by-default` to specify what refs should be published by default
  - got `--annex-wanted`, `--annex-groupwanted` and `--annex-group` settings which would be used to instruct annex about preferred content. `publish` then will publish data using those settings if `wanted` is set.
  - got `--inherit` option to automatically figure out url/wanted and other git/annex settings for new remote sub-dataset to be constructed
- `publish`

- `got --skip-failing` refactored into `--missing` option which could use new feature of `create-sibling --inherit`

## Fixes

- More consistent interaction through ssh - all ssh connections go through `sshrun` shim for a “single point of authentication”, etc.
- More robust `ls` operation outside of the datasets
- A number of fixes for direct and v6 mode of annex

## Enhancements and new features

- New `drop` and `remove` commands
- `clean`
  - `got --what` to specify explicitly what cleaning steps to perform and now could be invoked with `-r`
- `datalad` and `git-annex-remote*` scripts now do not use `setuptools` entry points mechanism and rely on simple `import` to shorten start up time
- `Dataset` is also now using `Flyweight pattern`, so the same instance is reused for the same dataset
- progressbars should not add more empty lines

## Internal refactoring

- Majority of the commands now go through `_prep` for arguments validation and pre-processing to avoid recursive invocations

## 1.1.34 0.4.1 (Nov 10, 2016) – CA release

Requires now `GitPython >= 2.1.0`

## Fixes

- `save`
  - to not save staged files if explicit paths were provided
- improved (but not yet complete) support for direct mode
- `update` to not crash if some sub-datasets are not installed
- do not log calls to `git config` to avoid leakage of possibly sensitive settings to the logs

## Enhancements and new features

- New `rfc822-compliant metadata` format
- `save`
  - `-S` to save the change also within all super-datasets
- `add` now has progress-bar reporting

- `create-sibling-github` to create a `:term:sibling` of a dataset on github
- `OpenfMRI` crawler and datasets were enriched with URLs to separate files where also available from `openfmri s3` bucket (if upgrading your `datalad` datasets, you might need to run `git annex enableremote datalad` to make them available)
- various enhancements to log messages
- web interface
  - populates “install” box first thus making UX better over slower connections

### 1.1.35 0.4 (Oct 22, 2016) – Paris is waiting

Primarily it is a bugfix release but because of significant refactoring of the `install` and `get` implementation, it gets a new minor release.

#### Fixes

- be able to `get` or `install` while providing paths while being outside of a dataset
- remote annex datasets get properly initialized
- robust detection of outdated `git-annex`

#### Enhancements and new features

- interface changes
  - `get --recursion-limit=existing` to not recurse into not-installed subdatasets
  - `get -n` to possibly install sub-datasets without getting any data
  - `install --jobs | -J` to specify number of parallel jobs for annex `get` call could use (ATM would not work when data comes from archives)
- more (unit-)testing
- documentation: see <http://docs.datalad.org/en/latest/basics.html> for basic principles and useful shortcuts in referring to datasets
- various webface improvements: breadcrumb paths, instructions how to install dataset, show version from the tags, etc.

### 1.1.36 0.3.1 (Oct 1, 2016) – what a wonderful week

Primarily bugfixes but also a number of enhancements and core refactorings

#### Fixes

- do not build manpages and examples during installation to avoid problems with possibly previously outdated dependencies
- `install` can be called on already installed dataset (with `-r` or `-g`)

## Enhancements and new features

- complete overhaul of datalad configuration settings handling (see [Configuration documentation](#)), so majority of the environment. Now uses git format and stores persistent configuration settings under `.datalad/config` and local within `.git/config` variables we have used were renamed to match configuration names
- `create-sibling` does not now by default upload web front-end
- `export` command with a plug-in interface and `tarball` plugin to export datasets
- in Python, `.api` functions with rendering of results in command line got a `_`-suffixed sibling, which would render results as well in Python as well (e.g., using `search_` instead of `search` would also render results, not only output them back as Python objects)
- `get`
  - `--jobs` option (passed to `annex get`) for parallel downloads
  - total and per-download (with `git-annex >= 6.20160923`) progress bars (note that if content is to be obtained from an archive, no progress will be reported yet)
- `install --reckless` mode option
- `search`
  - highlights locations and fieldmaps for better readability
  - supports `-d^` or `-d///` to point to top-most or centrally installed meta-datasets
  - “complete” paths to the datasets are reported now
  - `-s` option to specify which fields (only) to search
- various enhancements and small fixes to [meta-data](#) handling, `ls`, custom remotes, code-base formatting, downloaders, etc
- completely switched to `tqdm` library (`progressbar` is no longer used/supported)

### 1.1.37 0.3 (Sep 23, 2016) – winter is coming

Lots of everything, including but not limited to

- enhanced index viewer, as the one on <http://datasets.datalad.org>
- initial new data providers support: [Kaggle](#), [BALSA](#), [NDA](#), [NITRC](#)
- initial [meta-data support and management](#)
- new and/or improved crawler pipelines for [BALSA](#), [CRCNS](#), [OpenfMRI](#)
- refactored `install` command, now with separate `get`
- some other commands renaming/refactoring (e.g., `create-sibling`)
- datalad `search` would give you an option to install datalad’s super-dataset under `~/datalad` if ran outside of a dataset

### 0.2.3 (Jun 28, 2016) – busy OHBM

New features and bugfix release

- support of `///` urls to point to <http://datasets.datalad.org>
- variety of fixes and enhancements throughout

## **0.2.2 (Jun 20, 2016) – OHBM we are coming!**

New feature and bugfix release

- greatly improved documentation
- publish command API RFinng allows for custom options to annex, and uses `-to REMOTE` for consistent with annex invocation
- variety of fixes and enhancements throughout

## **0.2.1 (Jun 10, 2016)**

- variety of fixes and enhancements throughout

## **1.1.38 0.2 (May 20, 2016)**

Major RFinng to switch from relying on `rdf` to `git` native submodules etc

## **1.1.39 0.1 (Oct 14, 2015)**

Release primarily focusing on interface functionality including initial publishing

## **1.2 Acknowledgments**

DataLad development is being performed as part of a US-German collaboration in computational neuroscience (CR-CNS) project “DataGit: converging catalogues, warehouses, and deployment logistics into a federated ‘data distribution’” (Halchenko/Hanke), co-funded by the US National Science Foundation (NSF 1429999) and the German Federal Ministry of Education and Research (BMBF 01GQ1411). Additional support is provided by the German federal state of Saxony-Anhalt and the European Regional Development Fund (ERDF), Project: [Center for Behavioral Brain Sciences, Imaging Platform](#)

DataLad is built atop the `git-annex` software that is being developed and maintained by [Joey Hess](#).

## **1.3 Publications**

Further conceptual and technical information on DataLad, and applications built on DataLad, are available from the publications listed below.

### **The best of both worlds: Using semantic web with JSOB-LD. An example with NIDM Results & DataLad [poster]**

- Camille Maumet, Satrajit Ghosh, Yaroslav O. Halchenko, Dorota Jarecka, Nolan Nichols, Jean-Baptist Poline, Michael Hanke

### **One thing to bind them all: A complete raw data structure for auto-generation of BIDS datasets [poster]**

- Benjamin Poldrack, Kyle Meyer, Yaroslav O. Halchenko, Michael Hanke

### **Fantastic containers and how to tame them [poster]**

- Yaroslav O. Halchenko, Kyle Meyer, Matt Travers, Dorota Jarecka, Satrajit Ghosh, Jakub Kaczmarzyk, Michael Hanke



**YODA: YODA’s Organigram on Data Analysis [poster]**

- An outline of a simple approach to structuring and conducting data analyses that aims to tightly connect all their essential ingredients: data, code, and computational environments in a transparent, modular, accountable, and practical way.
- Michael Hanke, Kyle A. Meyer, Matteo Visconti di Oleggio Castello, Benjamin Poldrack, Yaroslav O. Halchenko
- F1000Research 2018, 7:1965 (<https://doi.org/10.7490/f1000research.1116363.1>)

**Go FAIR with DataLad [talk]**

- On DataLad’s capabilities to create and maintain Findable, Accessible, Interoperable, and Re-Usable (FAIR) resources.
- Michael Hanke, Yaroslav O. Halchenko
- Bernstein Conference 2018 workshop: Practical approaches to research data management and reproducibility ([slides](#))
- OpenNeuro kick-off meeting, 2018, Stanford ([slide sources](#))

## 1.4 Concepts and technologies

### 1.4.1 Background and motivation

#### Vision

Data is at the core of science, and unobstructed access promotes scientific discovery through collaboration between data producers and consumers. The last years have seen dramatic improvements in availability of data resources for collaborative research, and new data providers are becoming available all the time.

However, despite the increased availability of data, their accessibility is far from being optimal. Potential consumers of these public datasets have to manually browse various disconnected warehouses with heterogeneous interfaces. Once obtained, data is disconnected from its origin and data versioning is often ad-hoc or completely absent. If data consumers can be reliably informed about data updates at all, review of changes is difficult, and re-deployment is tedious and error-prone. This leads to wasteful friction caused by outdated or faulty data.

The vision for this project is to transform the state of data-sharing and collaborative work by providing uniform access to available datasets – independent of hosting solutions or authentication schemes – with reliable versioning and versatile deployment logistics. This is achieved by means of a *dataset* handle, a lightweight representation of a dataset that is capable of tracking the identity and location of a dataset’s content as well as carry meta-data. Together with associated software tools, scientists are able to obtain, use, extend, and share datasets (or parts thereof) in a way that is traceable back to the original data producer and is therefore capable of establishing a strong connection between data consumers and the evolution of a dataset by future extension or error correction.

Moreover, DataLad aims to provide all tools necessary to create and publish *data distributions* — an analog to software distributions or app-stores that provide logistics middleware for software deployment. Scientific communities can use these tools to gather, curate, and make publicly available specialized collections of datasets for specific research topics or data modalities. All of this is possible by leveraging existing data sharing platforms and institutional resources without the need for funding extra infrastructure of duplicate storage. Specifically, this project aims to provide a comprehensive, extensible data distribution for neuroscientific datasets that is kept up-to-date by an automated service.

## Technological foundation: git-annex

The outlined task is not unique to the problem of data-sharing in science. Logistical challenges such as delivering data, long-term storage and archiving, identity tracking, and synchronization between multiple sites are rather common. Consequently, solutions have been developed in other contexts that can be adapted to benefit scientific data-sharing.

The closest match is the software tool [git-annex](#). It combines the features of the distributed version control system (dVCS) [Git](#) — a technology that has revolutionized collaborative software development — with versatile data access and delivery logistics. Git-annex was originally developed to address use cases such as managing a collection of family pictures at home. With git-annex, any family member can obtain an individual copy of such a picture library — the *annex*. The annex in this example is essentially an image repository that presents individual pictures to users as files in a single directory structure, even though the actual image file contents may be distributed across multiple locations, including a home-server, cloud-storage, or even off-line media such as external hard-drives.

Git-annex provides functionality to obtain file contents upon request and can prompt users to make particular storage devices available when needed (e.g. a backup hard-drive kept in a fire-proof compartment). Git-annex can also remove files from a local copy of that image repository, for example to free up space on a laptop, while ensuring a configurable level of data redundancy across all known storage locations. Lastly, git-annex is able to synchronize the content of multiple distributed copies of this image repository, for example in order to incorporate images added with the git-annex on the laptop of another family member. It is important to note that git-annex is agnostic of the actual file types and is not limited to images.

We believe that the approach to data logistics taken by git-annex and the functionality it is currently providing are an ideal middleware for scientific data-sharing. Its data repository model *annex* readily provides the majority of principal features needed for a dataset handle such as history recording, identity tracking, and item-based resource locators. Consequently, instead of a from-scratch development, required features, such as dedicated support for existing data-sharing portals and dataset meta-information, can be added to a working solution that is already in production for several years. As a result, DataLad focuses on the expansion of git-annex's functionality and the development of tools that build atop Git and git-annex and enable the creation, management, use, and publication of dataset handles and collections thereof.

## Objective

Building atop git-annex, DataLad aims to provide a single, uniform interface to access data from various data-sharing initiatives and data providers, and functionality to create, deliver, update, and share datasets for individuals and portal maintainers. As a command-line tool, it provides an abstraction layer for the underlying Git-based middleware implementing the actual data logistics, and serves as a foundation for other future user front-ends, such as a web-interface.

### 1.4.2 Delineation from related solutions

To our knowledge, there is no other effort with a scope as broad as DataLad's. DataLad aims to unify access to vast arrays of (scientific) data in a domain and data modality agnostic fashion with as few and universally available software dependencies as possible.

The most comparable project regarding the idea of federating access to various data providers is the [iRODS-based INCF Dataspace](#) project. IRODS is a powerful, NSF-supported framework, but it requires non-trivial deployment and management procedures. As a representative of *data grid* technology, it is more suitable for an institutional deployment, as data access, authentication, permission management, and versioning are complex and not-feasible to be performed directly by researchers. DataLad on the other hand federates institutionally hosted data, but in addition enables individual researchers and small labs to contribute datasets to the federation with minimal cost and without the need for centralized coordination and permission management.

## Data catalogs

Existing data-portals, such as [DataDryad](#), or domain-specific ones (e.g. [Human Connectome](#), [OpenfMRI](#)), concentrate on collecting, cataloging, and making data available. They offer an abstraction from local data management peculiarities (organization, updates, sharing). Ad-hoc collections of pointers to available data, such as [reddit datasets](#) and [Inside-R datasets](#), do not provide any unified interface to assemble and manage such data. Data portals can be used as seed information and data providers for DataLad. These portals could in turn adopt DataLad to expose readily usable data collections via a federated infrastructure.

## Data delivery/management middleware

Even though there are projects to manage data directly with dVCS (e.g. Git), such as the [Rdatasets Git repository](#) this approach does not scale, for example to the amount of data typically observed in a scientific context. DataLad uses [git-annex](#) to support managing large amounts of data with Git, while avoiding the scalability issues of putting data directly into Git repositories.

In scientific software development, frequently using Git for source code management, many projects are also confronted with the problem of managing large data arrays needed, for example, for software testing. An exemplar project is [ITK Data](#) which is conceptually similar to git-annex: data content is referenced by unique keys (checksums), which are made redundantly available through multiple remote key-store farms and can be obtained using specialized functionality in the CMake software build system. However, the scope of this project is limited to software QA, and only provides an ad-hoc collection of guidelines and supporting scripts.

The git-annex website provides a [comparison](#) of Git-annex to other available distributed data management tools, such as [git-media](#), [git-fat](#), and others. None of the alternative frameworks provides all of the features of git-annex, such as integration with native Git workflows, distributed redundant storage, and partial checkouts in one project. Additional features of git-annex which are not necessarily needed by DataLad (git-annex assistant, encryption support, etc.) make it even more appealing for extended coverage of numerous scenarios. Moreover, neither of the alternative solutions has already reached a maturity, availability, and level of adoption that would be comparable to that of git-annex.

## Git/Git-annex/DataLad

Although it is possible, and intended, to use DataLad without ever invoking git or git-annex commands directly, it is useful to appreciate that DataLad is build atop of very flexible and powerful tools. Knowing basics of git and git-annex in addition to DataLad helps to not only make better use of DataLad but also to enable more advanced and more efficient data management scenarios. DataLad makes use of lower-level configuration and data structures as much as possible. Consequently, it is possible to manipulate DataLad datasets with low-level tools if needed. Moreover, DataLad datasets are compatible with tools and services designed to work with plain Git repositories, such as the popular [GitHub](#) service.

To better illustrate the different scopes, the following table provides an overview of the features that are contributed by each software technology layer.

| Feature   | Git           | Git-annex  | DataLad   |
|---|---------------|------------|-----------|
| Version control (text, code)  | ✓             | ✓ can mix  | ✓ can mix |
| Version control (binary data)   | (not advised) | ✓          | ✓         |
| Auto-crawling available resources   |               | ✓RSS feeds | ✓flexible |
| Unified dataset handling  |               |            | ✓         |
| <ul style="list-style-type: none"> <li>• recursive operation on datasets</li> </ul>               |               |            | ✓         |
| <ul style="list-style-type: none"> <li>• seamless operation across datasets boundaries</li> </ul> |               |            | ✓         |
| <ul style="list-style-type: none"> <li>• meta-data support</li> </ul>                             |               | ✓per-file  | ✓         |
| <ul style="list-style-type: none"> <li>• meta-data aggregation</li> </ul>                         |               |            | ✓flexible |
| Unified authentication interface  |               |            | ✓         |

### 1.4.3 Basic principles

DataLad is designed to be used both as a command-line tool, and as a Python module. The sections *Command line reference* and *Python module reference* provide detailed description of the commands and functions of the two interfaces. This section presents common concepts. Although examples will frequently be presented using command line interface commands, all functionality with identically named functions and options are available through Python API as well.

#### Datasets

A DataLad *dataset* is a Git repository that may or may not have a data *annex* that is used to manage data referenced in a dataset. In practice, most DataLad datasets will come with an annex.

#### Types of IDs used in datasets

Four types of unique identifiers are used by DataLad to enable identification of different aspects of datasets and their components.

**Dataset ID** A UUID that identifies a dataset as a whole across its entire history and flavors. This ID is stored in a dataset’s own configuration file (`<dataset root>/.datalad/config`) under the configuration key `datalad.dataset.id`. As this configuration is stored in a file that is part of the Git history of a dataset, this ID is identical for all “clones” of a dataset and across all its versions. If the purpose or scope of a dataset changes enough to warrant a new dataset ID, it can be changed by altering the dataset configuration setting.

**Annex ID** A UUID assigned to an annex of each individual clone of a dataset repository. Git-annex uses this UUID to track file content availability information. The UUID is available under the configuration key `annex.uuid`

and is stored in the configuration file of a local clone (<dataset root>/`.git/config`). A single dataset instance (i.e. clone) can only have a single annex UUID, but a dataset with multiple clones will have multiple annex UUIDs.

**Commit ID** A Git hexsha or tag that identifies a version of a dataset. This ID uniquely identifies the content and history of a dataset up to its present state. As the dataset history also includes the dataset ID, a commit ID of a DataLad dataset is unique to a particular dataset.

**Content ID** Git-annex key (typically a checksum) assigned to the content of a file in a dataset's annex. The checksum reflects the content of a file, not its name. Hence the content of multiple identical files in a single (or across) dataset(s) will have the same checksum. Content IDs are managed by Git-annex in a dedicated `annex` branch of the dataset's Git repository.

## Dataset nesting

Datasets can contain other datasets (*subdatasets*), which can in turn contain subdatasets, and so on. There is no limit to the depth of nesting datasets. Each dataset in such a hierarchy has its own annex and its own history. The parent or *superdataset* only tracks the specific state of a subdataset, and information on where it can be obtained. This is a powerful yet lightweight mechanism for combining multiple individual datasets for a specific purpose, such as the combination of source code repositories with other resources for a tailored application. In many cases DataLad can work with a hierarchy of datasets just as if it were a single dataset. Here is a demo:

```
~ % datalad create demo
[INFO ] Creating a new annex repo at /demo/demo
create(ok): /demo/demo (dataset)
~ % cd demo
```

A DataLad dataset is just a Git repo with some initial configuration

```
~/demo % git log --oneline
472e34b (HEAD -> master) [DATALAD] new dataset
f968257 [DATALAD] Set default backend for all files to be MD5E
```

We can generate nested datasets, by telling DataLad to register a new dataset in a parent dataset

```
~/demo % datalad create -d . sub1
[INFO ] Creating a new annex repo at /demo/demo/sub1
add(ok): sub1 (dataset) [added new subdataset]
add(notneeded): sub1 (dataset) [nothing to add from /demo/demo/sub1]
add(notneeded): .gitmodules (file) [already included in the dataset]
save(ok): /demo/demo (dataset)
create(ok): sub1 (dataset)
action summary:
  add (notneeded: 2, ok: 1)
  create (ok: 1)
  save (ok: 1)
```

A subdataset is nothing more than regular Git submodule

```
~/demo % git submodule
5f0cddf2026e3fb4864139f27e7415fd72c7d4d0 sub1 (heads/master)
```

Of course subdatasets can be nested

```
~/demo % datalad create -d . sub1/justadir/sub2
[INFO ] Creating a new annex repo at /demo/demo/sub1/justadir/sub2
```

(continues on next page)

(continued from previous page)

```
add(ok): sub1/justadir/sub2 (dataset) [added new subdataset]
add(notneeded): sub1/justadir/sub2 (dataset) [nothing to add from /demo/demo/sub1/
↳justadir/sub2]
add(notneeded): sub1/.gitmodules (file) [already included in the dataset]
add(notneeded): sub1 (dataset) [already known subdataset]
save(ok): /demo/demo/sub1 (dataset)
save(ok): /demo/demo (dataset)
create(ok): sub1/justadir/sub2 (dataset)
action summary:
  add (notneeded: 3, ok: 1)
  create (ok: 1)
  save (ok: 2)
```

Unlike Git, DataLad automatically takes care of committing all changes associated with the added subdataset up to the given parent dataset

```
~/demo % git status
On branch master
nothing to commit, working tree clean
```

Let's create some content in the deepest subdataset

```
~/demo % mkdir sub1/justadir/sub2/anotherdir
~/demo % touch sub1/justadir/sub2/anotherdir/afile
```

Git can only tell us that something underneath the top-most subdataset was modified

```
~/demo % git status
On branch master
Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git checkout -- <file>..." to discard changes in working directory)
  (commit or discard the untracked or modified content in submodules)

        modified:   sub1 (untracked content)

no changes added to commit (use "git add" and/or "git commit -a")
```

DataLad saves us from further investigation

```
~/demo % datalad diff -r
  modified(dataset): sub1
  modified(dataset): sub1/justadir/sub2
  untracked(directory): sub1/justadir/sub2/anotherdir
```

Like Git, it can report individual untracked files, but also across repository boundaries

```
~/demo % datalad diff -r --report-untracked all
  modified(dataset): sub1
  modified(dataset): sub1/justadir/sub2
  untracked(file): sub1/justadir/sub2/anotherdir/afile
```

Adding this new content with Git or git-annex would be an exercise

```
~/demo % git add sub1/justadir/sub2/anotherdir/afile
fatal: Pathspec 'sub1/justadir/sub2/anotherdir/afile' is in submodule 'sub1'
```

DataLad does not require users to determine the correct repository in the tree

```
~/demo % datalad add -d . sub1/justadir/sub2/anotherdir/afile
add(ok): sub1/justadir/sub2/anotherdir/afile (file)
save(ok): /demo/demo/sub1/justadir/sub2 (dataset)
save(ok): /demo/demo/sub1 (dataset)
save(ok): /demo/demo (dataset)
action summary:
  add (ok: 1)
  save (ok: 3)
```

Again, all associated changes in the entire dataset tree, up to the given parent dataset, were committed

```
~/demo % git status
On branch master
nothing to commit, working tree clean
```

DataLad's 'diff' is able to report the changes from these related commits throughout the repository tree

```
~/demo % datalad diff --revision @~1 -r
modified(dataset): sub1
modified(dataset): sub1/justadir/sub2
added(file): sub1/justadir/sub2/anotherdir/afile
```

## Dataset collections

A superdataset can also be seen as a curated collection of datasets, for example, for a certain data modality, a field of science, a certain author, or from one project (maybe the resource for a movie production). This lightweight coupling between super and subdatasets enables scenarios where individual datasets are maintained by a disjoint set of people, and the dataset collection itself can be curated by a completely independent entity. Any individual dataset can be part of any number of such collections.

Benefiting from Git's support for workflows based on decentralized "clones" of a repository, DataLad's datasets can be (re-)published to a new location without losing the connection between the "original" and the new "copy". This is extremely useful for collaborative work, but also in more mundane scenarios such as data backup, or temporary deployment of a dataset on a compute cluster, or in the cloud. Using git-annex, data can also get synchronized across different locations of a dataset (*siblings* in DataLad terminology). Using metadata tags, it is even possible to configure different levels of desired data redundancy across the network of dataset, or to prevent publication of sensitive data to publicly accessible repositories. Individual datasets in a hierarchy of (sub)datasets need not be stored at the same location. Continuing with an earlier example, it is possible to post a curated collection of datasets, as a superdataset, on Github, while the actual datasets live on different servers all around the world.

## Basic command line usage

All of DataLad's functionality is available through a single command: *datalad*

Running the *datalad* command without any arguments, gives a summary of basic options, and a list of available sub-commands.

```
~ % datalad
usage: datalad [-h] [-l LEVEL] [--pbs-runner {condor}] [-C PATH] [--version]
              [--dbg] [--idbg] [-c KEY=VALUE]
              [-f {default,json,json_pp,tailored,'<template>'}]
              [--report-status {success,failure,ok,notneeded,impossible,error}]
```

(continues on next page)

(continued from previous page)

```

        [--report-type {dataset,file}]
        [--on-failure {ignore,continue,stop}] [--cmd]
        {create,install,get,publish,uninstall,drop,remove,update,create-
↪sibling,create-sibling-github,unlock,save,search,metadata,aggregate-metadata,test,
↪ls,clean,add-archive-content,download-url,run,rerun,addurls,export-archive,extract-
↪metadata,export-to-figshare,no-annex,wtf,add-readme,annotate-paths,clone,create-
↪test-dataset,diff,siblings,sshrun,subdatasets}
        ...
[ERROR ] Please specify the command
~ % #

```

More comprehensive information is available via the `-help` long-option (we will truncate the output here)

```

~ % datalad --help | head -n20
Usage: datalad [global-opts] command [command-opts]

DataLad provides a unified data distribution with the convenience of git-annex
repositories as a backend. DataLad command line tools allow to manipulate
(obtain, create, update, publish, etc.) datasets and their collections.

*Commands for dataset operations*

create
    Create a new dataset from scratch
install
    Install a dataset from a (remote) source
get
    Get any dataset content (files/directories/subdatasets)
publish
    Publish a dataset to a known sibling
uninstall
    Uninstall subdatasets

```

Getting information on any of the available sub commands works in the same way – just pass `-help` AFTER the sub-command (output again truncated)

```

~ % datalad create --help | head -n20
Usage: datalad create [-h] [-f] [-D DESCRIPTION] [-d PATH] [--no-annex]
        [--nosave] [--annex-version ANNEX_VERSION]
        [--annex-backend ANNEX_BACKEND]
        [--native-metadata-type LABEL] [--shared-access MODE]
        [--git-opts STRING] [--annex-opts STRING]
        [--annex-init-opts STRING] [--text-no-annex]
        [PATH]

Create a new dataset from scratch.

This command initializes a new dataset at a given location, or the
current directory. The new dataset can optionally be registered in an
existing superdataset (the new dataset's path needs to be located
within the superdataset for that, and the superdataset needs to be given
explicitly). It is recommended to provide a brief description to label
the dataset's nature *and* location, e.g. "Michael's music on black
laptop". This helps humans to identify data locations in distributed
scenarios. By default an identifier comprised of user and machine name,
plus path will be generated.

```



## API principles

You can use DataLad's `install` command to download datasets. The command accepts URLs of different protocols (`http`, `ssh`) as an argument. Nevertheless, the easiest way to obtain a first dataset is downloading the default *superdataset* from <http://datasets.datalad.org/> using a shortcut.

### Downloading DataLad's default superdataset

<http://datasets.datalad.org> provides a super-dataset consisting of datasets from various portals and sites. Many of them were crawled, and periodically updated, using `datalad-crawler` extension. The argument `///` can be used as a shortcut that points to the superdataset located at <http://datasets.datalad.org/>. Here are three common examples in command line notation:

```
datalad install /// installs this superdataset (metadata without subdatasets) in a datasets.datalad.org/ subdirectory under the current directory
```

```
datalad install -r ///openfmri installs the openfmri superdataset into an openfmri/ subdirectory. Additionally, the -r flag recursively downloads all metadata of datasets available from http://openfmri.org as subdatasets into the openfmri/ subdirectory
```

```
datalad install -g -J3 -r ///labs/haxby installs the superdataset of datasets released by the lab of Dr. James V. Haxby and all subdatasets' metadata. The -g flag indicates getting the actual data, too. It does so by using 3 parallel download processes (-J3 flag).
```

`datalad search` command, if ran outside of any dataset, will install this default superdataset under a path specified in `datalad.locations.default-dataset` *configuration* variable (by default `$HOME/datalad`).

### Downloading datasets via http

In most places where DataLad accepts URLs as arguments these URLs can be regular `http` or `https` protocol URLs. For example:

```
datalad install https://github.com/psychoinformatics-de/studyforrest-data-phase2.git
```

### Downloading datasets via ssh

DataLad also supports SSH URLs, such as `ssh://me@localhost/path`.

```
datalad install ssh://me@localhost/path
```

Finally, DataLad supports SSH login style resource identifiers, such as `me@localhost:/path`.

```
datalad install me@localhost:/path
```

### `-dataset` argument

All commands which operate with/on datasets (practically all commands) have a `dataset` argument (`-d` or `--dataset` for the command line API) which takes a path to the dataset that the command should operate on. If a dataset is identified this way then any relative path that is provided as an argument to the command will be interpreted as being relative to the topmost directory of that dataset. If no dataset argument is provided, relative paths are considered to be relative to the current directory.

There are also some useful pre-defined “shortcut” values for dataset arguments:

`///` refers to the “default” dataset located under `$HOME/datalad/`. So running `datalad install -d/// crcns` will install the `crcns` subdataset under `$HOME/datalad/crcns`. This is the same as running `datalad install $HOME/datalad/crcns`.

`^` topmost superdataset containing the dataset the current directory is part of. For example, if you are in `$HOME/datalad/openfmri/ds000001/sub-01` and want to search metadata of the entire superdataset you are under (in this case `///`), run `datalad search -d^ [something to search]`.

`^.` the dataset the current directory is part of.

### Commands *install* vs *get*

The `install` and `get` commands might seem confusingly similar at first. Both of them could be used to install any number of subdatasets, and fetch content of the data files. Differences lie primarily in their default behaviour and outputs, and thus intended use. Both `install` and `get` take local paths as their arguments, but their default behavior and output might differ;

- **install** primarily operates and reports at the level of **datasets**, and returns as a result dataset(s) which either were just installed, or were installed previously already under specified locations. So result should be the same if the same `install` command ran twice on the same datasets. It **does not fetch** data files by default
- **get** primarily operates at the level of **paths** (datasets, directories, and/or files). As a result it returns only what was installed (datasets) or fetched (files). So result of rerunning the same `get` command should report that nothing new was installed or fetched. It **fetches** data files by default.

In how both commands operate on provided paths, it could be said that `install == get -n`, and `install -g == get`. But `install` also has ability to install new datasets from remote locations given their URLs (e.g., `http://datasets.datalad.org/` for our super-dataset) and SSH targets (e.g., `[login@]host:path`) if they are provided as the argument to its call or explicitly as `--source` option. If `datalad install --source URL DESTINATION` (command line example) is used, then dataset from URL gets installed under `PATH`. In case of `datalad install URL` invocation, `PATH` is taken from the last name within URL similar to how `git clone` does it. If former specification allows to specify only a single URL and a `PATH` at a time, later one can take multiple remote locations from which datasets could be installed.

So, as a rule of thumb – if you want to install from external URL or fetch a sub-dataset without downloading data files stored under annex – use `install`. In Python API `install` is also to be used when you want to receive in output the corresponding Dataset object to operate on, and be able to use it even if you rerun the script. In all other cases, use `get`.

## 1.4.4 Metadata

### Overview

DataLad has built-in, modular, and extensible support for metadata in various formats. Metadata is extracted from a dataset and its content by one or more extractors that have to be enabled in a dataset’s configuration. Extractors yield metadata in a **JSON-LD**-like structure that can be arbitrarily complex and deeply nested. Metadata from each extractor is kept unmodified, unmangled, and separate from metadata of other extractors. This design enables tailored applications using particular metadata that can use Datalad as a content-agnostic aggregation and transport layer without being limited or impacted by other metadata sources and schemas.

Extracted metadata is stored in a dataset in (compressed) files using a JSON stream format, separately for metadata describing a dataset as a whole, and metadata describing individual files in a dataset. This limits the amount of metadata that has to be obtained and processed for applications that do not require all available metadata.

DataLad provides a content-agnostic metadata aggregation mechanism that stores metadata of sub-datasets (with arbitrary nesting levels) in a superdataset, where it can then be queried without having the subdatasets locally present.

Lastly, DataLad comes with a *search* command that enable metadata queries via a flexible query language. However, alternative applications for metadata queries (e.g. graph-based queries) can be built on DataLad, by requesting a complete or partial dump of aggregated metadata available in a dataset.

## Supported metadata sources

This following sections provide an overview of included metadata extractors for particular types of data structures and file formats. Note that *DataLad extension packages*, such as the *neuroimaging extension*, can provide additional extractors for particular domains and formats.

Only *annex* and *datalad\_core* extractors are enabled by default. Any additional metadata extractor should be enabled by setting the `datalad.metadata.nativetype` *configuration* variable via the `git config` command or by editing `.datalad/config` directly. For example, `git config -f .datalad/config --add datalad.metadata.nativetype audio` would add *audio* metadata extractor to the list.

## Annex metadata (annex)

Content tracked by git-annex can have associated *metadata records*. From DataLad's perspective, git-annex metadata is just another source of metadata that can be extracted and aggregated.

You can use the *git-annex metadata* command to assign git-annex metadata. And, if you have a table or records that contain data sources and metadata, you can use *datalad addurls* to quickly populate a dataset with files and associated git-annex metadata. (<http://labs/openneurolab/metasearch> is an example of such a dataset.)

## Pros of git-annex level metadata

- Many git-annex commands, such as *git-annex get* and *git-annex copy*, can use metadata to decide which files (keys) to operate on, making it possible to automate file (re)distribution based on their metadata annotation
- Assigned metadata is available for use by git-annex right away without requiring any additional “aggregation” step
- *git-annex view* can be used to quickly generate completely new layouts of the repository solely based on the metadata fields associated with the files

## Cons of git-annex level metadata

- Metadata fields are actually stored per git-annex key rather than per file. If multiple files contain the same content, metadata will be shared among them.
- Files whose content is tracked directly by git cannot have git-annex metadata assigned.
- No per repository/directory metadata, and no mechanism to use/aggregate metadata from sub-datasets
- Field names cannot contain some symbols, such as ‘:’
- Metadata is stored within the *git-annex* branch, so it is distributed across all clones of the dataset, making it hard to scale for large metadata sizes or to work with sensitive metadata (not intended to be redistributed)
- It is a generic storage with no prescribed vocabulary, making it very flexible but also requiring consistency and harmonization to make the stored metadata useful for search

### Example uses of git-annex metadata

#### Annotating files for different purposes

FreeSurfer project uses *git-annex* for managing their source code+data base within a single git/git-annex repository. Files necessary for different scenarios (deployment, testing) are annotated and can be fetched selectively for the scenario at hand.

#### Automating “non-distribution” of sensitive files

In the [ReproIn](#) framework for automated conversion of BIDS dataset and in some manually prepared datasets (such as [///labs/gobbini/famface/data](#) and [///labs/haxby/raiders](#)), we annotated materials that must not be publicly shared with a git-annex metadata field *distribution-restrictions*. We used the following of values to describe why any particular file (content) should not be redistributed:

- **sensitive** - files which potentially contain participant sensitive information, such as non-defaced anatomicals
- **proprietary** - files which contain proprietary data, which we have no permissions to share (e.g., movie video files)

Having annotated files this way, we could instruct git-annex to *publish* all but those restricted files to our server: *git annex wanted datalad-public “not metadata=distribution-restrictions=\*”*.

#### Flexible directory layout

If you are maintaining a collection of music files or PDFs for the lab, you may want to display the files in an alternative or filtered hierarchy. *git-annex view* could be of help. Example:

```
datalad install ///labs/openneurolab/metasearch
cd metasearch
git annex view sex=* handedness=ambidextrous
```

would give you two directories (Male, Female) with only the files belonging to ambidextrous subjects.

#### Various audio file formats (audio)

This extractor uses the [mutagen](#) package to extract essential metadata from a range of audio file formats. For the most common metadata properties a constrained vocabulary, based on the [Music Ontology](#) is employed.

#### datacite.org compliant datasets (datacite)

This extractor can handle dataset-level metadata following the [datacite.org](#) specification. No constrained vocabulary is identified at the moment.

#### Datalad’s internal metadata storage (datalad\_core)

This extractor can express Datalad’s internal metadata representation, such as the relationship of a super- and a sub-dataset. It uses DataLad’s own constrained vocabulary.

## RFC822-compliant metadata (datalad\_rfc822)

This is a custom metadata format, inspired by the standard used for Debian software packages that is particularly suited for manual entry. This format is a good choice when metadata describing a dataset as a whole cannot be obtained from some other structured format. The syntax is **RFC 822**-compliant. In other words: this is a text-based format that uses the syntax of email headers. Metadata must be placed in `DATASETROOT/.datalad/meta.rfc822` for this format.

Here is an example:

```
Name: myamazingdataset
Version: 1.0.0-rc3
Description: Basic summary
  A text with arbitrary length and content that can span multiple
  .
  paragraphs (this is a new one)
License: CC0
  The person who associated a work with this deed has dedicated the work to the
  public domain by waiving all of his or her rights to the work worldwide under
  copyright law, including all related and neighboring rights, to the extent
  allowed by law.
  .
  You can copy, modify, distribute and perform the work, even for commercial
  purposes, all without asking permission.
Homepage: http://example.com
Funding: Grandma's and Grandpa's support
Issue-Tracker: https://github.com/datalad/datalad/issues
Cite-As: Mike Author (2016). We made it. The breakthrough journal of unlikely
  events. 1, 23-453.
DOI: 10.0000/nothere.48421
```

The following fields are supported:

**Audience:** A description of the target audience of the dataset.

**Author:** A comma-delimited list of authors of the dataset, preferably in the format. `Firstname Lastname <Email Adress>`

**Cite-as:** Instructions on how to cite the dataset, or a structured citation.

**Description:** Description of the dataset as a whole. The first line should represent a compact short description with no more than 6-8 words.

**DOI:** A [digital object identifier](#) for the dataset.

**Funding:** Information on potential funding for the creation of the dataset and/or its content. This field can also be used to acknowledge non-monetary support.

**Homepage:** A URL to a project website for the dataset.

**Issue-tracker:** A URL to an issue tracker where known problems are documented and/or new reports can be submitted.

**License:** A description of the license or terms of use for the dataset. The first lines should contain a list of license labels (e.g. CC0, PDDL) for standard licenses, if possible. Full license texts or term descriptions can be included.

**Maintainer:** Can be used in addition and analog to `Author`, when authors (creators of the data) need to be distinguished from maintainers of the dataset.

**Name:** A short name for the dataset. It may be beneficial to avoid special characters, umlauts, spaces, etc. to enable widespread use of this name for URL, catalog keys, etc. in unmodified form.

**Version:** A version for the dataset. This should be in a format that is alphanumerically sortable and lead to a “greater” version for an update of a dataset.

Metadata keys used by this extractor are defined in DataLad’s own constrained vocabulary.

### Friction-less data packages (`frictionless_datapackage`)

DataLad has basic support for extraction of essential dataset-level metadata from `friction-less data packages` (`datapackage.json`). file. Metadata keys are constrained to DataLad’s own vocabulary.

### Exchangeable Image File Format (`exif`)

The extractor yields EXIF metadata from any compatible file. It uses the W3C EXIF vocabulary (<http://www.w3.org/2003/12/exif/ns/>).

### Various image/photo formats (`image`)

Standard image metadata is extractor using the `Pillow` package. Core metadata is available using an adhoc vocabulary defined by the extractor.

### Extensible Metadata Platform (`xmp`)

This extractor yields any XMP-compliant metadata from any supported file (e.g. PDFs, photos). XMP metadata uses fully qualified terms from standard vocabularies that are simply passed through by the extractor. At the moment metadata extraction from side-car files is not supported, but would be easy to add.

### Metadata aggregation and query

Metadata aggregation can be performed with the `aggregate-metadata` command. Aggregation is done for two interrelated but distinct reasons:

- Fast uniform metadata access, independent of local data availability
- Comprehensive data discovery without access to or knowledge of individual datasets

In an individual dataset, metadata aggregation engages any number of enabled metadata extractors to build a JSON-LD based metadata representation that is separate from the original data files. These metadata objects are added to the dataset and are tracked with the same mechanisms that are used for any other dataset content. Based on this metadata, DataLad can provide fast and uniform access to metadata for any dataset component (individual files, subdatasets, the whole dataset itself), based on the relative path of a component within a dataset (available via the `metadata` command). This extracted metadata can be kept or made available locally for any such query, even when it is impossible or undesirable to keep the associated data files around (e.g. due to size constraints).

For any superdataset (a dataset that contains subdatasets as components), aggregation can go one step further. In this case, aggregation imports extracted metadata from subdatasets into the superdataset to offer the just described query feature for any aggregated subdataset too. This works across any number of levels of nesting. For example, a subdataset that contains the aggregated metadata for eight other datasets (that might have never been available locally) can be aggregated into a local superdataset with all its metadata. In that superdataset, a DataLad user is then able to query information on any content of any subdataset, regardless of their actual availability. This principle also allows any user to install the superdataset from <http://datasets.datalad.org> and perform *local and offline* queries about any dataset available online from this server.

Besides full access to all aggregated metadata by path (via the *metadata* command), DataLad also comes with a *search* command that provides different search modes to query the entirety of the locally available metadata. Its capabilities include simple keyword searches as well as more complex queries using date ranges or logical conjunctions.

## Internal metadata representation

**Warning:** The information in this section is meant to provide insight into how DataLad structures extracted and aggregated metadata. However, this representation is not considered stable or part of the public API, hence these data should not be accessed directly. Instead, all metadata access should happen via the **metadata** API command.

A dataset's metadata is stored in the *.datalad/metadata* directory. This directory contains two main elements:

- a metadata inventory or catalog
- a store for metadata “objects”

## The metadata inventory

The inventory is kept in a JSON file, presently named `aggregate_v1.json`. It contains a single top-level dictionary/object. Each element in this dictionary represents one subdataset from which metadata has been extracted and aggregated into the dataset at hand. Keys in this dictionary are paths to the respective (sub)datasets (relative to the root of the dataset). If a dataset has no subdataset and metadata extraction was performed, the dictionary will only have a single element under the key `". "`.

Here is an excerpt of an inventory dictionary showing the record of the root dataset itself.

```
{
  ". ": {
    "content_info":
      "objects/0c/cn-b046b2c3a5e2b9c5599c980c7b5fab.xz",
    "datalad_version":
      "0.10.0.rc4.dev191",
    "dataset_info":
      "objects/0c/ds-b046b2c3a5e2b9c5599c980c7b5fab",
    "extractors": [
      "datalad_core",
      "annex",
      "bids",
      "nifti1"
    ],
    "id":
      "00ce405e-6589-11e8-b749-a0369fb55db0",
    "refcommit":
      "d170979ef33a82c67e6fefe3084b9fe7391b422b"
  },
}
```

The record of each dataset contains the following elements:

**id** The DataLad dataset UUID of the dataset metadata was extracted and aggregated from.

**refcommit** The SHA sum of the last metadata-relevant commit in the history of the dataset metadata was extracted from. Metadata-relevant commits are any commits that modify dataset content that is not exclusively concerning DataLad's own internal status and configuration.

**datalad\_version** The version string of the DataLad version that was used to perform the metadata extraction (not necessarily the metadata aggregation, as pre-extracted metadata can be aggregated from other superdatasets for a dataset that is itself not available locally).

**extractors** A list with the names of all enabled metadata extractors for this dataset. This list may include names for extractors that are provided by extensions, and may not be available for any given DataLad installation.

**content\_info, dataset\_info** Path to the object files containing the actual metadata on the dataset as a whole, and on individual files in a dataset (content). Paths are to be interpreted relative to the inventory file, and point to the metadata object store.

Read-access to the metadata inventory is available via the `metadata` command and its `--get-aggregates` option.

### The metadata object store

The object store holds the files containing dataset and content metadata for each aggregated dataset. The object store is located in `.datalad/metadata/objects`. However, this directory itself and the subdirectory structure within it have no significance, they are completely defined and exclusively discoverable via the `content_info` and `dataset_info` values in the metadata inventory records.

Metadata objects for datasets and content use a slightly different internal format. Both files could be either compressed (XZ) or uncompressed. Current practice uses compression for content metadata, but not for dataset metadata. Any metadata object file could be directly committed to Git, or it could be tracked via Git-annex. Reasons to choose one over the other could be file size, or privacy concerns.

Read-access to the metadata objects of dataset and individual files is available via the `metadata` command. Importantly, metadata can be requested

### Metadata objects for datasets

These files have a single top-level JSON object/dictionary as content. A JSON-LD `@content` field is used to assign a semantic markup to allow for programmatic interpretation of metadata as linked data. Any other top-level key identifies the name of a metadata extractor, and the value stored under this key represents the output of the corresponding extractor.

Structure and content of an extractor's output are unconstrained and completely up to the implementation of that particular extractor. Extractor can report additional JSON-LD context information (but there is no requirement).

The output of one extractor does not interfere or collide with the output of any other extractor.

### Metadata objects for content/file

In contrast to metadata objects for entire datasets, these files use a JSON stream format, i.e. one JSON object/dictionary per line (no surrounding list). This makes it possible to process the content line-by-line instead of having to load an entire files (with potentially millions of records).

The only other difference to dataset metadata objects is an additional top-level key `path` that identifies the relative path (relative to the root of its parent dataset) of the file the metadata record is associated with.

Otherwise, the extractor-specific metadata structure and content is unconstrained.



Content metadata objects tend to contain massively redundant information (e.g. a dataset with a thousand 12 megapixel images will report the identical resolution information a thousand times). Therefore, content metadata objects are by default XZ compressed – as this compressor is particularly capable discovering such redundancy and yield a very compact file size.

The reason for gathering all metadata into a single file across all content files and metadata extractors is to limit the impact on the performance of the underlying Git repository. Large superdataset could otherwise quickly grow into dimensions where tens of thousands of files would be required just to manage the metadata. Such a configuration would also limit the compatibility of DataLad datasets with constrained storage environments (think e.g. inode limits on super computers), as these files are tracked in Git and would therefore be present in any copy, regardless of whether metadata access is desired or not.

## Vocabulary

The following sections describe details and changes in the metadata specifications implemented in datalad.

### v2.0

- Current development version that will be released together with DataLad v0.10.

### v1.0

- Original implementation that did not really see the light of the day.

## 1.4.5 Customization and extension of functionality

DataLad provides numerous commands that cover many use cases. However, there will always be a demand for further customization or extensions of built-in functionality at a particular site, or for an individual user. DataLad addresses this need with two mechanisms:

- *Plugins*
- *Extension packages*

Plugins are a quick'n'dirty way to implement a single additional command with very little overhead. They are, however, not the method of choice for extending particular Datalad functionality, such as metadata extractor, or providing entire command suites for a specialized purpose. For all these scenarios extension packages are the recommended method.

## Plugins

A number of plugins are shipped with DataLad. This includes plugins which operate on a particular dataset, but also general functionality that can be used outside the context of a specific dataset. The following table provides an overview of plugins included in this DataLad release.

|                                 |   |
|---------------------------------|---|
| <code>add_readme</code>         | add a README file to a dataset                    |
| <code>addurls</code>            | Create and update a dataset from a list of URLs.  |
| <code>check_dates</code>        | Extension for checking dates within repositories. |
| <code>export_archive</code>     | export a dataset as a compressed TAR/ZIP archive  |
| <code>export_to_figshare</code> | export a dataset as a TAR/ZIP archive to figshare |

Continued on next page

Table 1 – continued from previous page

|                 |   |
|-----------------|---|
| <i>no_annex</i> | configure which dataset parts to never put in the annex |
| <i>wtf</i>      | provide information about this DataLad installation     |

## **datalad.plugin.add\_readme**

add a README file to a dataset

**class** `datalad.plugin.add_readme.AddReadme`

Bases: `datalad.interface.base.Interface`

Add basic information about DataLad datasets to a README file

The README file is added to the dataset and the addition is saved in the dataset.

**class** `EnsureChoice (*values)`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is element of a set of possible values

**long\_description()**

**short\_description()**

**class** `EnsureDataset`

Bases: `datalad.support.constraints.Constraint`

Despite its name, this constraint does not actually ensure that the argument is a valid dataset, because for procedural reasons this would typically duplicate subsequent checks and processing. However, it can be used to achieve uniform documentation of *dataset* arguments.

**long\_description()**

**short\_description()**

**class** `EnsureNone`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is of value *None*

**long\_description()**

**short\_description()**

**class** `EnsureStr (min_len=0)`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is a string.

No automatic conversion is attempted.

**long\_description()**

**short\_description()**

**class** `Parameter (constraints=None, doc=None, args=None, **kwargs)`

Bases: `object`

This class shall serve as a representation of a parameter.

**get\_autodoc** (*name, indent=' ', width=70, default=None, has\_default=False*)

Docstring for the parameter to be used in lists of parameters

**Returns**

**Return type** string or list of strings (if indent is None)

**datasetmethod** (*name=None, dataset\_argname='dataset'*)

**eval\_results** ()

Decorator for return value evaluation of datalad commands.

Note, this decorator is only compatible with commands that return status dict sequences!

Two basic modes of operation are supported: 1) “generator mode” that *yields* individual results, and 2) “list mode” that returns a sequence of results. The behavior can be selected via the kwarg *return\_type*. Default is “list mode”.

This decorator implements common functionality for result rendering/output, error detection/handling, and logging.

Result rendering/output can be triggered via the *datalad.api.result-renderer* configuration variable, or the *result\_renderer* keyword argument of each decorated command. Supported modes are: ‘default’ (one line per result with action, status, path, and an optional message); ‘json’ (one object per result, like git-annex), ‘json\_pp’ (like ‘json’, but pretty-printed spanning multiple lines), ‘tailored’ custom output formatting provided by each command class (if any).

Error detection works by inspecting the *status* item of all result dictionaries. Any occurrence of a status other than ‘ok’ or ‘notneeded’ will cause an *IncompleteResultsError* exception to be raised that carries the failed actions’ status dictionaries in its *failed* attribute.

Status messages will be logged automatically, by default the following association of result status and log channel will be used: ‘ok’ (debug), ‘notneeded’ (debug), ‘impossible’ (warning), ‘error’ (error). Logger instances included in the results are used to capture the origin of a status report.

**Parameters** **func** (*function*) – `__call__` method of a subclass of `Interface`, i.e. a datalad command definition

## datalad.plugin.addurls

Create and update a dataset from a list of URLs.

**class** `datalad.plugin.addurls.Addurls`

Bases: `datalad.interface.base.Interface`

Create and update a dataset from a list of URLs.

*Format specification*

Several arguments take format strings. These are similar to normal Python format strings where the names from *URL-FILE* (column names for a CSV or properties for JSON) are available as placeholders. If *URL-FILE* is a CSV file, a positional index can also be used (i.e., “{0}” for the first column). Note that a placeholder cannot contain a ‘:’ or ‘!’.

In addition, the *FILENAME-FORMAT* arguments has a few special placeholders.

- `_reindex`

The constructed file names must be unique across all fields rows. To avoid collisions, the special placeholder “\_reindex” can be added to the formatter. Its value will start at 0 and increment every time a file name repeats.

- `_url_hostname, _urlN, _url_basename*`

Various parts of the formatted URL are available. Take “[http://datalad.org/asciicast/seamless\\_nested\\_repos.sh](http://datalad.org/asciicast/seamless_nested_repos.sh)” as an example.

“datalad.org” is stored as “\_url\_hostname”. Components of the URL’s path can be referenced as “\_urlN”. “\_url0” and “\_url1” would map to “asciicast” and “seamless\_nested\_repos.sh”, respectively. The final part of the path is also available as “\_url\_basename”.

This name is broken down further. “\_url\_basename\_root” and “\_url\_basename\_ext” provide access to the root name and extension. These values are similar to the result of `os.path.splitext`, but, in the case of multiple periods, the extension is identified using the same length heuristic that `git-annex` uses. As a result, the extension of “file.tar.gz” would be “.tar.gz”, not “.gz”. In addition, the fields “\_url\_basename\_root\_py” and “\_url\_basename\_ext\_py” provide access to the result of `os.path.splitext`.

- `_url_filename*`

These are similar to `_url_basename*` fields, but they are obtained with a server request. This is useful if the file name is set in the Content-Disposition header.

### Examples

Consider a file “avatars.csv” that contains:

```
who,ext,link
neurodebian,png,https://avatars3.githubusercontent.com/u/260793
datalad,png,https://avatars1.githubusercontent.com/u/8927200
```

To download each link into a file name composed of the ‘who’ and ‘ext’ fields, we could run:

```
$ datalad addurls -d avatar_ds --fast avatars.csv '{link}' '{who}.{ext}'
```

The `-d avatar_ds` is used to create a new dataset in “\$PWD/avatar\_ds”.

If we were already in a dataset and wanted to create a new subdataset in an “avatars” subdirectory, we could use “/” in the *FILENAME-FORMAT* argument:

```
$ datalad addurls --fast avatars.csv '{link}' 'avatars/{who}.{ext}'
```

---

**Note:** For users familiar with ‘git annex addurl’: A large part of this plugin’s functionality can be viewed as transforming data from *URL-FILE* into a “url filename” format that fed to ‘git annex addurl -batch -with-files’.

---

**class EnsureChoice** (\*values)

Bases: `datalad.support.constraints.Constraint`

Ensure an input is element of a set of possible values

`long_description()`

`short_description()`

**class EnsureDataset**

Bases: `datalad.support.constraints.Constraint`

Despite its name, this constraint does not actually ensure that the argument is a valid dataset, because for procedural reasons this would typically duplicate subsequent checks and processing. However, it can be used to achieve uniform documentation of *dataset* arguments.

`long_description()`

`short_description()`

**class EnsureNone**

Bases: `datalad.support.constraints.Constraint`

Ensure an input is of value *None*

```

    long_description()
    short_description()
class EnsureStr (min_len=0)
    Bases: datalad.support.constraints.Constraint
    Ensure an input is a string.
    No automatic conversion is attempted.
    long_description()
    short_description()
class Parameter (constraints=None, doc=None, args=None, **kwargs)
    Bases: object
    This class shall serve as a representation of a parameter.
    get_autodoc (name, indent=' ', width=70, default=None, has_default=False)
        Docstring for the parameter to be used in lists of parameters
        Returns
        Return type string or list of strings (if indent is None)
datasetmethod (name=None, dataset_argname='dataset')
eval_results ()
    Decorator for return value evaluation of datalad commands.
    Note, this decorator is only compatible with commands that return status dict sequences!
    Two basic modes of operation are supported: 1) “generator mode” that yields individual results, and 2)
    “list mode” that returns a sequence of results. The behavior can be selected via the kwarg return_type.
    Default is “list mode”.
    This decorator implements common functionality for result rendering/output, error detection/handling, and
    logging.
    Result rendering/output can be triggered via the datalad.api.result-renderer configuration variable, or the
    result_renderer keyword argument of each decorated command. Supported modes are: ‘default’ (one line
    per result with action, status, path, and an optional message); ‘json’ (one object per result, like git-annex),
    ‘json_pp’ (like ‘json’, but pretty-printed spanning multiple lines), ‘tailored’ custom output formatting
    provided by each command class (if any).
    Error detection works by inspecting the status item of all result dictionaries. Any occurrence of a status
    other than ‘ok’ or ‘notneeded’ will cause an IncompleteResultsError exception to be raised that carries the
    failed actions’ status dictionaries in its failed attribute.
    Status messages will be logged automatically, by default the following association of result status and log
    channel will be used: ‘ok’ (debug), ‘notneeded’ (debug), ‘impossible’ (warning), ‘error’ (error). Logger
    instances included in the results are used to capture the origin of a status report.
        Parameters func (function) – __call__ method of a subclass of Interface, i.e. a datalad
        command definition
class datalad.plugin.addurls.Formatter (idx_to_name=None, missing_value=None)
    Bases: string.Formatter
    Formatter that gives precedence to custom keys.
    The first positional argument to the format call should be a mapping whose keys are exposed as placeholders
    (e.g., “{key1}.py”).
    Parameters

```

- **idx\_to\_name** (*dict*) – A mapping from a positional index to a key. If not provided, “{N}” elements are not supported.
- **missing** (*str, optional*) – When column lookup results in an empty string, use this value in its place.

**convert\_field** (*value, conversion*)

**format** (*format\_string, \*args, \*\*kwargs*)

**get\_value** (*key, args, kwargs*)

Look for key’s value in *args[0]* mapping first.

**class** `datalad.plugin.addurls.RepFormatter` (*\*args, \*\*kwargs*)

Bases: `datalad.plugin.addurls.Formatter`

Extend Formatter to support a `{_reindex}` placeholder.

**format** (*\*args, \*\*kwargs*)

**get\_value** (*key, args, kwargs*)

Look for key’s value in *args[0]* mapping first.

`datalad.plugin.addurls.add_extra_filename_values` (*filename\_format, rows, urls, dry\_run*)

Extend *rows* with values for special formatting fields.

`datalad.plugin.addurls.clean_meta_args` (*args*)

Process metadata arguments.

**Parameters** *args* (*iterable of str*) – Formatted metadata arguments for ‘git-annex metadata –set’.

**Returns**

**Return type** A dict mapping field names to values.

`datalad.plugin.addurls.extract` (*stream, input\_type, url\_format='{0}', filename\_format='{1}', exclude\_autometa=None, meta=None, dry\_run=False, missing\_value=None*)

Extract and format information from *url\_file*.

**Parameters**

- **stream** (*file object*) – Items used to construct the file names and URLs.
- **input\_type** (*{'csv', 'json'}*) –
- **other parameters match those described in AddUrls.** (*All*) –

**Returns**

- *A tuple where the first item is a list with a dict of extracted information*
- *for each row in stream and the second item a list subdataset paths,*
- *sorted breadth-first.*

`datalad.plugin.addurls.filter_legal_metafield` (*fields*)

Remove illegal names from *fields*.

Note: This is like `filter(is_legal_metafield, fields)` but the dropped values are logged.

`datalad.plugin.addurls.fmt_to_name` (*format\_string, num\_to\_name*)

Try to map a format string to a single name.

**Parameters**

- **format\_string** (*string*) –
- **num\_to\_name** (*dict*) – A dictionary that maps from an integer to a column name. This enables mapping the format string to an integer to a name.

**Returns**

- A placeholder name if *format\_string* consists of a single
- *placeholder and no other text. Otherwise, None is returned.*

`datalad.plugin.addurls.get_file_parts` (*filename*, *prefix*='name')

Assign a name to various parts of a file.

**Parameters**

- **filename** (*str*) – A file name (no leading path is permitted).
- **prefix** (*str*) – Prefix to prepend to the key names.

**Returns**

**Return type** A dict mapping each part to a value.

`datalad.plugin.addurls.get_fmt_names` (*format\_string*)

Yield field names in *format\_string*.

`datalad.plugin.addurls.get_subpaths` (*filename*)

Convert “//” marker in *filename* to a list of subpaths.

```
>>> from datalad.plugin.addurls import get_subpaths
>>> get_subpaths("p1/p2//p3/p4//file")
('p1/p2/p3/p4/file', ['p1/p2', 'p1/p2/p3/p4'])
```

Note: With Python 3, the subpaths could be generated with  
`itertools.accumulate(filename.split("//")[:-1], os.path.join)`

**Parameters filename** (*str*) – File name with “//” marking subpaths.

**Returns**

- *A tuple of the filename with any “//” collapsed to a single*
- *separator and a list of subpaths (str).*

`datalad.plugin.addurls.get_url_parts` (*url*)

Assign a name to various parts of the URL.

**Parameters url** (*str*) –

**Returns**

- A dict with keys `_url_hostname` and, for a path with N+1 parts,
- `'_url0'` through `'_urlN'` . There is also a `_url_basename` key for
- *the rightmost part of the path.*

`datalad.plugin.addurls.is_legal_metafield` (*name*)

Test whether *name* is a valid metadata field.

The set of permitted characters is taken from `git-annex's MetaData.hs:legalField`.

`datalad.plugin.addurls.sort_paths` (*paths*)

Sort *paths* by directory level and then alphabetically.

**Parameters** `paths` (*iterable of str*)-

**Returns**

**Return type** Generator of sorted paths.

## datalad.plugin.check\_dates

Extension for checking dates within repositories.

**class** `datalad.plugin.check_dates.CheckDates`

Bases: `datalad.interface.base.Interface`

Find repository dates that are more recent than a reference date.

The main purpose of this tool is to find “leaked” real dates in repositories that are configured to use fake dates. It checks dates from three sources: (1) commit timestamps (author and committer dates), (2) timestamps within files of the “git-annex” branch, and (3) the timestamps of annotated tags.

**class** `EnsureChoice (*values)`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is element of a set of possible values

**long\_description** ()

**short\_description** ()

**class** `EnsureNone`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is of value *None*

**long\_description** ()

**short\_description** ()

**class** `EnsureStr (min_len=0)`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is a string.

No automatic conversion is attempted.

**long\_description** ()

**short\_description** ()

**class** `Parameter (constraints=None, doc=None, args=None, **kwargs)`

Bases: `object`

This class shall serve as a representation of a parameter.

**get\_autodoc** (*name, indent=' ', width=70, default=None, has\_default=False*)

Docstring for the parameter to be used in lists of parameters

**Returns**

**Return type** string or list of strings (if indent is None)

**ac** = `<module 'datalad.support.ansi_colors' from '/home/docs/checkouts/readthedocs.org/`

**static custom\_result\_renderer** (*res, \*\*kwargs*)

Like ‘json\_pp’, but skip non-error results without flagged objects.



**eval\_results ()**

Decorator for return value evaluation of datalad commands.

Note, this decorator is only compatible with commands that return status dict sequences!

Two basic modes of operation are supported: 1) “generator mode” that *yields* individual results, and 2) “list mode” that returns a sequence of results. The behavior can be selected via the kwarg *return\_type*. Default is “list mode”.

This decorator implements common functionality for result rendering/output, error detection/handling, and logging.

Result rendering/output can be triggered via the *datalad.api.result-renderer* configuration variable, or the *result\_renderer* keyword argument of each decorated command. Supported modes are: ‘default’ (one line per result with action, status, path, and an optional message); ‘json’ (one object per result, like git-annex), ‘json\_pp’ (like ‘json’, but pretty-printed spanning multiple lines), ‘tailored’ custom output formatting provided by each command class (if any).

Error detection works by inspecting the *status* item of all result dictionaries. Any occurrence of a status other than ‘ok’ or ‘notneeded’ will cause an *IncompleteResultsError* exception to be raised that carries the failed actions’ status dictionaries in its *failed* attribute.

Status messages will be logged automatically, by default the following association of result status and log channel will be used: ‘ok’ (debug), ‘notneeded’ (debug), ‘impossible’ (warning), ‘error’ (error). Logger instances included in the results are used to capture the origin of a status report.

**Parameters** *func* (*function*) – *\_\_call\_\_* method of a subclass of *Interface*, i.e. a datalad command definition

```
result_renderer = 'tailored'
```

**datalad.plugin.export\_archive**

export a dataset as a compressed TAR/ZIP archive

```
class datalad.plugin.export_archive.ExportArchive
```

Bases: *datalad.interface.base.Interface*

Export the content of a dataset as a TAR/ZIP archive.

```
class EnsureChoice (*values)
```

Bases: *datalad.support.constraints.Constraint*

Ensure an input is element of a set of possible values

```
long_description ()
```

```
short_description ()
```

```
class EnsureDataset
```

Bases: *datalad.support.constraints.Constraint*

Despite its name, this constraint does not actually ensure that the argument is a valid dataset, because for procedural reasons this would typically duplicate subsequent checks and processing. However, it can be used to achieve uniform documentation of *dataset* arguments.

```
long_description ()
```

```
short_description ()
```

```
class EnsureNone
```

Bases: *datalad.support.constraints.Constraint*

Ensure an input is of value *None*

`long_description()`

`short_description()`

**class EnsureStr** (*min\_len=0*)

Bases: `datalad.support.constraints.Constraint`

Ensure an input is a string.

No automatic conversion is attempted.

`long_description()`

`short_description()`

**class Parameter** (*constraints=None, doc=None, args=None, \*\*kwargs*)

Bases: `object`

This class shall serve as a representation of a parameter.

`get_autodoc` (*name, indent=' ', width=70, default=None, has\_default=False*)

Docstring for the parameter to be used in lists of parameters

**Returns**

**Return type** string or list of strings (if indent is None)

**datasetmethod** (*name=None, dataset\_argname='dataset'*)

**eval\_results** ()

Decorator for return value evaluation of datalad commands.

Note, this decorator is only compatible with commands that return status dict sequences!

Two basic modes of operation are supported: 1) “generator mode” that *yields* individual results, and 2) “list mode” that returns a sequence of results. The behavior can be selected via the kwarg *return\_type*. Default is “list mode”.

This decorator implements common functionality for result rendering/output, error detection/handling, and logging.

Result rendering/output can be triggered via the *datalad.api.result-renderer* configuration variable, or the *result\_renderer* keyword argument of each decorated command. Supported modes are: ‘default’ (one line per result with action, status, path, and an optional message); ‘json’ (one object per result, like git-annex), ‘json\_pp’ (like ‘json’, but pretty-printed spanning multiple lines), ‘tailored’ custom output formatting provided by each command class (if any).

Error detection works by inspecting the *status* item of all result dictionaries. Any occurrence of a status other than ‘ok’ or ‘notneeded’ will cause an `IncompleteResultsError` exception to be raised that carries the failed actions’ status dictionaries in its *failed* attribute.

Status messages will be logged automatically, by default the following association of result status and log channel will be used: ‘ok’ (debug), ‘notneeded’ (debug), ‘impossible’ (warning), ‘error’ (error). Logger instances included in the results are used to capture the origin of a status report.

**Parameters func** (*function*) – `__call__` method of a subclass of `Interface`, i.e. a datalad command definition

## datalad.plugin.export\_to\_figshare

export a dataset as a TAR/ZIP archive to figshare

**class** `datalad.plugin.export_to_figshare.ExportToFigshare`

Bases: `datalad.interface.base.Interface`

Export the content of a dataset as a ZIP archive to figshare

Very quick and dirty approach. Ideally figshare should be supported as a proper git annex special remote. Unfortunately, figshare does not support having directories, and can store only a flat list of files. That makes it impossible for any sensible publishing of complete datasets.

The only workaround is to publish dataset as a zip-ball, where the entire content is wrapped into a .zip archive for which figshare would provide a navigator.

**class** `EnsureChoice (*values)`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is element of a set of possible values

`long_description()`

`short_description()`

**class** `EnsureDataset`

Bases: `datalad.support.constraints.Constraint`

Despite its name, this constraint does not actually ensure that the argument is a valid dataset, because for procedural reasons this would typically duplicate subsequent checks and processing. However, it can be used to achieve uniform documentation of *dataset* arguments.

`long_description()`

`short_description()`

**class** `EnsureInt`

Bases: `datalad.support.constraints.EnsureDType`

Ensure that an input (or several inputs) are of a data type 'int'.

**class** `EnsureNone`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is of value *None*

`long_description()`

`short_description()`

**class** `EnsureStr (min_len=0)`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is a string.

No automatic conversion is attempted.

`long_description()`

`short_description()`

**class** `Parameter (constraints=None, doc=None, args=None, **kwargs)`

Bases: `object`

This class shall serve as a representation of a parameter.

`get_autodoc (name, indent=' ', width=70, default=None, has_default=False)`

Docstring for the parameter to be used in lists of parameters

**Returns**

**Return type** string or list of strings (if indent is None)

**datasetmethod** (*name=None, dataset\_argname='dataset'*)

**eval\_results** ()

Decorator for return value evaluation of datalad commands.

Note, this decorator is only compatible with commands that return status dict sequences!

Two basic modes of operation are supported: 1) “generator mode” that *yields* individual results, and 2) “list mode” that returns a sequence of results. The behavior can be selected via the kwarg *return\_type*. Default is “list mode”.

This decorator implements common functionality for result rendering/output, error detection/handling, and logging.

Result rendering/output can be triggered via the *datalad.api.result-renderer* configuration variable, or the *result\_renderer* keyword argument of each decorated command. Supported modes are: ‘default’ (one line per result with action, status, path, and an optional message); ‘json’ (one object per result, like git-annex), ‘json\_pp’ (like ‘json’, but pretty-printed spanning multiple lines), ‘tailored’ custom output formatting provided by each command class (if any).

Error detection works by inspecting the *status* item of all result dictionaries. Any occurrence of a status other than ‘ok’ or ‘notneeded’ will cause an *IncompleteResultsError* exception to be raised that carries the failed actions’ status dictionaries in its *failed* attribute.

Status messages will be logged automatically, by default the following association of result status and log channel will be used: ‘ok’ (debug), ‘notneeded’ (debug), ‘impossible’ (warning), ‘error’ (error). Logger instances included in the results are used to capture the origin of a status report.

**Parameters** **func** (*function*) – *\_\_call\_\_* method of a subclass of *Interface*, i.e. a datalad command definition

**class** `datalad.plugin.export_to_figshare.FigshareRESTLaison`

Bases: `object`

A little helper to provide minimal interface to interact with Figshare

**API\_URL** = `'https://api.figshare.com/v2'`

**create\_article** (*title*)

**get** (*\*args, \*\*kwargs*)

**get\_article\_ids** ()

**post** (*\*args, \*\*kwargs*)

**put** (*\*args, \*\*kwargs*)

**token**

**upload\_file** (*fname, files\_url*)

## **datalad.plugin.no\_annex**

configure which dataset parts to never put in the annex

**class** `datalad.plugin.no_annex.NoAnnex`

Bases: `datalad.interface.base.Interface`

Configure a dataset to never put some content into the dataset’s annex

This can be useful in mixed datasets that also contain textual data, such as source code, which can be efficiently and more conveniently managed directly in Git.

Patterns generally look like this:

```
code/*
```

which would match all file in the code directory. In order to match all files under `code/`, including all its subdirectories use such a pattern:

```
code/**
```

Note that the plugin works incrementally, hence any existing configuration (e.g. from a previous plugin run) is amended, not replaced.

### Parameters

- **ref\_dir** (*str*, *optional*) –
- **makedirs** (*bool*, *optional*) –

### class EnsureDataset

Bases: `datalad.support.constraints.Constraint`

Despite its name, this constraint does not actually ensure that the argument is a valid dataset, because for procedural reasons this would typically duplicate subsequent checks and processing. However, it can be used to achieve uniform documentation of *dataset* arguments.

**long\_description** ()

**short\_description** ()

### class EnsureNone

Bases: `datalad.support.constraints.Constraint`

Ensure an input is of value *None*

**long\_description** ()

**short\_description** ()

### class Parameter (*constraints=None*, *doc=None*, *args=None*, *\*\*kwargs*)

Bases: `object`

This class shall serve as a representation of a parameter.

**get\_autodoc** (*name*, *indent=' '*, *width=70*, *default=None*, *has\_default=False*)

Docstring for the parameter to be used in lists of parameters

#### Returns

**Return type** string or list of strings (if indent is None)

### datasetmethod (*name=None*, *dataset\_argname='dataset'*)

### eval\_results ()

Decorator for return value evaluation of datalad commands.

Note, this decorator is only compatible with commands that return status dict sequences!

Two basic modes of operation are supported: 1) “generator mode” that *yields* individual results, and 2) “list mode” that returns a sequence of results. The behavior can be selected via the kwarg *return\_type*. Default is “list mode”.

This decorator implements common functionality for result rendering/output, error detection/handling, and logging.

Result rendering/output can be triggered via the `datalad.api.result-renderer` configuration variable, or the `result_renderer` keyword argument of each decorated command. Supported modes are: ‘default’ (one line

per result with action, status, path, and an optional message); 'json' (one object per result, like git-annex), 'json\_pp' (like 'json', but pretty-printed spanning multiple lines), 'tailored' custom output formatting provided by each command class (if any).

Error detection works by inspecting the *status* item of all result dictionaries. Any occurrence of a status other than 'ok' or 'notneeded' will cause an `IncompleteResultsError` exception to be raised that carries the failed actions' status dictionaries in its *failed* attribute.

Status messages will be logged automatically, by default the following association of result status and log channel will be used: 'ok' (debug), 'notneeded' (debug), 'impossible' (warning), 'error' (error). Logger instances included in the results are used to capture the origin of a status report.

**Parameters** `func` (*function*) – `__call__` method of a subclass of `Interface`, i.e. a `datalad` command definition

## **datalad.plugin.wtf**

provide information about this DataLad installation

**class** `datalad.plugin.wtf.WTF`

Bases: `datalad.interface.base.Interface`

Generate a report about the DataLad installation and configuration

IMPORTANT: Sharing this report with untrusted parties (e.g. on the web) should be done with care, as it may include identifying information, and/or credentials or access tokens.

**class** `EnsureChoice` (*\*values*)

Bases: `datalad.support.constraints.Constraint`

Ensure an input is element of a set of possible values

`long_description`()

`short_description`()

**class** `EnsureDataset`

Bases: `datalad.support.constraints.Constraint`

Despite its name, this constraint does not actually ensure that the argument is a valid dataset, because for procedural reasons this would typically duplicate subsequent checks and processing. However, it can be used to achieve uniform documentation of *dataset* arguments.

`long_description`()

`short_description`()

**class** `EnsureNone`

Bases: `datalad.support.constraints.Constraint`

Ensure an input is of value *None*

`long_description`()

`short_description`()

**class** `Parameter` (*constraints=None, doc=None, args=None, \*\*kwargs*)

Bases: `object`

This class shall serve as a representation of a parameter.

`get_autodoc` (*name, indent=' ', width=70, default=None, has\_default=False*)

Docstring for the parameter to be used in lists of parameters

**Returns****Return type** string or list of strings (if indent is None)**static custom\_result\_renderer** (*res*, *\*\*kwargs*)**datasetmethod** (*name=None*, *dataset\_argname='dataset'*)**eval\_results** ()

Decorator for return value evaluation of datalad commands.

Note, this decorator is only compatible with commands that return status dict sequences!

Two basic modes of operation are supported: 1) “generator mode” that *yields* individual results, and 2) “list mode” that returns a sequence of results. The behavior can be selected via the kwarg *return\_type*. Default is “list mode”.

This decorator implements common functionality for result rendering/output, error detection/handling, and logging.

Result rendering/output can be triggered via the *datalad.api.result-renderer* configuration variable, or the *result\_renderer* keyword argument of each decorated command. Supported modes are: ‘default’ (one line per result with action, status, path, and an optional message); ‘json’ (one object per result, like git-annex), ‘json\_pp’ (like ‘json’, but pretty-printed spanning multiple lines), ‘tailored’ custom output formatting provided by each command class (if any).

Error detection works by inspecting the *status* item of all result dictionaries. Any occurrence of a status other than ‘ok’ or ‘notneeded’ will cause an *IncompleteResultsError* exception to be raised that carries the failed actions’ status dictionaries in its *failed* attribute.

Status messages will be logged automatically, by default the following association of result status and log channel will be used: ‘ok’ (debug), ‘notneeded’ (debug), ‘impossible’ (warning), ‘error’ (error). Logger instances included in the results are used to capture the origin of a status report.

**Parameters func** (*function*) – *\_\_call\_\_* method of a subclass of Interface, i.e. a datalad command definition

**result\_renderer** = ‘tailored’`datalad.plugin.wtf.get_max_path_length` (*top\_path=None*, *maxl=1000*)

Deduce the maximal length of the filename in a given path

In previous versions of DataLad, plugins were invoked differently than regular DataLad commands, but they can now be called like any other command. The `wtf` plugin, for example, is exposed as

```
% datalad wtf
```

**Plugin detection**

DataLad will discover plugins at three locations:

1. official plugins that are part of the local DataLad installation
2. system-wide plugins, provided by the local admin

The location where plugins need to be placed depends on the platform. On GNU/Linux systems this will be `/etc/xdg/datalad/plugins`, whereas on Windows it will be `C:\ProgramData\datalad.org\datalad\plugins`.

This default location can be overridden by setting the `datalad.locations.system-plugins` configuration variable in the local or global Git configuration.

### 3. user-supplied plugins, customizable by each user

Again, the location will depend on the platform. On GNU/Linux systems this will be `$HOME/.config/datalad/plugins`, whereas on Windows it will be `C:\Users\\AppData\Local\datalad.org\datalad\plugins`.

This default location can be overridden by setting the `datalad.locations.user-plugins` configuration variable in the local or global Git configuration.

Identically named plugins in latter location replace those in locations searched before. This can be used to alter the behavior of plugins provided with DataLad, and enables users to adjust a site-wide configuration.

## Writing own plugins

The best way to go about writing your own plugin, is to have a look at the [source code of those include in DataLad](#). Writing a plugin is rather simple when following the following rules.

## Language and location

Plugins are written in Python. In order for DataLad to be able to find them, plugins need to be placed in one of the supported locations described above. Plugin file names have to have a `.py` extensions and must not start with an underscore (`_`).

## Skeleton of a plugin

The basic structure of a plugin looks like this:

```
from datalad.interface.base import build_doc, Interface

@build_doc
class MyPlugin(Interface):
    """Help message description (parameters will be added automatically)"""
    from datalad.distribution.dataset import datasetmethod, EnsureDataset
    from datalad.interface.utils import eval_results
    from datalad.support.constraints import EnsureNone
    from datalad.support.param import Parameter

    _params_ = dict(
        dataset=Parameter(
            args=("-d", "--dataset"),
            doc="""specify the dataset to report on.
            no dataset is given, an attempt is made to identify the dataset
            based on the current working directory.""",
            constraints=EnsureDataset() | EnsureNone())

    @staticmethod
    @datasetmethod(name='my-plugin')
    @eval_results
    def __call__(dataset):
        # Do things and yield status dicts.
        pass

__datalad_plugin__ = MyPlugin
```



In this example, the plugin is called `my-plugin`. Any number of parameters can be added by extending both the `_params_` dictionary and the signature of `__call__`. The help message for the plugin command is generated using the docstring of the plugin class and the `_params_` dictionary.

### Expected behavior

The plugin's `__call__` method must yield its results as a Python generator. Results are DataLad status dictionaries. There are no constraints on the number of results, or the number and nature of result properties. However, conventions exist and must be followed for compatibility with the result evaluation and rendering performed by DataLad.

The following property keys must exist:

**“status”** { ‘ok’, ‘notneeded’, ‘impossible’, ‘error’ }

**“action”** label for the action performed by the plugin. In many cases this could be the plugin's name.

The following keys should exist if possible:

**“path”** absolute path to a result on the file system

**“type”** label indicating the nature of a result (e.g. ‘file’, ‘dataset’, ‘directory’, etc.)

**“message”** string message annotating the result, particularly important for non-ok results. This can be a tuple with ‘logging’-style string expansion.

### Extension packages

As the name suggests, a *DataLad extension* package is a proper Python package. Consequently, there is a significant amount of boilerplate code involved in the creation of a new Datalad extension. However, this overhead enables a number of useful features for extension developers:

- extensions can provide any number of additional commands that can be grouped into labeled command suites, and are automatically exposed via the standard DataLad commandline and Python API
- extensions can define *entry\_points* for any number of additional metadata extractors that become automatically available to DataLad
- extensions can define *entry\_points* for their test suites, such that the standard *datalad test* command will automatically run these tests in addition to the tests shipped with Datalad core
- extensions can ship additional dataset procedures by installing them into a directory `resources/procedures` underneath the extension module directory

### Using an extension

A *DataLad extension* is a standard Python package. Beyond installation of the package there is no additional setup required.

### Writing your own extensions

A good starting point for implementing a new extension is the “helloworld” demo extension available at <https://github.com/datalad/datalad-extension-template>. This repository can be cloned and adjusted to suit one's needs. It includes:

- a basic Python package setup
- simple demo command implementation

- Travis test setup

A more complex extension setup can be seen in the DataLad Neuroimaging extension: <https://github.com/datalad/datalad-neuroimaging>, including additional metadata extractors, test suite registration, and a sphinx-based documentation setup for a DataLad extension.

As a DataLad extension is a standard Python package, an extension should declare dependencies on an appropriate DataLad version, and possibly other extensions via the standard mechanisms.

## 1.4.6 Design patterns

DataLad is the result of a distributed and collaborative development effort over many years. During this time the scope of the project has changed multiple times. As a consequence, the API and employed technologies have been adjusted repeatedly. Depending on the age of a piece of code, a clear software design is not always immediately visible. This section documents a few design patterns that the project strives to adopt at present. Changes to existing code and new contributions should follow these guidelines.

### Generator methods in *Repo* classes

Substantial parts of DataLad are implemented to behave like Python generators in order to be maximally responsive when processing long-running tasks. This included methods of the core API classes *GitRepo* and *AnnexRepo*. By convention, such methods carry a trailing `_` in their name. In some cases, sibling methods with the same name, but without the trailing underscore are provided. These behave like their generator-equivalent, but eventually return an iterable once processing is fully completed.

### Calls to Git commands

DataLad is built on Git, so calls to Git commands are a key element of the code base. All such calls should be made through methods of the *GitRepo* class. This is necessary, as only there it is made sure that Git operates under the desired conditions (environment configuration, etc.).

For some functionality, for example querying and manipulating *gitattributes*, dedicated methods are provided. However, in many cases simple one-off calls to get specific information from Git, or trigger certain operations are needed. For these purposes the *GitRepo* class provides a set of convenience methods aiming to cover use cases requiring particular return values:

- test success of a command: `call_git_success()`
- obtain *stdout* of a command: `call_git()`
- obtain a single output line: `call_git_oneline()`
- obtain items from output split by a separator: `call_git_items_()`

All these methods take care of raising appropriate exceptions when expected conditions are not met. Whenever desired functionality can be achieved using simple custom calls to Git via these methods, their use is preferred over the implementation of additional, dedicated wrapper methods.

### Command examples

Examples of Python and commandline invocations of DataLad's user-oriented commands are defined in the class of the respective command as dictionaries within `_examples_`:

```

_examples_ = [
  dict(text="""Create a dataset 'mydataset' in the current directory""",
        code_py="create(path='mydataset')",
        code_cmd="datalad create mydataset",
  dict(text="""Apply the text2git procedure upon creation of a dataset""",
        code_py="create(path='mydataset', cfg_proc='text2git')",
        code_cmd="datalad create -c text2git mydataset")
]

```

The formatting of code lines is preserved. Changes to existing examples and new contributions should provide examples for Python and commandline API, as well as a concise description.

## 1.4.7 Glossary

DataLad purposefully uses a terminology that is different from the one used by its technological foundations [Git](#) and [git-annex](#). This glossary provides definitions for terms used in the [datalad](#) documentation and API, and relates them to the corresponding [Git](#)/[git-annex](#) concepts.

**annex** Extension to a [Git](#) repository, provided and managed by [git-annex](#) as means to track and distribute large (and small) files without having to inject them directly into a [Git](#) repository (which would slow [Git](#) operations significantly and impair handling of such repositories in general).

**DataLad extension** A Python package, developed outside of the core DataLad codebase, which (when installed) typically either provides additional top level *datalad* commands and/or additional metadata extractors. Visit [Handbook, Ch.2. DataLad's extensions](#) for a representative list of extensions and instructions on how to install them.

**dataset** A regular [Git](#) repository with an (optional) *annex*.

**sibling** A *dataset* (location) that is related to a particular dataset, by sharing content and history. In [Git](#) terminology, this is a *clone* of a dataset that is configured as a *remote*.

**subdataset** A *dataset* that is part of another dataset, by means of being tracked as a [Git](#) submodule. As such, a subdataset is also a complete dataset and not different from a standalone dataset.

**superdataset** A *dataset* that contains at least one *subdataset*.

## 1.5 Commands and API

### 1.5.1 Command line reference

#### Main command

#### **datalad**

#### Synopsis

```

datalad [-l LEVEL] [--pbs-runner {condor}] [-C PATH] [--version] [--dbg] [--idbg] [-c
↪ KEY=VALUE] [-f {default,json,json_pp,tailored,'<template>'}] [--report-status
↪ {success,failure,ok,notneeded,impossible,error}] [--report-type {dataset,file}] [--
↪ on-failure {ignore,continue,stop}] [--cmd] [-h] {create,install,get,publish,
↪ uninstall,drop,remove,update,create-sibling,create-sibling-github,create-sibling-
↪ gitlab,create-sibling-ria,unlock,save,search,metadata,aggregate-metadata,extract-
↪ metadata,wtf,test,ls,clean,add-archive-content,download-url,run,rerun,
↪ addurls,no-annex,export-to-figshare,add-readme,check-dates,export-archive,annotate-
↪ paths,clone,create-test-dataset,status,diff,siblings,sshrun,subdatasets} ...

```

---

## Description

Comprehensive data management solution

DataLad provides a unified data distribution system built on the Git and Git-annex. DataLad command line tools allow to manipulate (obtain, create, update, publish, etc.) datasets and provide a comprehensive toolbox for joint management of data and code. Compared to Git/annex it primarily extends their functionality to transparently and simultaneously work with multiple inter-related repositories.

## Options

**{create,install,get,publish,uninstall,drop,remove,update,create-sibling,create-sibling-github,create-sibling-gitlab,create-sibling-ria,unlock,save,search,metadata,aggregate-metadata,extract-metadata,wtf,test,ls,clean,add-archive-content,download-url,run,run-procedure,addurls,no-annex,export-to-figshare,add-readme,check-dates,export-archive,annotate-paths,clone,create-test-dataset,status,diff,siblings,sshrun,subdatasets}**

### **-l LEVEL, --log-level LEVEL**

set logging verbosity level. Choose among critical, error, warning, info, debug. Also you can specify an integer <10 to provide even more debugging information

### **--pbs-runner {condor}**

execute command by scheduling it via available PBS. For settings, config file will be consulted

### **-C PATH**

run as if datalad was started in <path> instead of the current working directory. When multiple -C options are given, each subsequent non-absolute -C <path> is interpreted relative to the preceding -C <path>. This option affects the interpretations of the path names in that they are made relative to the working directory caused by the -C option

### **--version**

show the program's version

### **--dbg**

enter Python debugger when uncaught exception happens

### **--idbg**

enter IPython debugger when uncaught exception happens

**-c KEY=VALUE**

configuration variable setting. Overrides any configuration read from a file, but is potentially overridden itself by configuration variables in the process environment.

**-f {default,json,json\_pp,tailored,'<template>'}, --output-format {default,json,json\_pp,tailored,'<template>'}**

select format for returned command results. 'default' give one line per result reporting action, status, path and an optional message; 'json' renders a JSON object with all properties for each result (one per line); 'json\_pp' pretty-prints JSON spanning multiple lines; 'tailored' enables a command-specific rendering style that is typically tailored to human consumption (no result output otherwise), '<template>' reports any value(s) of any result properties in any format indicated by the template (e.g. '{path}'; compare with JSON output for all key-value choices). The template syntax follows the Python "format() language". It is possible to report individual dictionary values, e.g. '{metadata[name]}'. If a 2nd-level key contains a colon, e.g. 'music:Genre', ':' must be substituted by '#' in the template, like so: '{metadata[music#Genre]}'.

**--report-status {success,failure,ok,notneeded,impossible,error}**

constrain command result report to records matching the given status. 'success' is a synonym for 'ok' OR 'notneeded', 'failure' stands for 'impossible' OR 'error'.

**--report-type {dataset,file}**

constrain command result report to records matching the given type. Can be given more than once to match multiple types.

**--on-failure {ignore,continue,stop}**

when an operation fails: 'ignore' and continue with remaining operations, the error is logged but does not lead to a non-zero exit code of the command; 'continue' works like 'ignore', but an error causes a non-zero exit code; 'stop' halts on first failure and yields non-zero exit code. A failure is any result with status 'impossible' or 'error'.

**-cmd**

syntactical helper that can be used to end the list of global command line options before the subcommand label. Options taking an arbitrary number of arguments may require to be followed by a single -cmd in order to enable identification of the subcommand.

**-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

"Be happy!"

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## Dataset operations

### datalad create

#### Synopsis

```
datalad create [-h] [-f] [-D DESCRIPTION] [-d DATASET] [--no-annex] [--fake-dates] [-c PROC] [PATH] ...
```

#### Description

Create a new dataset from scratch.

This command initializes a new dataset at a given location, or the current directory. The new dataset can optionally be registered in an existing superdataset (the new dataset's path needs to be located within the superdataset for that, and the superdataset needs to be given explicitly via `-dataset`). It is recommended to provide a brief description to label the dataset's nature *and* location, e.g. "Michael's music on black laptop". This helps humans to identify data locations in distributed scenarios. By default an identifier comprised of user and machine name, plus path will be generated.

This command only creates a new dataset, it does not add existing content to it, even if the target directory already contains additional files or directories.

Plain Git repositories can be created via the `-no-annex` flag. However, the result will not be a full dataset, and, consequently, not all features are supported (e.g. a description).

To create a local version of a remote dataset use the *install* command instead.

**NOTE** Power-user info: This command uses `git init` and `git annex init` to prepare the new dataset. Registering to a superdataset is performed via a `git submodule add` operation in the discovered superdataset.

#### Examples

Create a dataset 'mydataset' in the current directory:

```
% datalad create mydataset
```

Apply the `text2git` procedure upon creation of a dataset:

```
% datalad create -c text2git mydataset
```

Create a subdataset in the root of an existing dataset:

```
% datalad create -d . mysubdataset
```

Create a dataset in an existing, non-empty directory:

```
% datalad create --force
```

Create a plain Git repository:

```
% datalad create --no-annex mydataset
```

#### Options

## PATH

path where the dataset shall be created, directories will be created as necessary. If no location is provided, a dataset will be created in the current working directory. Either way the command will error if the target directory is not empty. Use FORCE to create a dataset in a non-empty directory. Constraints: value must be a string, or Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

## INIT OPTIONS

options to pass to git init. Any argument specified after the destination path of the repository will be passed to git-init as-is. Note that not all options will lead to viable results. For example ‘-bare’ will not yield a repository where DataLad can adjust files in its working tree.

### **-h, -help, -help-np**

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

### **-f, -force**

enforce creation of a dataset in a non-empty directory.

### **-D DESCRIPTION, -description DESCRIPTION**

short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. Constraints: value must be a string

### **-d DATASET, -dataset DATASET**

specify the dataset to perform the create operation on. If a dataset is given, a new subdataset will be created in it. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-no-annex**

if set, a plain Git repository will be created without any annex.

### **-fake-dates**

Configure the repository to use fake dates. The date for a new commit will be set to one second later than the latest commit in the repository. This can be used to anonymize dates.

### **-c PROC, -cfg-proc PROC**

Run `cfg_PROC` procedure(s) (can be specified multiple times) on the created dataset. Use `run_procedure -discover` to get a list of available procedures, such as `cfg_text2git`.

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad create-sibling

### Synopsis

```
datalad create-sibling [-h] [-s [NAME]] [--target-dir PATH] [--target-url URL] [--
↳target-pushurl URL] [--dataset DATASET] [-r] [-R LEVELS] [--existing MODE] [--
↳shared {false|true|umask|group|all|world|everybody|0xxx}] [--group GROUP] [--ui
↳{false|true|html_filename}] [--as-common-datasrc NAME] [--publish-by-default
↳REFSPEC] [--publish-depends SIBLINGNAME] [--annex-wanted EXPR] [--annex-group EXPR]
↳[--annex-groupwanted EXPR] [--inherit] [--since SINCE] [SSHURL]
```

### Description

Create a dataset sibling on a UNIX-like SSH-accessible machine

Given a local dataset, and SSH login information this command creates a remote dataset repository and configures it as a dataset sibling to be used as a publication target (see PUBLISH command).

Various properties of the remote sibling can be configured (e.g. name location on the server, read and write access URLs, and access permissions).

Optionally, a basic web-viewer for DataLad datasets can be installed at the remote location.

This command supports recursive processing of dataset hierarchies, creating a remote sibling for each dataset in the hierarchy. By default, remote siblings are created in hierarchical structure that reflects the organization on the local file system. However, a simple templating mechanism is provided to produce a flat list of datasets (see `-target-dir`).

### Options

#### SSHURL

Login information for the target server. This can be given as a URL (`ssh://host/path`) or SSH-style (`user@host:path`). Unless overridden, this also serves the future dataset's access URL and path on the server. Constraints: value must be a string

#### **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

#### **-s [NAME], -name [NAME]**

sibling name to create for this publication target. If RECURSIVE is set, the same name will be used to label all the subdatasets' siblings. When creating a target dataset fails, no sibling is added. Constraints: value must be a string



**-target-dir PATH**

path to the directory *on the server* where the dataset shall be created. By default the SSH access URL is used to identify this directory. If a relative path is provided here, it is interpreted as being relative to the user's home directory on the server. Additional features are relevant for recursive processing of datasets with subdatasets. By default, the local dataset structure is replicated on the server. However, it is possible to provide a template for generating different target directory names for all (sub)datasets. Templates can contain certain placeholder that are substituted for each (sub)dataset. For example: `"/mydirectory/dataset%RELNAME"`. Supported placeholders: `%RELNAME` - the name of the datasets, with any slashes replaced by dashes. Constraints: value must be a string

**-target-url URL**

"public" access URL of the to-be-created target dataset(s) (default: SSHURL). Accessibility of this URL determines the access permissions of potential consumers of the dataset. As with `TARGET_DIR`, templates (same set of placeholders) are supported. Also, if specified, it is provided as the annex description. Constraints: value must be a string

**-target-pushurl URL**

In case the `TARGET_URL` cannot be used to publish to the dataset, this option specifies an alternative URL for this purpose. As with `TARGET_URL`, templates (same set of placeholders) are supported. Constraints: value must be a string

**-dataset DATASET, -d DATASET**

specify the dataset to create the publication target for. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-r, --recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**-existing MODE**

action to perform, if a sibling is already configured under the given name and/or a target (non-empty) directory already exists. In this case, a dataset can be skipped ('skip'), the sibling configuration be updated ('reconfigure'), or process interrupts with error ('error'). **DANGER ZONE:** If 'replace' is used, an existing target directory will be forcefully removed, re-initialized, and the sibling (re-)configured (thus implies 'reconfigure'). **REPLACE** could lead to data loss, so use with care. To minimize possibility of data loss, in interactive mode DataLad will ask for confirmation, but it would just issue a warning and proceed in non-interactive mode. Constraints: value must be one of ('skip', 'error', 'reconfigure', 'replace') [Default: 'error']

**–shared {false|true|umask|group|all|world|everybody|0xxx}**

if given, configures the access permissions on the server for multi-users (this could include access by a webserver!). Possible values for this option are identical to those of *git init –shared* and are described in its documentation. Constraints: value must be a string, or value must be convertible to type bool

**–group GROUP**

Filesystem group for the repository. Specifying the group is particularly important when *–shared=group*. Constraints: value must be a string

**–ui {false|true|html\_filename}**

publish a web interface for the dataset with an optional user-specified name for the html at publication target. defaults to INDEX.HTML at dataset root. Constraints: value must be convertible to type bool, or value must be a string [Default: False]

**–as-common-datasrc NAME**

configure the created sibling as a common data source of the dataset that can be automatically used by all consumers of the dataset (technical: git-annex auto- enabled special remote).

**–publish-by-default REFSPEC**

add a refsPEC to be published to this sibling by default if nothing specified. Constraints: value must be a string

**–publish-depends SIBLINGNAME**

add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item ‘remote.SIBLINGNAME.datalad-publish-depends’. This option can be given more than once to configure multiple dependencies. Constraints: value must be a string

**–annex-wanted EXPR**

expression to specify ‘wanted’ content for the repository/sibling. See <https://git-annex.branchable.com/git-annex-wanted/> for more information. Constraints: value must be a string

**–annex-group EXPR**

expression to specify a group for the repository. See <https://git-annex.branchable.com/git-annex-group/> for more information. Constraints: value must be a string

### –annex-groupwanted EXPR

expression for the groupwanted. Makes sense only if –annex-wanted=”groupwanted” and annex-group is given too. See <https://git-annex.branchable.com/git-annex-groupwanted/> for more information. Constraints: value must be a string

### –inherit

if sibling is missing, inherit settings (git config, git annex wanted/group/groupwanted) from its super-dataset.

### –since SINCE

limit processing to datasets that have been changed since a given state (by tag, branch, commit, etc). This can be used to create siblings for recently added subdatasets. Constraints: value must be a string

## Authors

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## datalad create-sibling-github

### Synopsis

```
datalad create-sibling-github [-h] [--dataset DATASET] [-r] [-R LEVELS] [-s NAME] [--
↪existing MODE] [--github-login NAME] [--github-passwd PASSWORD] [--github-
↪organization NAME] [--access-protocol {https|ssh}] [--publish-depends SIBLINGNAME]
↪ [--dryrun] REPONAME
```

### Description

Create dataset sibling on Github.

An existing GitHub project, or a project created via the GitHub website can be configured as a sibling with the siblings command. Alternatively, this command can create a repository under a user’s Github account, or any organization a user is a member of (given appropriate permissions). This is particularly helpful for recursive sibling creation for subdatasets. In such a case, a dataset hierarchy is represented as a flat list of GitHub repositories.

Github cannot host dataset content. However, in combination with other data sources (and siblings), publishing a dataset to Github can facilitate distribution and exchange, while still allowing any dataset consumer to obtain actual data content from alternative sources.

For Github authentication user credentials can be given as arguments. Alternatively, they are obtained interactively or queried from the systems credential store. Lastly, an *oauth* token stored in the Git configuration under variable *hub.oauthtoken* will be used automatically. Such a token can be obtained, for example, using the commandline Github interface (<https://github.com/sociomantic/git-hub>) by running: git hub setup (if no 2FA is used).

## Options

### **REPONAME**

Github repository name. When operating recursively, a suffix will be appended to this name for each subdataset. Constraints: value must be a string

### **-h, -help, -help-np**

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

### **-dataset DATASET, -d DATASET**

specify the dataset to create the publication target for. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-r, -recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, -recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

### **-s NAME, -name NAME**

name to represent the Github repository in the local dataset installation. Constraints: value must be a string [Default: 'github']

### **-existing MODE**

desired behavior when already existing or configured siblings are discovered. 'skip': ignore; 'error': fail immediately; 'reconfigure': use the existing repository and reconfigure the local dataset to use it as a sibling. Constraints: value must be one of ('skip', 'error', 'reconfigure') [Default: 'error']

### **-github-login NAME**

Github user name or access token. Constraints: value must be a string

### **-github-passwd PASSWORD**

Github user password. Constraints: value must be a string

### **–github-organization NAME**

If provided, the repository will be created under this Github organization. The respective Github user needs appropriate permissions. Constraints: value must be a string

### **–access-protocol {https|ssh}**

Which access protocol/URL to configure for the sibling. Constraints: value must be one of ('https', 'ssh') [Default: 'https']

### **–publish-depends SIBLINGNAME**

add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item 'remote.SIBLINGNAME.datalad-publish-depends'. This option can be given more than once to configure multiple dependencies. Constraints: value must be a string

### **–dryrun**

If this flag is set, no communication with Github is performed, and no repositories will be created. Instead would-be repository names are reported for all relevant datasets.

## **Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## **datalad create-sibling-gitlab**

### **Synopsis**

```
datalad create-sibling-gitlab [-h] [--site SITENAME] [--project NAME/LOCATION] [--
↳ layout {hierarchy|collection|flat}] [--dataset DATASET] [-r] [-R LEVELS] [-s NAME]
↳ [--existing {skip|error|reconfigure}] [--access {http|ssh|ssh+http}] [--publish-
↳ depends SIBLINGNAME] [--description DESCRIPTION] [--dryrun] [PATH [PATH ...]]
```

### **Description**

Create dataset sibling at a GitLab site

An existing GitLab project, or a project created via the GitLab web interface can be configured as a sibling with the siblings command. Alternatively, this command can create a GitLab project at any location/path a given user has appropriate permissions for. This is particularly helpful for recursive sibling creation for subdatasets. API access and authentication are implemented via python-gitlab, and all its features are supported. A particular GitLab site must be configured in a named section of a python-gitlab.cfg file (see <https://python-gitlab.readthedocs.io/en/stable/cli.html#configuration> for details), such as:

```
[mygit]
url = https://git.example.com
api_version = 4
private_token = abcdefghijklmnopqrst
```

Subsequently, this site is identified by its name (‘mygit’ in the example above).

(Recursive) sibling creation for all, or a selected subset of subdatasets is supported with three different project layouts (see `-layout`):

**“hierarchy”** Each dataset is placed into its own group, and the actual GitLab project for a dataset is put in a project named “\_repo\_” inside this group. Using this layout, arbitrarily deep hierarchies of nested datasets can be represented, while the hierarchical structure is reflected in the project path. This is the default layout, if no project path is specified.

**“flat”** All datasets are placed in the same group. The name of a project is its relative path within the root dataset, with all path separator characters replaced by ‘-’.

**“collection”** This is a hybrid layout, where the root dataset is placed in a “\_repo\_” project inside a group, and all nested subdatasets are represented inside the group using a “flat” layout.

GitLab cannot host dataset content. However, in combination with other data sources (and siblings), publishing a dataset to GitLab can facilitate distribution and exchange, while still allowing any dataset consumer to obtain actual data content from alternative sources.

### *Configuration*

All configuration switches and options for GitLab sibling creation can be provided arguments to the command. However, it is also possible to specify a particular setup in a dataset’s configuration. This is particularly important when managing large collections of datasets. Configuration options are:

**“datalad.gitlab-default-site”** Name of the default GitLab site (see `-site`)

**“datalad.gitlab-SITENAME-siblingname”** Name of the sibling configured for the local dataset that points to the GitLab instance SITENAME (see `-name`)

**“datalad.gitlab-SITENAME-layout”** Project layout used at the GitLab instance SITENAME (see `-layout`)

**“datalad.gitlab-SITENAME-access”** Access method used for the GitLab instance SITENAME (see `-access`)

**“datalad.gitlab-SITENAME-project”** Project location/path used for a datasets at GitLab instance SITENAME (see `-project`). Configuring this is useful for deriving project paths for subdatasets, relative to superdataset.

## Options

### PATH

selectively create siblings for any datasets underneath a given path. By default only the root dataset is considered.

### **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

**-site SITENAME**

name of the GitLab site to create a sibling at. Must match an existing python- gitlab configuration section with location and authentication settings (see <https://python-gitlab.readthedocs.io/en/stable/cli.html#configuration>). By default the dataset configuration is consulted. Constraints: value must be NONE, or value must be a string

**-project NAME/LOCATION**

project path at the GitLab site. If a subdataset of the reference dataset is processed, its project path is automatically determined by the LAYOUT configuration, by default. Constraints: value must be NONE, or value must be a string

**-layout {hierarchy|collection|flat}**

layout of projects at the GitLab site, if a collection, or a hierarchy of datasets and subdatasets is to be created. By default the dataset configuration is consulted. Constraints: value must be one of (None, 'hierarchy', 'collection', 'flat')

**-dataset DATASET, -d DATASET**

reference or root dataset. If no path constraints are given, a sibling for this dataset will be created. In this and all other cases, the reference dataset is also consulted for the GitLab configuration, and desired project layout. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-r, --recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**-s NAME, --name NAME**

name to represent the GitLab sibling remote in the local dataset installation. If not specified a name is looked up in the dataset configuration, or defaults to the SITE name. Constraints: value must be a string

**-existing {skip|error|reconfigure}**

desired behavior when already existing or configured siblings are discovered. 'skip': ignore; 'error': fail, if access URLs differ; 'reconfigure': use the existing repository and reconfigure the local dataset to use it as a sibling. Constraints: value must be one of ('skip', 'error', 'reconfigure') [Default: 'error']

### **–access {http|ssh|ssh+http}**

access method used for data transfer to and from the sibling. ‘ssh’: read and write access used the SSH protocol; ‘http’: read and write access use HTTP requests; ‘ssh+http’: read access is done via HTTP and write access performed with SSH. Dataset configuration is consulted for a default, ‘http’ is used otherwise. Constraints: value must be one of (None, ‘http’, ‘ssh’, ‘ssh+http’)

### **–publish-depends SIBLINGNAME**

add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item ‘remote.SIBLINGNAME.datalad-publish-depends’. This option can be given more than once to configure multiple dependencies. Constraints: value must be a string

### **–description DESCRIPTION**

brief description for the GitLab project (displayed on the site). Constraints: value must be a string

### **–dryrun**

If this flag is set, no communication with GitLab is performed, and no repositories will be created. Instead would-be repository names and configurations are reported for all relevant datasets.

## **Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## **datalad create-sibling-ria**

### **Synopsis**

```
datalad create-sibling-ria [-h] -s NAME [-d DATASET] [--ria-remote-name NAME] [--post-  
↪update-hook] [--shared {false|true|umask|group|all|world|everybody|0xxx}] [--group_  
↪GROUP] [--no-ria-remote] [--existing MODE] [--trust-level TRUST-LEVEL] [-r] [-R_  
↪LEVELS] ria+<ssh|file>://<host>[/path]
```

### **Description**

Creates a sibling to a dataset in a RIA store

This creates a representation of a dataset in a ria-remote compliant storage location. For access to it two siblings are configured for the dataset by default. A “regular” one and a RIA remote (git-annex special remote). Furthermore, the former is configured to have a publication dependency on the latter. If not given a default name for the RIA remote is derived from the sibling’s name by appending “-ria”.

The store’s base path currently is expected to either:

- not yet exist or
- be empty or



- have a valid “ria-layout-version” file and an “error\_logs” directory.

In the first two cases, said file and directory are created by this command. Alternatively you can manually create the third case, of course. Please note, that “ria-layout-version” needs to contain a line stating the version (currently “1”) and optionally enable error logging (append a pipe symbol and an “l” in that case). Currently, this line **MUST** end with a newline!

Error logging will create files in the “error\_log” directory whenever the RIA special remote (storage sibling) raises an exception, storing the python traceback of it. The logfiles are named according to the scheme <dataset id>.<annex uuid of the remote>.log showing ‘who’ ran into this issue with what dataset. Since this logging can potentially leak personal data (like local file paths for example) it can be disabled from the client side via “annex.ria-remote.<RIAREMOTE>.ignore-remote-config”.

## Todo

Where to put the description of a RIA store (see below)?

The targeted layout of such a store is a tree of directories containing datasets at the lowest level, starting at the configured base path. First level of subdirectories are named for the first three characters of the datasets’ id, second level is the remainder of those ids. The thereby created dataset directories contain a bare git repository. Those bare repositories are slightly different from plain git-annex bare repositories in that they use the standard dirhashmixed layout beneath annex/objects as opposed to dirhashlower, which is git-annex’s default for bare repositories. Furthermore, there is an additional directory ‘archives’ within the dataset directories, which may or may not contain archives with annexed content. Note, that this helps to reduce the number of inodes consumed (no checkout + potential archive) as well as it allows to resolve dependencies (that is (sub)datasets) merely by their id. Finally, there is a file “ria-layout-version” put beneath the store’s base path, determining the version of the dataset tree layout and a file of the same name per each dataset directory determining object tree layout version (we already switch from dirhashlower to dirhashmixed for example) and an additional directory “error\_logs” at the toplevel.

## Options

**ria+<ssh|file>://<host>[/path]**

URL identifying the target RIA store and access protocol. Constraints: value must be a string

**-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

**-s NAME, --name NAME**

Name of the sibling. With RECURSIVE, the same name will be used to label all the subdatasets’ siblings. Constraints: value must be a string

**-d DATASET, --dataset DATASET**

specify the dataset to process. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**–ria-remote-name NAME**

Name of the RIA remote (a git-annex special remote). Must not be identical to the sibling name. If not specified, defaults to the sibling name plus a ‘-ria’ suffix. Constraints: value must be a string

**–post-update-hook**

Enable git’s default post-update-hook for the created sibling.

**–shared {false|true|umask|group|all|world|everybody|0xxx}**

If given, configures the permissions in the RIA store for multi-users access. Possible values for this option are identical to those of *git init –shared* and are described in its documentation. Constraints: value must be a string, or value must be convertible to type bool

**–group GROUP**

Filesystem group for the repository. Specifying the group is crucial when *–shared=group*. Constraints: value must be a string

**–no-ria-remote**

Flag to disable establishing remote indexed archive (RIA) capabilities for the created sibling. If enabled, git-annex special remote access will be configured to enable regular git-annex key storage, and also retrieval of keys from (compressed) 7z archives that might be provided by the dataset store. If disabled, git-annex is instructed to ignore the sibling.

**–existing MODE**

Action to perform, if a sibling or ria-remote is already configured under the given name and/or a target already exists. In this case, a dataset can be skipped (‘skip’), an existing target repository be forcefully re-initialized, and the sibling (re-)configured (‘reconfigure’), or the command be instructed to fail (‘error’). Constraints: value must be one of (‘skip’, ‘error’, ‘reconfigure’) [Default: ‘error’]

**–trust-level TRUST-LEVEL**

specify a trust level for the RIA sibling. Internally this will call the respective git-annex command. If not specified nothing will be explicitly done, thereby defaulting to git-annex’ default. Constraints: value must be one of (‘trust’, ‘semitrust’, ‘untrust’)

**-r, –recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

**datalad drop****Synopsis**

```
datalad drop [-h] [-d DATASET] [-r] [-R LEVELS] [--nocheck] [--if-dirty {fail,save-
↪before,ignore}] [PATH [PATH ...]]
```

**Description**

Drop file content from datasets

This command takes any number of paths of files and/or directories. If a common (super)dataset is given explicitly, the given paths are interpreted relative to this dataset.

Recursion into subdatasets needs to be explicitly enabled, while recursion into subdirectories within a dataset is done automatically. An optional recursion limit is applied relative to each given input path.

By default, the availability of at least one remote copy is verified before file content is dropped. As these checks could lead to slow operation (network latencies, etc), they can be disabled.

*Examples*

Drop single file content:

```
% datalad drop <path/to/file>
```

Drop all file content in the current dataset:

```
% datalad drop
```

Drop all file content in a dataset and all its subdatasets:

```
% datalad drop --dataset <path/to/dataset> --recursive
```

Disable check to ensure the configured minimum number of remote sources for dropped data:

```
% datalad drop <path/to/content> --nocheck
```

**Options****PATH**

path/name of the component to be dropped. Constraints: value must be a string

**-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

**-d DATASET, --dataset DATASET**

specify the dataset to perform the operation on. If no dataset is given, an attempt is made to identify a dataset based on the PATH given. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-r, --recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**--nocheck**

whether to perform checks to assure the configured minimum number (remote) source for data. Give this option to skip checks.

**--if-dirty {fail,save-before,ignore}**

desired behavior if a dataset with unsaved changes is discovered: 'fail' will trigger an error and further processing is aborted; 'save-before' will save all changes prior any further action; 'ignore' let's datalad proceed as if the dataset would not have unsaved changes. [Default: 'save-before']

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad get**

**Synopsis**

```
datalad get [-h] [-s LABEL] [-d PATH] [-r] [-R LEVELS] [-n] [-D DESCRIPTION] [--  
↪reckless [{auto|ephemeral}]] [-J NJOBS] [PATH [PATH ...]]
```

**Description**

Get any dataset content (files/directories/subdatasets).

This command only operates on dataset content. To obtain a new independent dataset from some source use the `INSTALL` command.

By default this command operates recursively within a dataset, but not across potential subdatasets, i.e. if a directory is provided, all files in the directory are obtained. Recursion into subdatasets is supported too. If enabled, relevant subdatasets are detected and installed in order to fulfill a request.

Known data locations for each requested file are evaluated and data are obtained from some available location (according to git-annex configuration and possibly assigned remote priorities), unless a specific source is specified.

**NOTE** Power-user info: This command uses git annex get to fulfill file handles.

### Examples

Get a single file:

```
% datalad get <path/to/file>
```

Get contents of a directory:

```
% datalad get <path/to/dir/>
```

Get all contents of the current dataset and its subdatasets:

```
% datalad get . --recursive
```

Get (clone) a registered subdataset, but don't retrieve data:

```
% datalad get -n <path/to/subds>
```

## Options

### PATH

path/name of the requested dataset component. The component must already be known to a dataset. To add new components to a dataset use the ADD command. Constraints: value must be a string

### -h, --help, --help-np

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

### -s LABEL, --source LABEL

label of the data source to be used to fulfill requests. This can be the name of a dataset sibling or another known source. Constraints: value must be a string

### -d PATH, --dataset PATH

specify the dataset to perform the add operation on, in which case PATH arguments are interpreted as being relative to this dataset. If no dataset is given, an attempt is made to identify a dataset for each input PATH. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### -r, --recursive

if set, recurse into potential subdataset.

### **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Alternatively, ‘existing’ will limit recursion to subdatasets that already existed on the filesystem at the start of processing, and prevent new subdatasets from being obtained recursively. Constraints: value must be convertible to type ‘int’, or value must be one of (‘existing’,)

### **-n, --no-data**

whether to obtain data for all file handles. If disabled, GET operations are limited to dataset handles. This option prevents data for file handles from being obtained.

### **-D DESCRIPTION, --description DESCRIPTION**

short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. Constraints: value must be a string

### **--reckless [{auto|ephemeral}]**

Set up the dataset to be able to obtain content in the cheapest/fastest possible way, even if this poses a potential risk the data integrity (‘auto’: hardlink files from a local clone of the dataset, ‘ephemeral’: symlink annex to origin’s annex and discard local availability info via git-annex-dead ‘here’. Please note, that with a symlinked annex you share the annex with origin w/o git-annex knowing this. In case of a change in origin you need to update the clone before you’re able to save new content on your end.). Use with care, and limit to “read-only” use cases. With this flag the installed dataset will be marked as untrusted. The reckless mode is stored in a dataset’s local configuration under ‘datalad.clone.reckless’, and will be inherited to any of its subdatasets. Constraints: value must be one of (None, True, False, ‘auto’, ‘ephemeral’)

### **-J NJOBS, --jobs NJOBS**

how many parallel jobs (where possible) to use. Constraints: value must be convertible to type ‘int’, or value must be one of (‘auto’,) [Default: ‘auto’]

## **Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## **datalad install**

### **Synopsis**

```
datalad install [-h] [-s SOURCE] [-d DATASET] [-g] [-D DESCRIPTION] [-r] [-R LEVELS]   
↔ [--nosave] [--reckless [{auto|ephemeral}]] [-J NJOBS] [PATH [PATH ...]]
```

## Description

Install a dataset from a (remote) source.

This command creates a local sibling of an existing dataset from a (remote) location identified via a URL or path. Optional recursion into potential subdatasets, and download of all referenced data is supported. The new dataset can be optionally registered in an existing superdataset by identifying it via the DATASET argument (the new dataset's path needs to be located within the superdataset for that).

It is recommended to provide a brief description to label the dataset's nature *and* location, e.g. "Michael's music on black laptop". This helps humans to identify data locations in distributed scenarios. By default an identifier comprised of user and machine name, plus path will be generated.

When only partial dataset content shall be obtained, it is recommended to use this command without the GET-DATA flag, followed by a *get* operation to obtain the desired data.

**NOTE** Power-user info: This command uses git clone, and git annex init to prepare the dataset. Registering to a superdataset is performed via a git submodule add operation in the discovered superdataset.

### Examples

Install a dataset from Github into the current directory:

```
% datalad install https://github.com/datalad-datasets/longnow-podcasts.git
```

Install a dataset as a subdataset into the current dataset:

```
% datalad install -d . \
  --source='https://github.com/datalad-datasets/longnow-podcasts.git'
```

Install a dataset, and get all content right away:

```
% datalad install --get-data \
  --source https://github.com/datalad-datasets/longnow-podcasts.git
```

Install a dataset with all its subdatasets:

```
% datalad install --recursive \
  https://github.com/datalad-datasets/longnow-podcasts.git
```

## Options

### PATH

path/name of the installation target. If no PATH is provided a destination path will be derived from a source URL similar to git clone.

### **-h, --help, --help-np**

show this help message. **--help-np** forcefully disables the use of a pager for displaying the help message

### **-s SOURCE, --source SOURCE**

URL or local path of the installation source. Constraints: value must be a string

**-d DATASET, --dataset DATASET**

specify the dataset to perform the install operation on. If no dataset is given, an attempt is made to identify the dataset in a parent directory of the current working directory and/or the PATH given. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-g, --get-data**

if given, obtain all data content too.

**-D DESCRIPTION, --description DESCRIPTION**

short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. Constraints: value must be a string

**-r, --recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type ‘int’

**--nosave**

by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior.

**--reckless [{auto|ephemeral}]**

Set up the dataset to be able to obtain content in the cheapest/fastest possible way, even if this poses a potential risk the data integrity (‘auto’: hardlink files from a local clone of the dataset, ‘ephemeral’: symlink annex to origin’s annex and discard local availability info via git-annex-dead ‘here’. Please note, that with a symlinked annex you share the annex with origin w/o git-annex knowing this. In case of a change in origin you need to update the clone before you’re able to save new content on your end.). Use with care, and limit to “read-only” use cases. With this flag the installed dataset will be marked as untrusted. The reckless mode is stored in a dataset’s local configuration under ‘datalad.clone.reckless’, and will be inherited to any of its subdatasets. Constraints: value must be one of (None, True, False, ‘auto’, ‘ephemeral’)

**-J NJOBS, --jobs NJOBS**

how many parallel jobs (where possible) to use. Constraints: value must be convertible to type ‘int’, or value must be one of (‘auto’,) [Default: ‘auto’]



## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad publish

### Synopsis

```
datalad publish [-h] [-d DATASET] [--to LABEL] [--since SINCE] [--missing MODE] [-f]
↪ [--transfer-data {auto|none|all}] [-r] [-R LEVELS] [--git-opts STRING] [--annex-
↪ opts STRING] [--annex-copy-opts STRING] [-J NJOBS] [PATH [PATH ...]]
```

### Description

Publish a dataset to a known sibling.

This makes the last saved state of a dataset available to a sibling or special remote data store of a dataset. Any target sibling must already exist and be known to the dataset.

Optionally, it is possible to limit publication to change sets relative to a particular point in the version history of a dataset (e.g. a release tag). By default, the state of the local dataset is evaluated against the last known state of the target sibling. Actual publication is only attempted if there was a change compared to the reference state, in order to speed up processing of large collections of datasets. Evaluation with respect to a particular “historic” state is only supported in conjunction with a specified reference dataset. Change sets are also evaluated recursively, i.e. only those subdatasets are published where a change was recorded that is reflected in to current state of the top-level reference dataset. See “since” option for more information.

Only publication of saved changes is supported. Any unsaved changes in a dataset (hierarchy) have to be saved before publication.

**NOTE** Power-user info: This command uses git push, and git annex copy to publish a dataset. Publication targets are either configured remote Git repositories, or git-annex special remotes (if they support data upload).

### Options

#### PATH

path(s), that may point to file handle(s) to publish including their actual content or to subdataset(s) to be published. If a file handle is published with its data, this implicitly means to also publish the (sub)dataset it belongs to. ‘.’ as a path is treated in a special way in the sense, that it is passed to subdatasets in case RECURSIVE is also given. Constraints: value must be a string

#### -h, -help, -help-np

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

#### -d DATASET, -dataset DATASET

specify the (top-level) dataset to be published. If no dataset is given, the datasets are determined based on the input arguments. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-to LABEL**

name of the target sibling. If no name is given an attempt is made to identify the target based on the dataset's configuration (i.e. a configured tracking branch, or a single sibling that is configured for publication). Constraints: value must be a string

**-since SINCE**

When publishing dataset(s), specifies commit (treeish, tag, etc) from which to look for changes to decide whether updated publishing is necessary for this and which children. If empty argument is provided, then we would take from the previously published to that remote/sibling state (for the current branch). Constraints: value must be a string

**-missing MODE**

action to perform, if a sibling does not exist in a given dataset. By default it would fail the run ('fail' setting). With 'inherit' a 'create-sibling' with '- inherit-settings' will be used to create sibling on the remote. With 'skip' - it simply will be skipped. Constraints: value must be one of ('fail', 'inherit', 'skip') [Default: 'fail']

**-f, -force**

enforce doing publish activities (git push etc) regardless of the analysis if they seemed needed.

**-transfer-data {auto|none|all}**

ADDME. Constraints: value must be one of ('auto', 'none', 'all') [Default: 'auto']

**-r, -recursive**

if set, recurse into potential subdataset.

**-R LEVELS, -recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**-git-opts STRING**

option string to be passed to git calls. Constraints: value must be a string

**-annex-opts STRING**

option string to be passed to git annex calls. Constraints: value must be a string

**-annex-copy-opts STRING**

option string to be passed to git annex copy calls. Constraints: value must be a string

## -J NJOBS, -jobs NJOBS

how many parallel jobs (where possible) to use. Constraints: value must be convertible to type ‘int’, or value must be one of (‘auto’,)

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad remove

### Synopsis

```
datalad remove [-h] [-d DATASET] [-r] [--nocheck] [--nosave] [-m MESSAGE] [--if-dirty
↪{fail,save-before,ignore}] [PATH [PATH ...]]
```

### Description

Remove components from datasets

This command can remove subdatasets and paths, including non-empty directories, from datasets. Removing a component implies dropping present content and uninstalling associated subdatasets. Subsequently, the component is “unregistered” from the respective dataset. This means that the component is no longer present on the file system.

By default, the availability of at least one remote copy is verified before file content is dropped. As these checks could lead to slow operation (network latencies, etc), they can be disabled.

#### Examples

Permanently remove a subdataset from a dataset and wipe out the subdataset association too:

```
% datalad remove -d <path/to/dataset> <path/to/subds>
```

Permanently remove a dataset and all subdatasets:

```
% datalad remove -d <path/to/dataset/> --recursive
```

Permanently remove a dataset and all subdatasets even if there are fewer than the configured minimum number of (remote) sources for data:

```
% datalad remove -d <path/to/dataset/> --recursive --nocheck
```

### Options

#### PATH

path/name of the component to be removed. Constraints: value must be a string

**-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

**-d DATASET, --dataset DATASET**

specify the dataset to perform the operation on. If no dataset is given, an attempt is made to identify a dataset based on the `PATH` given. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-r, --recursive**

if set, recurse into potential subdataset.

**--nocheck**

whether to perform checks to assure the configured minimum number (remote) source for data. Give this option to skip checks.

**--nosave**

by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior.

**-m MESSAGE, --message MESSAGE**

a description of the state or the changes made to a dataset. Constraints: value must be a string

**--if-dirty {fail,save-before,ignore}**

desired behavior if a dataset with unsaved changes is discovered: ‘fail’ will trigger an error and further processing is aborted; ‘save-before’ will save all changes prior any further action; ‘ignore’ let’s datalad proceed as if the dataset would not have unsaved changes. [Default: ‘save-before’]

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad save**

**Synopsis**

```
datalad save [-h] [-m MESSAGE] [-d DATASET] [-t ID] [-r] [-R LEVELS] [-u] [-F MESSAGE_
↪FILE] [--to-git] [PATH [PATH ...]]
```

## Description

Save the current state of a dataset

Saving the state of a dataset records changes that have been made to it. This change record is annotated with a user-provided description. Optionally, an additional tag, such as a version, can be assigned to the saved state. Such tag enables straightforward retrieval of past versions at a later point in time.

**NOTE** Before Git v2.22, any Git repository without an initial commit located inside a Dataset is ignored, and content underneath it will be saved to the respective superdataset. DataLad datasets always have an initial commit, hence are not affected by this behavior.

### Examples

Save any content underneath the current directory, without altering any potential subdataset:

```
% datalad save .
```

Save specific content in the dataset:

```
% datalad save myfile.txt
```

Attach a commit message to save:

```
% datalad save -m 'add file' myfile.txt
```

Save any content underneath the current directory, and recurse into any potential subdatasets:

```
% datalad save . --recursive
```

Save any modification of known dataset content in the current directory, but leave untracked files (e.g. temporary files) untouched:

```
% datalad save -u .
```

Tag the most recent saved state of a dataset:

```
% datalad save --version-tag 'bestyet'
```

## Options

### PATH

path/name of the dataset component to save. If given, only changes made to those components are recorded in the new state. Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-m MESSAGE, --message MESSAGE**

a description of the state or the changes made to a dataset. Constraints: value must be a string

**-d DATASET, --dataset DATASET**

“specify the dataset to save. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-t ID, --version-tag ID**

an additional marker for that state. Every dataset that is touched will receive the tag. Constraints: value must be a string

**-r, --recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type ‘int’

**-u, --updated**

if given, only saves previously tracked paths.

**-F MESSAGE\_FILE, --message-file MESSAGE\_FILE**

take the commit message from this file. This flag is mutually exclusive with -m. Constraints: value must be a string

**--to-git**

flag whether to add data directly to Git, instead of tracking data identity only. Usually this is not desired, as it inflates dataset sizes and impacts flexibility of data transport. If not specified - it will be up to git-annex to decide, possibly on .gitattributes options. Use this flag with a simultaneous selection of paths to save. In general, it is better to pre-configure a dataset to track particular paths, file types, or file sizes with either Git or git-annex. See <https://git-annex.branchable.com/tips/largefiles/>.

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad update**

**Synopsis**

```
datalad update [-h] [-s SIBLING] [--merge [ALLOWED]] [--follow {sibling|parentds}] [-  
↪d DATASET] [-r] [-R LEVELS] [--fetch-all] [--reobtain-data] [PATH [PATH ...]]
```

## Description

Update a dataset from a sibling.

### Examples

Update from a particular sibling:

```
% datalad update -s <siblingname>
```

Update from a particular sibling and merge the changes from a configured or matching branch from the sibling (see `--follow` for details):

```
% datalad update --merge -s <siblingname>
```

Update from the sibling 'origin', traversing into subdatasets. For subdatasets, merge the revision registered in the parent dataset into the current branch:

```
% datalad update -s origin --merge --follow=parentds --recursive
```

## Options

### PATH

constrain to-be-updated subdatasets to the given path for recursive operation. Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-s *SIBLING*, --sibling *SIBLING***

name of the sibling to update from. If no sibling is given, updates from all siblings are obtained. Constraints: value must be a string

### **--merge [ALLOWED]**

merge obtained changes from the sibling. If a sibling is not explicitly given and there is only a single known sibling, that sibling is used. Otherwise, an unspecified sibling defaults to the configured remote for the current branch. By default, changes are fetched from the sibling but not merged into the current branch. With `--merge` or `--merge=any`, the changes will be merged into the current branch. A value of 'ff-only' restricts the allowed merges to fast-forwards. Constraints: value must be convertible to type bool, or value must be one of ('any', 'ff-only') [Default: False]

### **--follow {sibling|parentds}**

source of updates for subdatasets. For 'sibling', the update will be done by merging in a branch from the (specified or inferred) sibling. The branch brought in will either be the current branch's configured branch, if it points to a branch that belongs to the sibling, or a sibling branch with a name that matches the current branch. For 'parentds', the revision registered in the parent dataset of the subdataset is merged in. Note that the current dataset is always updated according

to 'sibling'. This option has no effect unless a merge is requested and `--recursive` is specified. Constraints: value must be one of ('sibling', 'parentds') [Default: 'sibling']

### **-d DATASET, --dataset DATASET**

specify the dataset to update. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-r, --recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

### **--fetch-all**

this option has no effect and will be removed in a future version. When no siblings are given, an all-sibling update will be performed.

### **--reobtain-data**

if enabled, file content that was present before an update will be re-obtained in case a file was changed by the update.

## **Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## **datalad uninstall**

### **Synopsis**

```
datalad uninstall [-h] [-d DATASET] [-r] [--nocheck] [--if-dirty {fail,save-before,  
→ignore}] [PATH [PATH ...]]
```

### **Description**

Uninstall subdatasets

This command can be used to uninstall any number of installed subdatasets. This command will error if individual files or non-dataset directories are given as input (use the drop or remove command depending on the desired goal), nor will it uninstall top-level datasets (i.e. datasets that are not a subdataset in another dataset; use the remove command for this purpose).



By default, the availability of at least one remote copy for each currently available file in any dataset is verified. As these checks could lead to slow operation (network latencies, etc), they can be disabled.

Any number of paths to process can be given as input. Recursion into subdatasets needs to be explicitly enabled, while recursion into subdirectories within a dataset is done automatically. An optional recursion limit is applied relative to each given input path.

### Examples

Uninstall a subdataset (undo installation):

```
% datalad uninstall <path/to/subds>
```

Uninstall a subdataset and all potential subdatasets:

```
% datalad uninstall -r <path/to/subds>
```

Skip checks that ensure a minimal number of (remote) sources:

```
% datalad uninstall <path/to/subds> --nocheck
```

## Options

### PATH

path/name of the component to be uninstalled. Constraints: value must be a string

### -h, --help, --help-np

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

### -d DATASET, --dataset DATASET

specify the dataset to perform the operation on. If no dataset is given, an attempt is made to identify a dataset based on the PATH given. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### -r, --recursive

if set, recurse into potential subdataset.

### --nocheck

whether to perform checks to assure the configured minimum number (remote) source for data. Give this option to skip checks.

### --if-dirty {fail,save-before,ignore}

desired behavior if a dataset with unsaved changes is discovered: 'fail' will trigger an error and further processing is aborted; 'save-before' will save all changes prior any further action; 'ignore' let's datalad proceed as if the dataset would not have unsaved changes. [Default: 'save-before']

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad unlock

### Synopsis

```
datalad unlock [-h] [-d DATASET] [-r] [-R LEVELS] [path [path ...]]
```

### Description

Unlock file(s) of a dataset

Unlock files of a dataset in order to be able to edit the actual content

#### *Examples*

Unlock a single file:

```
% datalad unlock <path/to/file>
```

Unlock all contents in the dataset:

```
% datalad unlock .
```

### Options

#### **path**

file(s) to unlock. Constraints: value must be a string

#### **-h, -help, -help-np**

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

#### **-d DATASET, -dataset DATASET**

“specify the dataset to unlock files in. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

#### **-r, -recursive**

if set, recurse into potential subdataset.

## **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

## **Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## **Metadata handling**

### **datalad search**

## **Synopsis**

```
datalad search [-h] [-d DATASET] [--reindex] [--max-nresults MAX_NRESULTS] [--mode
↪ {egrep, textblob, autofield}] [--full-record] [--show-keys {name, short, full}] [--show-
↪ query] [QUERY [QUERY ...]]
```

## **Description**

Search dataset metadata

DataLad can search metadata extracted from a dataset and/or aggregated into a superdataset (see the AGGREGATE-METADATA command). This makes it possible to discover datasets, or individual files in a dataset even when they are not available locally.

Ultimately DataLad metadata are a graph of linked data structures. However, this command does not (yet) support queries that can exploit all information stored in the metadata. At the moment the following search modes are implemented that represent different trade-offs between the expressiveness of a query and the computational and storage resources required to execute a query.

- `egrep` (default)
- `egrepcs` [case-sensitive `egrep`]
- `textblob`
- `autofield`

An alternative default mode can be configured by tuning the configuration variable 'datalad.search.default-mode':

```
[datalad "search"]
  default-mode = egrepcs
```

Each search mode has its own default configuration for what kind of documents to query. The respective default can be changed via configuration variables:

```
[datalad "search"]
  index-<mode_name>-documenttype = (all|datasets|files)
```

*Mode: `egrep/egrepcs`*

These search modes are largely ignorant of the metadata structure, and simply perform matching of a search pattern against a flat string-representation of metadata. This is advantageous when the query is simple and the metadata

structure is irrelevant, or precisely known. Moreover, it does not require a search index, hence results can be reported without an initial latency for building a search index when the underlying metadata has changed (e.g. due to a dataset update). By default, these search modes only consider datasets and do not investigate records for individual files for speed reasons. Search results are reported in the order in which they were discovered.

Queries can make use of Python regular expression syntax (<https://docs.python.org/3/library/re.html>). In EGREP mode, matching is case-insensitive when the query does not contain upper case characters, but is case-sensitive when it does. In EGREPCS mode, matching is always case-sensitive. Expressions will match anywhere in a metadata string, not only at the start.

When multiple queries are given, all queries have to match for a search hit (AND behavior).

It is possible to search individual metadata key/value items by prefixing the query with a metadata key name, separated by a colon (':'). The key name can also be a regular expression to match multiple keys. A query match happens when any value of an item with a matching key name matches the query (OR behavior). See examples for more information.

Examples:

Query for (what happens to be) an author:

```
% datalad search haxby
```

Queries are case-INsensitive when the query contains no upper case characters, and can be regular expressions. Use EGREPCS mode when it is desired to perform a case-sensitive lowercase match:

```
% datalad search --mode egrepcs halchenko.*haxby
```

This search mode performs NO analysis of the metadata content. Therefore queries can easily fail to match. For example, the above query implicitly assumes that authors are listed in alphabetical order. If that is the case (which may or may not be true), the following query would yield NO hits:

```
% datalad search Haxby.*Halchenko
```

The TEXTBLOB search mode represents an alternative that is more robust in such cases.

For more complex queries multiple query expressions can be provided that all have to match to be considered a hit (AND behavior). This query discovers all files (non-default behavior) that match 'bids.type=T1w' AND 'nifti1.qform\_code=scanner':

```
% datalad -c datalad.search.index-egrep-documenttype=all search bids.  
↪type:T1w nifti1.qform_code:scanner
```

Key name selectors can also be expressions, which can be used to select multiple keys or construct “fuzzy” queries. In such cases a query matches when any item with a matching key matches the query (OR behavior). However, multiple queries are always evaluated using an AND conjunction. The following query extends the example above to match any files that have either 'nifti1.qform\_code=scanner' or 'nifti1.sform\_code=scanner':

```
% datalad -c datalad.search.index-egrep-documenttype=all search bids.  
↪type:T1w nifti1.(q|s)form_code:scanner
```

*Mode: textblob*

This search mode is very similar to the EGREP mode, but with a few key differences. A search index is built from the string-representation of metadata records. By default, only datasets are included in this index, hence the indexing is usually completed within a few seconds, even for hundreds of datasets. This mode uses its own query language (not regular expressions) that is similar to other search engines. It supports logical conjunctions and fuzzy search terms. More information on this is available from the Whoosh project (search engine implementation):

- Description of the Whoosh query language: <http://whoosh.readthedocs.io/en/latest/querylang.html>

- Description of a number of query language customizations that are enabled in DataLad, such as, fuzzy term matching: <http://whoosh.readthedocs.io/en/latest/parsing.html#common-customizations>

Importantly, search hits are scored and reported in order of descending relevance, hence limiting the number of search results is more meaningful than in the ‘egrep’ mode and can also reduce the query duration.

Examples:

Search for (what happens to be) two authors, regardless of the order in which those names appear in the metadata:

```
% datalad search --mode textblob halchenko haxby
```

Fuzzy search when you only have an approximate idea what you are looking for or how it is spelled:

```
% datalad search --mode textblob haxbi~
```

Very fuzzy search, when you are basically only confident about the first two characters and how it sounds approximately (or more precisely: allow for three edits and require matching of the first two characters):

```
% datalad search --mode textblob haksbi~3/2
```

Combine fuzzy search with logical constructs:

```
% datalad search --mode textblob 'haxbi~ AND (hanke OR halchenko)'
```

### Mode: autofield

This mode is similar to the ‘textblob’ mode, but builds a vastly more detailed search index that represents individual metadata variables as individual fields. By default, this search index includes records for datasets and individual fields, hence it can grow very quickly into a huge structure that can easily take an hour or more to build and require more than a GB of storage. However, limiting it to documents on datasets (see above) retains the enhanced expressiveness of queries while dramatically reducing the resource demands.

Examples:

List names of search index fields (auto-discovered from the set of indexed datasets):

```
% datalad search --mode autofield --show-keys name
```

Fuzzy search for datasets with an author that is specified in a particular metadata field:

```
% datalad search --mode autofield bids.author:haxbi~ type:dataset
```

Search for individual files that carry a particular description prefix in their ‘nifti1’ metadata:

```
% datalad search --mode autofield nifti1.description:FSL* type:file
```

### Reporting

Search hits are returned as standard DataLad results. On the command line the ‘–output-format’ (or ‘-f’) option can be used to tweak results for further processing.

Examples:

Format search hits as a JSON stream (one hit per line):

```
% datalad -f json search haxby
```

Custom formatting: which terms matched the query of particular results. Useful for investigating fuzzy search results:

```
$ datalad -f '{path}: {query_matched}' search --mode autofield bids.  
↪author:haxbi~
```

## Options

### QUERY

query string, supported syntax and features depends on the selected search mode (see documentation).

#### **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

#### **-d DATASET, -dataset DATASET**

specify the dataset to perform the query operation on. If no dataset is given, an attempt is made to identify the dataset based on the current working directory and/or the PATH given. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

#### **-reindex**

force rebuilding the search index, even if no change in the dataset's state has been detected, for example, when the index documenttype configuration has changed.

#### **-max-nresults MAX\_NRESULTS**

maximum number of search results to report. Setting this to 0 will report all search matches. Depending on the mode this can search substantially slower. If not specified, a mode-specific default setting will be used. Constraints: value must be convertible to type 'int'

#### **-mode {egrep,textblob,autofield}**

Mode of search index structure and content. See section SEARCH MODES for details.

#### **-full-record, -f**

If set, return the full metadata record for each search hit. Depending on the search mode this might require additional queries. By default, only data that is available to the respective search modes is returned. This always includes essential information, such as the path and the type.

**–show-keys {name,short,full}**

if given, a list of known search keys is shown. If ‘name’ - only the name is printed one per line. If ‘short’ or ‘full’, statistics (in how many datasets, and how many unique values) are printed. ‘short’ truncates the listing of unique values. No other action is performed (except for reindexing), even if other arguments are given. Each key is accompanied by a term definition in parenthesis (TODO). In most cases a definition is given in the form of a URL. If an ontology definition for a term is known, this URL can resolve to a webpage that provides a comprehensive definition of the term. However, for speed reasons term resolution is solely done on information contained in a local dataset’s metadata, and definition URLs might be outdated or point to no longer existing resources.

**–show-query**

if given, the formal query that was generated from the given query string is shown, but not actually executed. This is mostly useful for debugging purposes.

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

**datalad metadata****Synopsis**

```
datalad metadata [-h] [-d DATASET] [--get-aggregates] [--reporton TYPE] [-r] [PATH_
↪ [PATH ...]]
```

**Description**

Metadata reporting for files and entire datasets

Two types of metadata are supported:

1. metadata describing a dataset as a whole (dataset-global metadata), and
2. metadata for files in a dataset (content metadata).

Both types can be accessed with this command.

Examples:

Report the metadata of a single file, as aggregated into the closest locally available dataset, containing the query path:

```
% datalad metadata somedir/subdir/thisfile.dat
```

Sometimes it is helpful to get metadata records formatted in a more accessible form, here as pretty-printed JSON:

```
% datalad -f json_pp metadata somedir/subdir/thisfile.dat
```

Same query as above, but specify which dataset to query (must be containing the query path):

```
% datalad metadata -d . somedir/subdir/thisfile.dat
```

Report any metadata record of any dataset known to the queried dataset:

```
% datalad metadata --recursive --reporton datasets
```

Get a JSON-formatted report of aggregated metadata in a dataset, incl. information on enabled metadata extractors, dataset versions, dataset IDs, and dataset paths:

```
% datalad -f json metadata --get-aggregates
```

## Options

### PATH

path(s) to query metadata for. Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-d DATASET, --dataset DATASET**

dataset to query. If given, metadata will be reported as stored in this dataset. Otherwise, the closest available dataset containing a query path will be consulted. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **--get-aggregates**

if set, yields all (sub)datasets for which aggregate metadata are available in the dataset. No other action is performed, even if other arguments are given. The reported results contain a datasets's ID, the commit hash at which metadata aggregation was performed, and the location of the object file(s) containing the aggregated metadata.

### **--reporton TYPE**

choose on what type result to report on: 'datasets', 'files', 'all' (both datasets and files), or 'none' (no report). Constraints: value must be one of ('all', 'datasets', 'files', 'none') [Default: 'all']

### **-r, --recursive**

if set, recurse into potential subdataset.

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.



## datalad aggregate-metadata

### Synopsis

```
datalad aggregate-metadata [-h] [-d DATASET] [-r] [-R LEVELS] [--update-mode
→{all|target}] [--incremental] [--force-extraction] [--nosave] [PATH [PATH ...]]
```

### Description

Aggregate metadata of one or more datasets for later query.

Metadata aggregation refers to a procedure that extracts metadata present in a dataset into a portable representation that is stored a single standardized format. Moreover, metadata aggregation can also extract metadata in this format from one dataset and store it in another (super)dataset. Based on such collections of aggregated metadata it is possible to discover particular datasets and specific parts of their content, without having to obtain the target datasets first (see the DataLad ‘search’ command).

To enable aggregation of metadata that are contained in files of a dataset, one has to enable one or more metadata extractor for a dataset. DataLad supports a number of common metadata standards, such as the Exchangeable Image File Format (EXIF), Adobe’s Extensible Metadata Platform (XMP), and various audio file metadata systems like ID3. DataLad extension packages can provide metadata data extractors for additional metadata sources. For example, the neuroimaging extension provides extractors for scientific (meta)data standards like BIDS, DICOM, and NIfTI1. Some metadata extractors depend on particular 3rd-party software. The list of metadata extractors available to a particular DataLad installation is reported by the ‘wtf’ command (‘datalad wtf’).

Enabling a metadata extractor for a dataset is done by adding its name to the ‘datalad.metadata.nativetype’ configuration variable – typically in the dataset’s configuration file (‘.datalad/config’), e.g.:

```
[datalad "metadata"]
  nativetype = exif
  nativetype = xmp
```

If an enabled metadata extractor is not available in a particular DataLad installation, metadata extraction will not succeed in order to avoid inconsistent aggregation results.

Enabling multiple extractors is supported. In this case, metadata are extracted by each extractor individually, and stored alongside each other. Metadata aggregation will also extract DataLad’s own metadata (extractors ‘datalad\_core’, and ‘annex’).

Metadata aggregation can be performed recursively, in order to aggregate all metadata across all subdatasets, for example, to be able to search across any content in any dataset of a collection. Aggregation can also be performed for subdatasets that are not available locally. In this case, pre-aggregated metadata from the closest available superdataset will be considered instead.

Depending on the versatility of the present metadata and the number of dataset or files, aggregated metadata can grow prohibitively large. A number of configuration switches are provided to mitigate such issues.

**datalad.metadata.aggregate-content-<extractor-name>** If set to false, content metadata aggregation will not be performed for the named metadata extractor (a potential underscore ‘\_’ in the extractor name must be replaced by a dash ‘-’). This can substantially reduce the runtime for metadata extraction, and also reduce the size of the generated metadata aggregate. Note, however, that some extractors may not produce any metadata when this is disabled, because their metadata might come from individual file headers only. ‘datalad.metadata.store-aggregate-content’ might be a more appropriate setting in such cases.

**datalad.metadata.aggregate-ignore-fields** Any metadata key matching any regular expression in this configuration setting is removed prior to generating the dataset-level metadata summary (keys and their unique values across

all dataset content), and from the dataset metadata itself. This switch can also be used to filter out sensitive information prior aggregation.

**datalad.metadata.generate-unique-<extractor-name>** If set to false, DataLad will not auto-generate a summary of unique content metadata values for a particular extractor as part of the dataset-global metadata (a potential underscore '\_' in the extractor name must be replaced by a dash '-'). This can be useful if such a summary is bloated due to minor uninformative (e.g. numerical) differences, or when a particular extractor already provides a carefully designed content metadata summary.

**datalad.metadata.maxfieldsize** Any metadata value that exceeds the size threshold given by this configuration setting (in bytes/characters) is removed.

**datalad.metadata.store-aggregate-content** If set, extracted content metadata are still used to generate a dataset-level summary of present metadata (all keys and their unique values across all files in a dataset are determined and stored as part of the dataset-level metadata aggregate, see `datalad.metadata.generate-unique-<extractor-name>`), but metadata on individual files are not stored. This switch can be used to avoid prohibitively large metadata files. Discovery of datasets containing content matching particular metadata properties will still be possible, but such datasets would have to be obtained first in order to discover which particular files in them match these properties.

## Options

### PATH

path to datasets that shall be aggregated. When a given path is pointing into a dataset, the metadata of the containing dataset will be aggregated. If no paths given, current dataset metadata is aggregated. Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-d DATASET, --dataset DATASET**

topmost dataset metadata will be aggregated into. All dataset between this dataset and any given path will receive updated aggregated metadata from all given paths. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-r, --recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**-update-mode {all|target}**

which datasets to update with newly aggregated metadata: all datasets from any leaf dataset to the top-level target dataset including all intermediate datasets (all), or just the top-level target dataset (target). Constraints: value must be one of ('all', 'target') [Default: 'target']

**-incremental**

If set, all information on metadata records of subdatasets that have not been (re-)aggregated in this run will be kept unchanged. This is useful when (re-)aggregation only a subset of a dataset hierarchy, for example, because not all subdatasets are locally available.

**-force-extraction**

If set, all enabled extractors will be engaged regardless of whether change detection indicates that metadata has already been extracted for a given dataset state.

**-nosave**

by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior.

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

**datalad extract-metadata****Synopsis**

```
datalad extract-metadata [-h] --type NAME [-d DATASET] [FILE [FILE ...]]
```

**Description**

Run one or more of DataLad's metadata extractors on a dataset or file.

The result(s) are structured like the metadata DataLad would extract during metadata aggregation. There is one result per dataset/file.

Examples:

Extract metadata with two extractors from a dataset in the current directory and also from all its files:

```
$ datalad extract-metadata -d . --type frictionless_datapackage --type ↵
↵datalad_core
```

Extract XMP metadata from a single PDF that is not part of any dataset:

```
$ datalad extract-metadata --type xmp Downloads/freshfromtheweb.pdf
```

## Options

### FILE

Path of a file to extract metadata from. Constraints: value must be a string

### **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

### **-type NAME**

Name of a metadata extractor to be executed. This option can be given more than once.

### **-d DATASET, -dataset DATASET**

“Dataset to extract metadata from. If no FILE is given, metadata is extracted from all files of the dataset. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

## Authors

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## Reproducible execution

### **datalad run**

## Synopsis

```
datalad run [-h] [-d DATASET] [-i PATH] [-o PATH] [--expand {inputs|outputs|both}] [--  
↪explicit] [-m MESSAGE] [--sidecar {yes|no}] ...
```

## Description

Run an arbitrary shell command and record its impact on a dataset.

It is recommended to craft the command such that it can run in the root directory of the dataset that the command will be recorded in. However, as long as the command is executed somewhere underneath the dataset root, the exact location will be recorded relative to the dataset root.

If the executed command did not alter the dataset in any way, no record of the command execution is made.

If the given command errors, a `COMMANDERROR` exception with the same exit code will be raised, and no modifications will be saved.

### *Command format*

A few placeholders are supported in the command via Python format specification. “`{pwd}`” will be replaced with the full path of the current working directory. “`{dspath}`” will be replaced with the full path of the dataset that run

is invoked on. “{tmpdir}” will be replaced with the full path of a temporary directory. “{inputs}” and “{outputs}” represent the values specified by `-input` and `-output`. If multiple values are specified, the values will be joined by a space. The order of the values will match that order from the command line, with any globs expanded in alphabetical order (like bash). Individual values can be accessed with an integer index (e.g., “{inputs[0]}”).

Note that the representation of the inputs or outputs in the formatted command string depends on whether the command is given as a list of arguments or as a string (quotes surrounding the command). The concatenated list of inputs or outputs will be surrounded by quotes when the command is given as a list but not when it is given as a string. This means that the string form is required if you need to pass each input as a separate argument to a preceding script (i.e., write the command as `./script {inputs}`, quotes included). The string form should also be used if the input or output paths contain spaces or other characters that need to be escaped.

To escape a brace character, double it (i.e., “{{” or “}}”).

Custom placeholders can be added as configuration variables under “`datalad.run.substitutions`”. As an example:

Add a placeholder “name” with the value “joe”:

```
% git config --file=.datalad/config datalad.run.substitutions.name joe
% datalad add -m "Configure name placeholder" .datalad/config
```

Access the new placeholder in a command:

```
% datalad run "echo my name is {name} >me"
```

### Examples

Run an executable script and record the impact on a dataset:

```
% datalad run -m 'run my script' 'code/script.sh'
```

Run a command and specify a directory as a dependency for the run. The contents of the dependency will be retrieved prior to running the script:

```
% datalad run -m 'run my script' --input 'data/*' 'code/script.sh'
```

Run an executable script and specify output files of the script to be unlocked prior to running the script:

```
% datalad run -m 'run my script' --input 'data/*' \
  --output 'output_dir/*' 'code/script.sh'
```

Specify multiple inputs and outputs:

```
% datalad run -m 'run my script' --input 'data/*' \
  --input 'datafile.txt' --output 'output_dir/*' --output \
  'outfile.txt' 'code/script.sh'
```

## Options

### COMMAND

command for execution. A leading ‘-’ can be used to disambiguate this command from the preceding options to DataLad.

**-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

**-d DATASET, --dataset DATASET**

specify the dataset to record the command results in. An attempt is made to identify the dataset based on the current working directory. If a dataset is given, the command will be executed in the root directory of this dataset. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-i PATH, --input PATH**

A dependency for the run. Before running the command, the content of this file will be retrieved. A value of “.” means “run datalad get .”. The value can also be a glob. This option can be given more than once.

**-o PATH, --output PATH**

Prepare this file to be an output file of the command. A value of “.” means “run datalad unlock .” (and will fail if some content isn’t present). For any other value, if the content of this file is present, unlock the file. Otherwise, remove it. The value can also be a glob. This option can be given more than once.

**--expand {inputs|outputs|both}**

Expand globs when storing inputs and/or outputs in the commit message. Constraints: value must be one of (None, ‘inputs’, ‘outputs’, ‘both’)

**--explicit**

Consider the specification of inputs and outputs to be explicit. Don’t warn if the repository is dirty, and only save modifications to the listed outputs.

**-m MESSAGE, --message MESSAGE**

a description of the state or the changes made to a dataset. Constraints: value must be a string

**--sidecar {yes|no}**

By default, the configuration variable ‘datalad.run.record-sidecar’ determines whether a record with information on a command’s execution is placed into a separate record file instead of the commit message (default: off). This option can be used to override the configured behavior on a case-by-case basis. Sidecar files are placed into the dataset’s ‘.datalad/runinfo’ directory (customizable via the ‘datalad.run.record-directory’ configuration variable). Constraints: value must be NONE, or value must be convertible to type bool

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad rerun

### Synopsis

```
datalad rerun [-h] [--since SINCE] [-d DATASET] [-b NAME] [-m MESSAGE] [--onto base]
↳ [--script FILE] [--report] [--explicit] [REVISION]
```

### Description

Re-execute previous *datalad run* commands.

This will unlock any dataset content that is on record to have been modified by the command in the specified revision. It will then re-execute the command in the recorded path (if it was inside the dataset). Afterwards, all modifications will be saved.

#### Report mode

When called with `--report`, this command reports information about what would be re-executed as a series of records. There will be a record for each revision in the specified revision range. Each of these will have one of the following “`rerun_action`” values:

- `run`: the revision has a recorded command that would be re-executed
- `skip-or-pick`: the revision does not have a recorded command and would be either skipped or cherry picked
- `merge`: the revision is a merge commit and a corresponding merge would be made

The decision to skip rather than cherry pick a revision is based on whether the revision would be reachable from `HEAD` at the time of execution.

In addition, when a starting point other than `HEAD` is specified, there is a `rerun_action` value “`checkout`”, in which case the record includes information about the revision the would be checked out before rerunning any commands.

**NOTE** Currently the “`onto`” feature only sets the working tree of the current dataset to a previous state. The working trees of any subdatasets remain unchanged.

#### Examples

Re-execute the command from the previous commit:

```
% datalad rerun
```

Re-execute any commands in the last five commits:

```
% datalad rerun --since=HEAD~5
```

Do the same as above, but re-execute the commands on top of `HEAD~5` in a detached state:

```
% datalad rerun --onto= --since=HEAD~5
```

Re-execute all previous commands and compare the old and new results:

```
% # on master branch
% datalad rerun --branch=verify --since=
% # now on verify branch
% datalad diff --revision=master..
% git log --oneline --left-right --cherry-pick master...
```

## Options

### REVISION

rerun command(s) in REVISION. By default, the command from this commit will be executed, but `--since` can be used to construct a revision range. Constraints: value must be a string [Default: 'HEAD']

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **--since SINCE**

If SINCE is a commit-ish, the commands from all commits that are reachable from REVISION but not SINCE will be re-executed (in other words, the commands in `git log SINCE..REVISION`). If SINCE is an empty string, it is set to the parent of the first commit that contains a recorded command (i.e., all commands in `git log REVISION` will be re-executed). Constraints: value must be a string

### **-d DATASET, --dataset DATASET**

specify the dataset from which to rerun a recorded command. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. If a dataset is given, the command will be executed in the root directory of this dataset. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-b NAME, --branch NAME**

create and checkout this branch before rerunning the commands. Constraints: value must be a string

### **-m MESSAGE, --message MESSAGE**

use MESSAGE for the reran commit rather than the recorded commit message. In the case of a multi-commit rerun, all the reran commits will have this message. Constraints: value must be a string

### **--onto base**

start point for rerunning the commands. If not specified, commands are executed at HEAD. This option can be used to specify an alternative start point, which will be checked out with the branch name specified by `--branch` or in a detached state otherwise. As a special case, an empty value for this option means the parent of the first run commit in the specified revision list. Constraints: value must be a string

### **--script FILE**

extract the commands into FILE rather than rerunning. Use `-` to write to stdout instead. This option implies `--report`. Constraints: value must be a string



**–report**

Don't actually re-execute anything, just display what would be done. Note: If you give this option, you most likely want to set `–output-format` to `'json'` or `'json_pp'`.

**–explicit**

Consider the specification of inputs and outputs in the run record to be explicit. Don't warn if the repository is dirty, and only save modifications to the outputs from the original record. Note that when several run commits are specified, this applies to every one. Care should also be taken when using `–onto` because checking out a new HEAD can easily fail when the working tree has modifications.

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad run-procedure****Synopsis**

```
datalad run-procedure [-h] [-d PATH] [--discover] [--help-proc] ...
```

**Description**

Run prepared procedures (DataLad scripts) on a dataset

*Concept*

A “procedure” is an algorithm with the purpose to process a dataset in a particular way. Procedures can be useful in a wide range of scenarios, like adjusting dataset configuration in a uniform fashion, populating a dataset with particular content, or automating other routine tasks, such as synchronizing dataset content with certain siblings.

Implementations of some procedures are shipped together with DataLad, but additional procedures can be provided by 1) any DataLad extension, 2) any (sub-)dataset, 3) a local user, or 4) a local system administrator. DataLad will look for procedures in the following locations and order:

Directories identified by the configuration settings

- `'datalad.locations.user-procedures'` (determined by `appdirs.user_config_dir`; defaults to `'$HOME/.config/datalad/procedures'` on GNU/Linux systems)
- `'datalad.locations.system-procedures'` (determined by `appdirs.site_config_dir`; defaults to `'/etc/xdg/datalad/procedures'` on GNU/Linux systems)
- `'datalad.locations.dataset-procedures'`

and subsequently in the `'resources/procedures/'` directories of any installed extension, and, lastly, of the DataLad installation itself.

Please note that a dataset that defines `'datalad.locations.dataset-procedures'` provides its procedures to any dataset it is a subdataset of. That way you can have a collection of such procedures in a dedicated dataset and install it as a subdataset into any dataset you want to use those procedures with. In case of a naming conflict with such a dataset hierarchy, the dataset you're calling run-procedures on will take precedence over its subdatasets and so on.

Each configuration setting can occur multiple times to indicate multiple directories to be searched. If a procedure matching a given name is found (filename without a possible extension), the search is aborted and this implementation will be executed. This makes it possible for individual datasets, users, or machines to override externally provided procedures (enabling the implementation of customizable processing “hooks”).

### *Procedure implementation*

A procedure can be any executable. Executables must have the appropriate permissions and, in the case of a script, must contain an appropriate “shebang” line. If a procedure is not executable, but its filename ends with ‘.py’, it is automatically executed by the ‘python’ interpreter (whichever version is available in the present environment). Likewise, procedure implementations ending on ‘.sh’ are executed via ‘bash’.

Procedures can implement any argument handling, but must be capable of taking at least one positional argument (the absolute path to the dataset they shall operate on).

For further customization there are two configuration settings per procedure available:

- ‘datalad.procedures.<NAME>.call-format’ fully customizable format string to determine how to execute procedure NAME (see also `datalad-run`). It currently requires to include the following placeholders:
  - ‘{script}’: will be replaced by the path to the procedure
  - ‘{ds}’: will be replaced by the absolute path to the dataset the procedure shall operate on
  - ‘{args}’: (not actually required) will be replaced by all additional arguments passed into `run-procedure` after NAME

As an example the default format string for a call to a python script is: “python {script} {ds} {args}”

- ‘datalad.procedures.<NAME>.help’ will be shown on `datalad run-procedure -help-proc NAME` to provide a description and/or usage info for procedure NAME

### *Examples*

Find out which procedures are available on the current system:

```
% datalad run-procedure --discover
```

Run the ‘yoda’ procedure in the current dataset:

```
% datalad run-procedure cfg_yoda
```

## **Options**

### **NAME [ARGS]**

Name and possibly additional arguments of the to-be-executed procedure. Note, that all options to `run-procedure` need to be put before NAME, since all ARGS get assigned to NAME.

### **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

### **-d PATH, -dataset PATH**

specify the dataset to run the procedure on. An attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

## –discover

if given, all configured paths are searched for procedures and one result record per discovered procedure is yielded, but no procedure is executed.

## –help-proc

if given, get a help message for procedure NAME from config setting datalad.procedures.NAME.help.

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## Miscellaneous commands

### datalad add-archive-content

## Synopsis

```
datalad add-archive-content [-h] [--annex ANNEX] [--add-archive-leading-dir] [--strip-
↳ leading-dirs] [--leading-dirs-depth LEADING_DIRS_DEPTH] [--leading-dirs-consider_
↳ LEADING_DIRS_CONSIDER] [--use-current-dir] [-d] [--key] [-e EXCLUDE] [-r RENAME] [--
↳ existing {fail,overwrite,archive-suffix,numeric-suffix}] [-o ANNEX_OPTIONS] [--
↳ copy] [--no-commit] [--allow-dirty] [--stats STATS] [--drop-after] [--delete-after]_
↳ archive
```

## Description

Add content of an archive under git annex control.

This results in the files within archive (which must be already under annex control itself) added under annex referencing original archive via custom special remotes mechanism

Example:

```
annex-repo$ datalad add-archive-content my_big_tarball.tar.gz
```

## Options

### archive

archive file or a key (if –key specified). Constraints: value must be a string

### -h, –help, –help-np

show this help message. –help-np forcefully disables the use of a pager for displaying the help message

**-annex *ANNEX***

annex instance to use.

**-add-archive-leading-dir**

flag to place extracted content under a directory which would correspond to archive name with suffix stripped. E.g. for archive EXAMPLE.ZIP its content will be extracted under a directory EXAMPLE/.

**-strip-leading-dirs**

flag to move all files directories up, from how they were stored in an archive, if that one contained a number (possibly more than 1 down) single leading directories.

**-leading-dirs-depth *LEADING\_DIRS\_DEPTH***

maximal depth to strip leading directories to. If not specified (None), no limit.

**-leading-dirs-consider *LEADING\_DIRS\_CONSIDER***

regular expression(s) for directories to consider to strip away. Constraints: value must be a string

**-use-current-dir**

flag to extract archive under the current directory, not the directory where archive is located. Note that it will be of no effect if `-key` is given.

**-d, -delete**

flag to delete original archive from the filesystem/git in current tree. Note that it will be of no effect if `-key` is given.

**-key**

flag to signal if provided archive is not actually a filename on its own but an annex key.

**-e *EXCLUDE*, -exclude *EXCLUDE***

regular expressions for filenames which to exclude from being added to annex. Applied after `-rename` if that one is specified. For exact matching, use anchoring. Constraints: value must be a string

**-r *RENAME*, -rename *RENAME***

regular expressions to rename files before being added under git. First letter defines how to split provided string into two parts: Python regular expression (with groups), and replacement string. Constraints: value must be a string

**–existing {fail,overwrite,archive-suffix,numeric-suffix}**

what operation to perform a file from archive tries to overwrite an existing file with the same name. ‘fail’ (default) leads to RuntimeError exception. ‘overwrite’ silently replaces existing file. ‘archive-suffix’ instructs to add a suffix (prefixed with a ‘-’) matching archive name from which file gets extracted, and if that one present, ‘numeric-suffix’ is in effect in addition, when incremental numeric suffix (prefixed with a ‘.’) is added until no name collision is longer detected. [Default: ‘fail’]

**-o ANNEX\_OPTIONS, –annex-options ANNEX\_OPTIONS**

additional options to pass to git-annex. Constraints: value must be a string

**–copy**

flag to copy the content of the archive instead of moving.

**–no-commit**

flag to not commit upon completion.

**–allow-dirty**

flag that operating on a dirty repository (uncommitted or untracked content) is ok.

**–stats STATS**

ActivityStats instance for global tracking.

**–drop-after**

drop extracted files after adding to annex.

**–delete-after**

extract under a temporary directory, git-annex add, and delete after. To be used to “index” files within annex without actually creating corresponding files under git. Note that *annex dropunused* would later remove that load.

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## **datalad clean**

### **Synopsis**

```
datalad clean [-h] [-d DATASET] [--what [{cached-archives,annex-tmp,annex-transfer,
↪search-index} [{cached-archives,annex-tmp,annex-transfer,search-index} ...]]] [-r]
↪ [-R LEVELS]
```

### **Description**

Clean up after DataLad (possible temporary files etc.)

Removes extracted temporary archives, etc.

Examples:

```
$ datalad clean
```

### **Options**

**-h, --help, --help-np**

show this help message. **--help-np** forcefully disables the use of a pager for displaying the help message

**-d DATASET, --dataset DATASET**

specify the dataset to perform the clean operation on. If no dataset is given, an attempt is made to identify the dataset in current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**--what [{cached-archives,annex-tmp,annex-transfer,search-index} [{cached-archives,annex-tmp,annex-transfer,search-index} ...]]**

What to clean. If none specified – all known targets are cleaned.

**-r, --recursive**

if set, recurse into potential subdataset.

**-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

### **Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## datalad download-url

### Synopsis

```
datalad download-url [-h] [-d PATH] [-O PATH] [-o] [--archive] [--nosave] [-m MESSAGE] url [url ...]
```

### Description

Download content

It allows for a uniform download interface to various supported URL schemes, re-using or asking for authentication details maintained by datalad.

#### Examples

Download files from an http and S3 URL:

```
% datalad download-url http://example.com/file.dat s3://bucket/file2.dat
```

Download a file to a path and provide a commit message:

```
% datalad download-url -m 'added a file' -O myfile.dat \
  s3://bucket/file2.dat
```

### Options

#### url

URL(s) to be downloaded. Constraints: value must be a string

#### -h, -help, -help-np

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

#### -d PATH, -dataset PATH

specify the dataset to add files to. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Use --nosave to prevent adding files to the dataset. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

#### -O PATH, -path PATH

target for download. If the path has a trailing separator, it is treated as a directory, and each specified URL is downloaded under that directory to a base name taken from the URL. Without a trailing separator, the value specifies the name of the downloaded file (file name extensions inferred from the URL may be added to it, if they are not yet present) and only a single URL should be given. In both cases, leading directories will be created if needed. This argument defaults to the current directory. Constraints: value must be a string

### **-O, --overwrite**

flag to overwrite it if target file exists.

### **--archive**

pass the downloaded files to `datalad add-archive-content --delete`.

### **--nosave**

by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior.

### **-m MESSAGE, --message MESSAGE**

a description of the state or the changes made to a dataset. Constraints: value must be a string

## **Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## **datalad ls**

### **Synopsis**

```
datalad ls [-h] [-r] [-F] [-a] [-L] [--config-file CONFIG_FILE] [--list-content {None,
↪first10,md5,full}] [--json {file,display,delete}] [PATH/URL [PATH/URL ...]]
```

### **Description**

List summary information about URLs and dataset(s)

ATM only `s3://` URLs and datasets are supported

Examples:

```
$ datalad ls s3://openfmri/tarballs/ds202 # to list S3 bucket $ datalad ls # to list current dataset
```

### **Options**

#### **PATH/URL**

URL or path to list, e.g. `s3://...` Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message



**-r, --recursive**

recurse into subdirectories.

**-F, --fast**

only perform fast operations. Would be overridden by `--all`.

**-a, --all**

list all (versions of) entries, not e.g. only latest entries in case of S3.

**-L, --long**

list more information on entries (e.g. acl, urls in s3, annex sizes etc).

**--config-file *CONFIG\_FILE***

path to config file which could help the 'ls'. E.g. for s3:// URLs could be some `~/s3cfg` file which would provide credentials. Constraints: value must be a string

**--list-content {None,first10,md5,full}**

list also the content or only first 10 bytes (first10), or md5 checksum of an entry. Might require expensive transfer and dump binary output to your screen. Do not enable unless you know what you are after. [Default: False]

**--json {file,display,delete}**

metadata json of dataset for creating web user interface. `display`: prints jsons to stdout or `file`: writes each subdir metadata to json file in subdir of dataset or `delete`: deletes all metadata json files in dataset.

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad test**

**Synopsis**

```
datalad test [-h] [-v] [-s] [--pdb] [-x] [module [module ...]]
```

## Description

Run internal DataLad (unit)tests.

This can be used to verify correct operation on the system. It is just a thin wrapper around a call to nose, so number of exposed options is minimal

## Options

### module

test name(s), by default all tests of DataLad core and any installed extensions are executed.

### -h, --help, --help-np

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

### -v, --verbose

be verbose - list test names.

### -s, --nocapture

do not capture stdout.

### --pdb

drop into debugger on failures or errors.

### -x, --stop

stop running tests after the first error or failure.

## Authors

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## Plugin commands

### datalad add-readme

## Synopsis

```
datalad add-readme [-h] [-d DATASET] [--existing {skip|append|replace}] [PATH]
```

## Description

Add basic information about DataLad datasets to a README file

The README file is added to the dataset and the addition is saved in the dataset.

## Options

### PATH

Path of the README file within the dataset. Constraints: value must be a string [Default: 'README.md']

### -h, -help, -help-np

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

### -d DATASET, -dataset DATASET

Dataset to add information to. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### -existing {skip|append|replace}

How to react if a file with the target name already exists: 'skip': do nothing; 'append': append information to the existing file; 'replace': replace the existing file with new content. Constraints: value must be one of ('skip', 'append', 'replace') [Default: 'skip']

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad addurls

### Synopsis

```
datalad addurls [-h] [-d DATASET] [-t TYPE] [-x REGEXP] [-m FORMAT] [--message_
↪MESSAGE] [-n] [--fast] [--ifexists {overwrite|skip}] [--missing-value VALUE] [--
↪nosave] [--version-urls] [-c PROC] URL-FILE URL-FORMAT FILENAME-FORMAT
```

## Description

Create and update a dataset from a list of URLs.

*Format specification*

Several arguments take format strings. These are similar to normal Python format strings where the names from URL-FILE (column names for a CSV or properties for JSON) are available as placeholders. If URL-FILE is a CSV file, a positional index can also be used (i.e., “{0}” for the first column). Note that a placeholder cannot contain a ‘:’ or ‘!’.

In addition, the FILENAME-FORMAT arguments has a few special placeholders.

- `_reindex`

The constructed file names must be unique across all fields rows. To avoid collisions, the special placeholder “\_reindex” can be added to the formatter. Its value will start at 0 and increment every time a file name repeats.

- `_url_hostname`, `_urlN`, `_url_basename*`

Various parts of the formatted URL are available. Take “[http://datalad.org/asciicast/seamless\\_nested\\_repos.sh](http://datalad.org/asciicast/seamless_nested_repos.sh)” as an example.

“datalad.org” is stored as “\_url\_hostname”. Components of the URL’s path can be referenced as “\_urlN”. “\_url0” and “\_url1” would map to “asciicast” and “seamless\_nested\_repos.sh”, respectively. The final part of the path is also available as “\_url\_basename”.

This name is broken down further. “\_url\_basename\_root” and “\_url\_basename\_ext” provide access to the root name and extension. These values are similar to the result of `os.path.splitext`, but, in the case of multiple periods, the extension is identified using the same length heuristic that `git-annex` uses. As a result, the extension of “file.tar.gz” would be “.tar.gz”, not “.gz”. In addition, the fields “\_url\_basename\_root\_py” and “\_url\_basename\_ext\_py” provide access to the result of `os.path.splitext`.

- `_url_filename*`

These are similar to `_url_basename*` fields, but they are obtained with a server request. This is useful if the file name is set in the Content-Disposition header.

### Examples

Consider a file “avatars.csv” that contains:

```
who,ext,link
neurodebian,png,https://avatars3.githubusercontent.com/u/260793
datalad,png,https://avatars1.githubusercontent.com/u/8927200
```

To download each link into a file name composed of the ‘who’ and ‘ext’ fields, we could run:

```
$ datalad addurls -d avatar_ds --fast avatars.csv '{link}' '{who}.{ext}'
```

The `-d avatar_ds` is used to create a new dataset in “\$PWD/avatar\_ds”.

If we were already in a dataset and wanted to create a new subdataset in an “avatars” subdirectory, we could use “//” in the FILENAME-FORMAT argument:

```
$ datalad addurls --fast avatars.csv '{link}' 'avatars://{who}.{ext}'
```

### NOTE

For users familiar with ‘git annex addurl’: A large part of this plugin’s functionality can be viewed as transforming data from URL-FILE into a “url filename” format that fed to ‘git annex addurl -batch -with-files’.

## Options

## URL-FILE

A file that contains URLs or information that can be used to construct URLs. Depending on the value of `-input-type`, this should be a CSV file (with a header as the first row) or a JSON file (structured as a list of objects with string values).

## URL-FORMAT

A format string that specifies the URL for each entry. See the ‘Format Specification’ section above.

## FILENAME-FORMAT

Like URL-FORMAT, but this format string specifies the file to which the URL’s content will be downloaded. The name should be a relative path and will be taken as relative to the top-level dataset, regardless of whether it is specified via `-dataset`) or inferred. The file name may contain directories. The separator “//” can be used to indicate that the left-side directory should be created as a new subdataset. See the ‘Format Specification’ section above.

## **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

## **-d DATASET, -dataset DATASET**

Add the URLs to this dataset (or possibly subdatasets of this dataset). An empty or non-existent directory is passed to create a new dataset. New subdatasets can be specified with FILENAME-FORMAT. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

## **-t TYPE, -input-type TYPE**

Whether URL-FILE should be considered a CSV file or a JSON file. The default value, “ext”, means to consider URL-FILE as a JSON file if it ends with “.json”. Otherwise, treat it as a CSV file. Constraints: value must be one of (‘ext’, ‘csv’, ‘json’) [Default: ‘ext’]

## **-x REGEXP, -exclude-autometa REGEXP**

By default, metadata field=value pairs are constructed with each column in URL- FILE, excluding any single column that is specified via URL-FORMAT. This argument can be used to exclude columns that match a regular expression. If set to ‘\*’ or an empty string, automatic metadata extraction is disabled completely. This argument does not affect metadata set explicitly with `-meta`.

## **-m FORMAT, -meta FORMAT**

A format string that specifies metadata. It should be structured as “<field>=<value>”. As an example, “location={3}” would mean that the value for the “location” metadata field should be set the value of the fourth column. This option can be given multiple times.

**-message MESSAGE**

Use this message when committing the URL additions. Constraints: value must be NONE, or value must be a string

**-n, -dry-run**

Report which URLs would be downloaded to which files and then exit.

**-fast**

If True, add the URLs, but don't download their content. Underneath, this passes the `-fast` flag to *git annex addurl*.

**-ifexists {overwrite|skip}**

What to do if a constructed file name already exists. The default behavior is to proceed with the *git annex addurl*, which will fail if the file size has changed. If set to 'overwrite', remove the old file before adding the new one. If set to 'skip', do not add the new file. Constraints: value must be one of (None, 'overwrite', 'skip')

**-missing-value VALUE**

When an empty string is encountered, use this value instead. Constraints: value must be NONE, or value must be a string

**-nosave**

by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior.

**-version-urls**

Try to add a version ID to the URL. This currently only has an effect on HTTP URLs for AWS S3 buckets. `s3://` URL versioning is not yet supported, but any URL that already contains a "versionId=" parameter will be used as is.

**-c PROC, -cfg-proc PROC**

Pass this `-cfg_proc` value when calling `CREATE` to make datasets.

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## datalad check-dates

### Synopsis

```
datalad check-dates [-h] [-D DATE] [--rev REVISION] [--annex {all|tree|none}] [--no-
→tags] [--older] [PATH [PATH ...]]
```

### Description

Find repository dates that are more recent than a reference date.

The main purpose of this tool is to find “leaked” real dates in repositories that are configured to use fake dates. It checks dates from three sources: (1) commit timestamps (author and committer dates), (2) timestamps within files of the “git-annex” branch, and (3) the timestamps of annotated tags.

### Options

#### PATH

Root directory in which to search for Git repositories. The current working directory will be used by default. Constraints: value must be a string

#### -h, -help, -help-np

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

#### -D DATE, -reference-date DATE

Compare dates to this date. If dateutil is installed, this value can be any format that its parser recognizes. Otherwise, it should be a unix timestamp that starts with a “@”. The default value corresponds to 01 Jan, 2018 00:00:00 -0000. Constraints: value must be a string [Default: '@1514764800']

#### --rev REVISION

Search timestamps from commits that are reachable from REVISION. Any revision specification supported by git log, including flags like -all and -tags, can be used. This option can be given multiple times.

#### --annex {all|tree|none}

Mode for “git-annex” branch search. If ‘all’, all blobs within the branch are searched. ‘tree’ limits the search to blobs that are referenced by the tree at the tip of the branch. ‘none’ disables search of “git-annex” blobs. Constraints: value must be one of (‘all’, ‘tree’, ‘none’) [Default: ‘all’]

#### --no-tags

Don’t check the dates of annotated tags.

## **-older**

Find dates which are older than the reference date rather than newer.

## **Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## **datalad export-archive**

### **Synopsis**

```
datalad export-archive [-h] [-d DATASET] [-t {tar|zip}] [-c {gz|bz2|}] [--missing-  
→content {error|continue|ignore}] [PATH]
```

### **Description**

Export the content of a dataset as a TAR/ZIP archive.

### **Options**

#### **PATH**

File name of the generated TAR archive. If no file name is given the archive will be generated in the current directory and will be named: `datalad_<dataset_uuid>.(tar.*|zip)`. To generate that file in a different directory, provide an existing directory as the file name. Constraints: value must be a string

#### **-h, -help, -help-np**

show this help message. `-help-np` forcefully disables the use of a pager for displaying the help message

#### **-d DATASET, -dataset DATASET**

“specify the dataset to export. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

#### **-t {tar|zip}, -archivetype {tar|zip}**

Type of archive to generate. Constraints: value must be one of ('tar', 'zip') [Default: 'tar']

#### **-c {gz|bz2|}, -compression {gz|bz2|}**

Compression method to use. 'bz2' is not supported for ZIP archives. No compression is used when an empty string is given. Constraints: value must be one of ('gz', 'bz2', '') [Default: 'gz']



**-missing-content {error|continue|ignore}**

By default, any discovered file with missing content will result in an error and the export is aborted. Setting this to 'continue' will issue warnings instead of failing on error. The value 'ignore' will only inform about problem at the 'debug' log level. The latter two can be helpful when generating a TAR archive from a dataset where some file content is not available locally. Constraints: value must be one of ('error', 'continue', 'ignore') [Default: 'error']

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

**datalad export-to-figshare****Synopsis**

```
datalad export-to-figshare [-h] [-d DATASET] [--missing-content
↪{error|continue|ignore}] [--no-annex] [--article-id ID] [PATH]
```

**Description**

Export the content of a dataset as a ZIP archive to figshare

Very quick and dirty approach. Ideally figshare should be supported as a proper git annex special remote. Unfortunately, figshare does not support having directories, and can store only a flat list of files. That makes it impossible for any sensible publishing of complete datasets.

The only workaround is to publish dataset as a zip-ball, where the entire content is wrapped into a .zip archive for which figshare would provide a navigator.

**Options****PATH**

File name of the generated ZIP archive. If no file name is given the archive will be generated in the top directory of the dataset and will be named: datalad\_<dataset\_uuid>.zip. Constraints: value must be a string

**-h, -help, -help-np**

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

**-d DATASET, -dataset DATASET**

“specify the dataset to export. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **–missing-content {error|continue|ignore}**

By default, any discovered file with missing content will result in an error and the plugin is aborted. Setting this to ‘continue’ will issue warnings instead of failing on error. The value ‘ignore’ will only inform about problem at the ‘debug’ log level. The latter two can be helpful when generating a TAR archive from a dataset where some file content is not available locally. Constraints: value must be one of (‘error’, ‘continue’, ‘ignore’) [Default: ‘error’]

### **–no-annex**

By default the generated .zip file would be added to annex, and all files would get registered in git-annex to be available from such a tarball. Also upon upload we will register for that archive to be a possible source for it in annex. Setting this flag disables this behavior.

### **–article-id ID**

Which article to publish to. Constraints: value must be convertible to type ‘int’

## **Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## **datalad no-annex**

### **Synopsis**

```
datalad no-annex [-h] [-d DATASET] [--pattern PATTERN [PATTERN ...]] [--ref-dir REF_↵DIR] [--makedirs]
```

### **Description**

Configure a dataset to never put some content into the dataset’s annex

This can be useful in mixed datasets that also contain textual data, such as source code, which can be efficiently and more conveniently managed directly in Git.

Patterns generally look like this:

```
code/ *
```

which would match all file in the code directory. In order to match all files under CODE/, including all its subdirectories use such a pattern:

```
code/**
```

Note that the plugin works incrementally, hence any existing configuration (e.g. from a previous plugin run) is amended, not replaced.

## Parameters

ref\_dir : str, optional makedirs : bool, optional

## Options

**-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

**-d DATASET, --dataset DATASET**

“specify the dataset to configure. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**--pattern PATTERN [PATTERN ...]**

list of path patterns. Any content whose path is matching any pattern will not be annexed when added to a dataset, but instead will be tracked directly in Git. Path patterns have to be relative to the directory given by the REF\_DIR option. By default, patterns should be relative to the root of the dataset.

**--ref-dir REF\_DIR**

Relative path (within the dataset) to the directory that is to be configured. All patterns are interpreted relative to this path, and configuration is written to a .GITATTRIBUTES file in this directory. [Default: ‘.’]

**--makedirs**

If set, any missing directories will be created in order to be able to place a file into --REF-DIR.

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad wtf

### Synopsis

```
datalad wtf [-h] [-d DATASET] [-s {some|all}] [-S SECTION] [-D DECOR] [-c]
```

## Description

Generate a report about the DataLad installation and configuration

IMPORTANT: Sharing this report with untrusted parties (e.g. on the web) should be done with care, as it may include identifying information, and/or credentials or access tokens.

## Options

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-d DATASET, --dataset DATASET**

“specify the dataset to report on. no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-s {some|all}, --sensitive {some|all}**

if set to ‘some’ or ‘all’, it will display sections such as config and metadata which could potentially contain sensitive information (credentials, names, etc.). If ‘some’, the fields which are known to be sensitive will still be masked out. Constraints: value must be one of (None, ‘some’, ‘all’)

### **-S SECTION, --section SECTION**

section to include. If not set, all sections. This option can be given multiple times. Constraints: value must be one of (‘configuration’, ‘datalad’, ‘dataset’, ‘dependencies’, ‘environment’, ‘extensions’, ‘git-annex’, ‘location’, ‘meta-data\_extractors’, ‘python’, ‘system’)

### **-D DECOR, --decor DECOR**

decoration around the rendering to facilitate embedding into issues etc, e.g. use ‘html\_details’ for posting collapsable entry to GitHub issues. Constraints: value must be one of (‘html\_details’,)

### **-c, --clipboard**

if set, do not print but copy to clipboard (requires pyperclip module).

## Authors

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

## Plumbing commands

### datalad annotate-paths

#### Synopsis

```
datalad annotate-paths [-h] [-d DATASET] [-r] [-R LEVELS] [--action LABEL] [--
↪unavailable-path-status LABEL] [--unavailable-path-msg message] [--nondataset-path-
↪status LABEL] [--no-parentds-discovery] [--no-subds-discovery] [--revision-change-
↪discovery] [--no-untracked-discovery] [--modified [MODIFIED]] [PATH [PATH ...]]
```

#### Description

Analyze and act upon input paths

Given paths (or more generally location requests) are inspected and annotated with a number of properties. A list of recognized properties is provided below.

##### *Recognized path properties*

“**action**” label of the action that triggered the path annotation

“**annexkey**” annex key for the content of a file

“**logger**” logger for reporting a message

“**message**” message (plus possible tstring expansion arguments)

“**orig\_request**” original input by which a path was determined

“**parentds**” path of dataset containing the annotated path (superdataset for subdatasets)

“**path**” absolute path that is annotated

“**process\_content**” flag that content underneath the path is to be processed

“**process\_updated\_only**” flag that only known dataset components are to be processed

“**raw\_input**” flag whether this path was given as raw (non-annotated) input

“**refds**” path of a reference/base dataset the annotated path is part of

“**registered\_subds**” flag whether a dataset is known to be a true subdataset of PARENTDS

“**revision**” the recorded commit for a subdataset in a superdataset

“**revision\_descr**” a human-readable description of REVISION

“**source\_url**” URL a dataset was installed from

“**staged**” flag whether a path is known to be “staged” in its containing dataset

“**state**” state indicator for a path in its containing dataset (clean, modified, absent (also for files), conflict)

“**status**” action result status (ok, notneeded, impossible, error)

“**type**” nature of the path (file, directory, dataset)

“**url**” registered URL for a subdataset in a superdataset

In the case of enabled modification detection the results may contain additional properties regarding the nature of the modification. See the documentation of the DIFF command for details.

## Options

### **PATH**

path to be annotated. Constraints: value must be a string

### **-h, -help, -help-np**

show this help message. -help-np forcefully disables the use of a pager for displaying the help message

### **-d DATASET, -dataset DATASET**

an optional reference/base dataset for the paths. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-r, -recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, -recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

### **-action LABEL**

an "action" property value to include in the path annotation. Constraints: value must be a string

### **-unavailable-path-status LABEL**

a "status" property value to include in the annotation for paths that are underneath a dataset, but do not exist on the filesystem. Constraints: value must be a string [Default: '']

### **-unavailable-path-msg message**

a "message" property value to include in the annotation for paths that are underneath a dataset, but do not exist on the filesystem. Constraints: value must be a string

### **-nondataset-path-status LABEL**

a "status" property value to include in the annotation for paths that are not underneath any dataset. Constraints: value must be a string [Default: 'error']

**–no-parentds-discovery**

Flag to disable reports of parent dataset information for any path, in particular dataset root paths. Disabling saves on command run time, if this information is not needed.

**–no-subds-discovery**

Flag to disable reporting type='dataset' for subdatasets, even when they are not installed, or their mount point directory doesn't exist. Disabling saves on command run time, if this information is not needed.

**–revision-change-discovery**

Flag to disable discovery of changes which were not yet committed. Disabling saves on command run time, if this information is not needed.

**–no-untracked-discovery**

Flag to disable discovery of untracked changes. Disabling saves on command run time, if this information is not needed.

**–modified [*MODIFIED*]**

comparison reference specification for modification detection. This can be (mostly) anything that *git diff* understands (commit, treeish, tag, etc). See the documentation of *datalad diff –revision* for details. Unmodified paths will not be annotated. If a requested path was not modified but some content underneath it was, then the request is replaced by the modified paths and those are annotated instead. This option can be used without an argument to test against changes that have been made, but have not yet been staged for a commit. Constraints: value must be a string, or value must be convertible to type bool

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad clone****Synopsis**

```
datalad clone [-h] [-d DATASET] [-D DESCRIPTION] [--reckless [{auto|ephemeral}]]
↳SOURCE [PATH]
```

**Description**

Obtain a dataset (copy) from a URL or local directory

The purpose of this command is to obtain a new clone (copy) of a dataset and place it into a not-yet-existing or empty directory. As such CLONE provides a strict subset of the functionality offered by INSTALL. Only a single dataset can

be obtained, and immediate recursive installation of subdatasets is not supported. However, once a (super)dataset is installed via CLONE, any content, including subdatasets can be obtained by a subsequent GET command.

Primary differences over a direct *git clone* call are 1) the automatic initialization of a dataset annex (pure Git repositories are equally supported); 2) automatic registration of the newly obtained dataset as a subdataset (submodule), if a parent dataset is specified; and 3) support for additional resource identifiers (DataLad resource identifiers as used on datasets.datalad.org, and RIA store URLs as used for store.datalad.org; see examples); and 4) automatic configurable generation of alternative access URL for common cases (such as appending ‘.git’ to the URL in case the accessing the base URL failed).

SEEALSO

**handbook:3-001** (<http://handbook.datalad.org/symbols>) More information on Remote Indexed Archive (RIA) stores

*Examples*

Install a dataset from Github into the current directory:

```
% datalad clone https://github.com/datalad-datasets/longnow-podcasts.git
```

Install a dataset into a specific directory:

```
% datalad clone https://github.com/datalad-datasets/longnow-podcasts.git \
myfavpodcasts
```

Install a dataset as a subdataset into the current dataset:

```
% datalad clone -d . https://github.com/datalad-datasets/longnow-podcasts.git
```

Install the main superdataset from datasets.datalad.org:

```
% datalad clone ///
```

Install a dataset identified by its ID from store.datalad.org:

```
% datalad clone ria+http://store.datalad.org#76b6ca66-36b1-11ea-a2e6-f0d5bf7b5561
```

## Options

### SOURCE

URL, DataLad resource identifier, local path or instance of dataset to be cloned. Constraints: value must be a string

### PATH

path to clone into. If no PATH is provided a destination path will be derived from a source URL similar to git clone.

### -h, -help, -help-np

show this help message. -help-np forcefully disables the use of a pager for displaying the help message



**-d DATASET, --dataset DATASET**

(parent) dataset to clone into. If given, the newly cloned dataset is registered as a subdataset of the parent. Also, if given, relative paths are interpreted as being relative to the parent dataset, and not relative to the working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

**-D DESCRIPTION, --description DESCRIPTION**

short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. Constraints: value must be a string

**--reckless [{auto|ephemeral}]**

Set up the dataset to be able to obtain content in the cheapest/fastest possible way, even if this poses a potential risk the data integrity (‘auto’: hardlink files from a local clone of the dataset, ‘ephemeral’: symlink annex to origin’s annex and discard local availability info via git-annex-dead ‘here’. Please note, that with a symlinked annex you share the annex with origin w/o git-annex knowing this. In case of a change in origin you need to update the clone before you’re able to save new content on your end.). Use with care, and limit to “read-only” use cases. With this flag the installed dataset will be marked as untrusted. The reckless mode is stored in a dataset’s local configuration under ‘datalad.clone.reckless’, and will be inherited to any of its subdatasets. Constraints: value must be one of (None, True, False, ‘auto’, ‘ephemeral’)

**Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

**datalad create-test-dataset****Synopsis**

```
datalad create-test-dataset [-h] [--spec SPEC] [--seed SEED] path
```

**Description**

Create test (meta-)dataset.

**Options****path**

path/name where to create (if specified, must not exist). Constraints: value must be a string

**-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

### **-spec** *SPEC*

spec for hierarchy, defined as a min-max (min could be omitted to assume 0) defining how many (random number from min to max) of sub-datasets to generate at any given level of the hierarchy. Each level separated from each other with /. Example: 1-3/-2 would generate from 1 to 3 subdatasets at the top level, and up to two within those at the 2nd level. Constraints: value must be a string

### **-seed** *SEED*

seed for rng. Constraints: value must be convertible to type 'int'

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad diff

### Synopsis

```
datalad diff [-h] [-f REVISION] [-t REVISION] [-d DATASET] [--annex [MODE]] [--  
↪untracked MODE] [-r] [-R LEVELS] [PATH [PATH ...]]
```

### Description

Report differences between two states of a dataset (hierarchy)

The two to-be-compared states are given via the `-from` and `-to` options. These state identifiers are evaluated in the context of the (specified or detected) dataset. In the case of a recursive report on a dataset hierarchy, corresponding state pairs for any subdataset are determined from the subdataset record in the respective superdataset. Only changes recorded in a subdataset between these two states are reported, and so on.

Any paths given as additional arguments will be used to constrain the difference report. As with Git's diff, it will not result in an error when a path is specified that does not exist on the filesystem.

Reports are very similar to those of the `STATUS` command, with the distinguished content types and states being identical.

#### *Examples*

Show unsaved changes in a dataset:

```
% datalad diff
```

Compare a previous dataset state identified by shasum against current worktree:

```
% datalad diff --from <SHASUM>
```

Compare two branches against each other:

```
% datalad diff --from branch1 --to branch2
```

Show unsaved changes in the dataset and potential subdatasets:

```
% datalad diff --recursive
```

Show unsaved changes made to a particular file:

```
% datalad diff <path/to/file>
```

## Options

### PATH

path to constrain the report to. Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-f REVISION, --from REVISION**

original state to compare to, as given by any identifier that Git understands. Constraints: value must be a string [Default: 'HEAD']

### **-t REVISION, --to REVISION**

state to compare against the original state, as given by any identifier that Git understands. If none is specified, the state of the working tree will be compared. Constraints: value must be a string

### **-d DATASET, --dataset DATASET**

specify the dataset to query. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **--annex [MODE]**

Switch whether to include information on the annex content of individual files in the status report, such as recorded file size. By default no annex information is reported (faster). Three report modes are available: basic information like file size and key name ('basic'); additionally test whether file content is present in the local annex ('availability'; requires one or two additional file system stat calls, but does not call `git-annex`), this will add the result properties `'has_content'` (boolean flag) and `'objloc'` (absolute path to an existing annex object file); or 'all' which will report all available information (presently identical to 'availability'). The 'basic' mode will be assumed when this option is given, but no mode is specified. Constraints: value must be one of (None, 'basic', 'availability', 'all')

### **--untracked MODE**

If and how untracked content is reported when comparing a revision to the state of the working tree. 'no': no untracked content is reported; 'normal': untracked files and entire untracked directories are reported as such; 'all': report

individual files even in fully untracked directories. Constraints: value must be one of ('no', 'normal', 'all') [Default: 'normal']

### **-r, --recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

## **Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## **datalad sshrun**

### **Synopsis**

```
datalad sshrun [-h] [-p PORT] [-4] [-6] [-o OPTION] [-n] login cmd
```

### **Description**

Run command on remote machines via SSH.

This is a replacement for a small part of the functionality of SSH. In addition to SSH alone, this command can make use of datalad's SSH connection management. Its primary use case is to be used with Git as 'core.sshCommand' or via "GIT\_SSH\_COMMAND".

Configure DATALAD.SSH.IDENTITYFILE to pass a file to the ssh's -i option.

### **Options**

#### **login**

[user@]hostname.

#### **cmd**

command for remote execution.

### **-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

**-p *PORT*, --port *PORT***

port to connect to on the remote host.

**-4**

use IPv4 addresses only.

**-6**

use IPv6 addresses only.

**-o *OPTION***

configuration option passed to SSH.

**-n**

Do not connect stdin to the process.

**Authors**

datalad is developed by The DataLad Team and Contributors <[team@datalad.org](mailto:team@datalad.org)>.

**datalad siblings****Synopsis**

```
datalad siblings [-h] [-d DATASET] [-s NAME] [--url [URL]] [--pushurl PUSHURL] [-D_
↪DESCRIPTION] [--fetch] [--as-common-datasrc NAME] [--publish-depends SIBLINGNAME] [-
↪publish-by-default REFSPEC] [--annex-wanted EXPR] [--annex-required EXPR] [--annex-
↪group EXPR] [--annex-groupwanted EXPR] [--inherit] [--no-annex-info] [-r] [-R_
↪LEVELS] [ACTION]
```

**Description****Manage sibling configuration**

This command offers four different actions: ‘query’, ‘add’, ‘remove’, ‘configure’, ‘enable’. ‘query’ is the default action and can be used to obtain information about (all) known siblings. ‘add’ and ‘configure’ are highly similar actions, the only difference being that adding a sibling with a name that is already registered will fail, whereas re-configuring a (different) sibling under a known name will not be considered an error. ‘enable’ can be used to complete access configuration for non-Git sibling (aka git-annex special remotes). Lastly, the ‘remove’ action allows for the removal (or de-configuration) of a registered sibling.

For each sibling (added, configured, or queried) all known sibling properties are reported. This includes:

“**name**” Name of the sibling

“**path**” Absolute path of the dataset

“**url**” For regular siblings at minimum a “fetch” URL, possibly also a “pushurl”

Additionally, any further configuration will also be reported using a key that matches that in the Git configuration.

By default, sibling information is rendered as one line per sibling following this scheme:

```
<dataset_path>: <sibling_name>(<+|->) [<access_specification>]
```

where the + and - labels indicate the presence or absence of a remote data annex at a particular remote, and ACCESS\_SPECIFICATION contains either a URL and/or a type label for the sibling.

## Options

### ACTION

command action selection (see general documentation). Constraints: value must be one of ('query', 'add', 'remove', 'configure', 'enable') [Default: 'query']

### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

### **-d DATASET, --dataset DATASET**

specify the dataset to configure. If no dataset is given, an attempt is made to identify the dataset based on the input and/or the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **-s NAME, --name NAME**

name of the sibling. For addition with path “URLs” and sibling removal this option is mandatory, otherwise the hostname part of a given URL is used as a default. This option can be used to limit ‘query’ to a specific sibling. Constraints: value must be a string

### **--url [URL]**

the URL of or path to the dataset sibling named by NAME. For recursive operation it is required that a template string for building subdataset sibling URLs is given. List of currently available placeholders: %NAME the name of the dataset, where slashes are replaced by dashes. Constraints: value must be a string

### **--pushurl PUSHURL**

in case the URL cannot be used to publish to the dataset sibling, this option specifies a URL to be used instead. If no URL is given, PUSHURL serves as URL as well. Constraints: value must be a string

**-D DESCRIPTION, -description DESCRIPTION**

short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. Constraints: value must be a string

**-fetch**

fetch the sibling after configuration.

**-as-common-datasrc NAME**

configure the created sibling as a common data source of the dataset that can be automatically used by all consumers of the dataset (technical: git-annex auto- enabled special remote).

**-publish-depends SIBLINGNAME**

add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item ‘remote.SIBLINGNAME.datalad-publish-depends’. This option can be given more than once to configure multiple dependencies. Constraints: value must be a string

**-publish-by-default REFSPEC**

add a refsPEC to be published to this sibling by default if nothing specified. Constraints: value must be a string

**-annex-wanted EXPR**

expression to specify ‘wanted’ content for the repository/sibling. See <https://git-annex.branchable.com/git-annex-wanted/> for more information. Constraints: value must be a string

**-annex-required EXPR**

expression to specify ‘required’ content for the repository/sibling. See <https://git-annex.branchable.com/git-annex-required/> for more information. Constraints: value must be a string

**-annex-group EXPR**

expression to specify a group for the repository. See <https://git-annex.branchable.com/git-annex-group/> for more information. Constraints: value must be a string

**-annex-groupwanted EXPR**

expression for the groupwanted. Makes sense only if `-annex-wanted="groupwanted"` and `annex-group` is given too. See <https://git-annex.branchable.com/git-annex-groupwanted/> for more information. Constraints: value must be a string

### **-inherit**

if sibling is missing, inherit settings (git config, git annex wanted/group/groupwanted) from its super-dataset.

### **-no-annex-info**

Whether to query all information about the annex configurations of siblings. Can be disabled if speed is a concern.

### **-r, -recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, -recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

## **Authors**

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## **datalad status**

### **Synopsis**

```
datalad status [-h] [-d DATASET] [--annex [MODE]] [--untracked MODE] [-r] [-R LEVELS]   
↳ [-e {no|commit|full}] [-t {raw|eval}] [PATH [PATH ...]]
```

### **Description**

Report on the state of dataset content.

This is an analog to *git status* that is simultaneously crippled and more powerful. It is crippled, because it only supports a fraction of the functionality of its counter part and only distinguishes a subset of the states that Git knows about. But it is also more powerful as it can handle status reports for a whole hierarchy of datasets, with the ability to report on a subset of the content (selection of paths) across any number of datasets in the hierarchy.

#### *Path conventions*

All reports are guaranteed to use absolute paths that are underneath the given or detected reference dataset, regardless of whether query paths are given as absolute or relative paths (with respect to the working directory, or to the reference dataset, when such a dataset is given explicitly). Moreover, so-called “explicit relative paths” (i.e. paths that start with ‘.’ or ‘..’) are also supported, and are interpreted as relative paths with respect to the current working directory regardless of whether a reference dataset with specified.

When it is necessary to address a subdataset record in a superdataset without causing a status query for the state `_within_` the subdataset itself, this can be achieved by explicitly providing a reference dataset and the path to the root of the subdataset like so:



```
datalad status --dataset . subdspath
```

In contrast, when the state of the subdataset within the superdataset is not relevant, a status query for the content of the subdataset can be obtained by adding a trailing path separator to the query path (rsync-like syntax):

```
datalad status --dataset . subdspath/
```

When both aspects are relevant (the state of the subdataset content and the state of the subdataset within the superdataset), both queries can be combined:

```
datalad status --dataset . subdspath subdspath/
```

When performing a recursive status query, both status aspects of subdataset are always included in the report.

### *Content types*

The following content types are distinguished:

- ‘dataset’ – any top-level dataset, or any subdataset that is properly registered in superdataset
- ‘directory’ – any directory that does not qualify for type ‘dataset’
- ‘file’ – any file, or any symlink that is placeholder to an annexed file
- ‘symlink’ – any symlink that is not used as a placeholder for an annexed file

### *Content states*

The following content states are distinguished:

- ‘clean’
- ‘added’
- ‘modified’
- ‘deleted’
- ‘untracked’

### *Examples*

Report on the state of a dataset:

```
% datalad status
```

Report on the state of a dataset and all subdatasets:

```
% datalad status --recursive
```

Address a subdataset record in a superdataset without causing a status query for the state `_within_` the subdataset itself:

```
% datalad status --dataset . mysubdataset
```

Get a status query for the state within the subdataset without causing a status query for the superdataset (using trailing path separator in the query path)::

```
% datalad status --dataset . mysubdataset/
```

Report on the state of a subdataset in a superdataset and on the state within the subdataset:

```
% datalad status --dataset . mysubdataset mysubdataset/
```

Report the file size of annexed content in a dataset:

```
% datalad status --annex
```

### Options

#### PATH

path to be evaluated. Constraints: value must be a string

#### **-h, --help, --help-np**

show this help message. `--help-np` forcefully disables the use of a pager for displaying the help message

#### **-d DATASET, --dataset DATASET**

specify the dataset to query. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

#### **--annex [MODE]**

Switch whether to include information on the annex content of individual files in the status report, such as recorded file size. By default no annex information is reported (faster). Three report modes are available: basic information like file size and key name ('basic'); additionally test whether file content is present in the local annex ('availability'; requires one or two additional file system stat calls, but does not call `git-annex`), this will add the result properties 'has\_content' (boolean flag) and 'objloc' (absolute path to an existing annex object file); or 'all' which will report all available information (presently identical to 'availability'). The 'basic' mode will be assumed when this option is given, but no mode is specified. Constraints: value must be one of (None, 'basic', 'availability', 'all')

#### **--untracked MODE**

If and how untracked content is reported when comparing a revision to the state of the working tree. 'no': no untracked content is reported; 'normal': untracked files and entire untracked directories are reported as such; 'all': report individual files even in fully untracked directories. Constraints: value must be one of ('no', 'normal', 'all') [Default: 'normal']

#### **-r, --recursive**

if set, recurse into potential subdataset.

#### **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

**-e {no|commit|full}, --eval-subdataset-state {no|commit|full}**

Evaluation of subdataset state (clean vs. modified) can be expensive for deep dataset hierarchies as subdataset have to be tested recursively for uncommitted modifications. Setting this option to ‘no’ or ‘commit’ can substantially boost performance by limiting what is being tested. With ‘no’ no state is evaluated and subdataset result records typically do not contain a ‘state’ property. With ‘commit’ only a discrepancy of the HEAD commit shasum of a subdataset and the shasum recorded in the superdataset’s record is evaluated, and the ‘state’ result property only reflects this aspect. With ‘full’ any other modification is considered too (see the ‘untracked’ option for further tailoring modification testing). Constraints: value must be one of (‘no’, ‘commit’, ‘full’) [Default: ‘full’]

**-t {raw|eval}, --report-filetype {raw|eval}**

Report mode for file types. With ‘eval’ each symlink is inspected whether it is a pointer to an annex’ed file, and is reported as ‘type=file’ in this case, and ‘type=symlink’ otherwise. With ‘raw’ no type inspection is performed, and symlinks representing annex’ed files are indistinguishable from other symlinks. Type inspection is relatively expensive and can lead to slow operation in datasets with a large number of files. Constraints: value must be one of (‘raw’, ‘eval’) [Default: ‘eval’]

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## datalad subdatasets

### Synopsis

```
datalad subdatasets [-h] [-d DATASET] [--fulfilled FULFILLED] [-r] [-R LEVELS] [--contains PATH] [--bottomup] [--set-property NAME VALUE] [--delete-property NAME]
↳ [PATH [PATH ...]]
```

### Description

Report subdatasets and their properties.

The following properties are reported (if possible) for each matching subdataset record.

“**name**” Name of the subdataset in the parent (often identical with the relative path in the parent dataset)

“**path**” Absolute path to the subdataset

“**parentds**” Absolute path to the parent dataset

“**gitshasum**” SHA1 of the subdataset commit recorded in the parent dataset

“**state**” Condition of the subdataset: ‘clean’, ‘modified’, ‘absent’, ‘conflict’ as reported by *git submodule*

“**gitmodule\_url**” URL of the subdataset recorded in the parent

“**gitmodule\_name**” Name of the subdataset recorded in the parent

“**gitmodule\_<label>**” Any additional configuration property on record.

Performance note: Property modification, requesting BOTTOMUP reporting order, or a particular numerical RECURSION\_LIMIT implies an internal switch to an alternative query implementation for recursive query that is more flexible, but also notably slower (performs one call to Git per dataset versus a single call for all combined).

The following properties for subdatasets are recognized by DataLad (without the 'gitmodule\_' prefix that is used in the query results):

**“datalad-recursiveinstall”** If set to 'skip', the respective subdataset is skipped when DataLad is recursively installing its superdataset. However, the subdataset remains installable when explicitly requested, and no other features are impaired.

## Options

### PATH

path/name to query for subdatasets. Defaults to the current directory. Constraints: value must be a string

### **-h, --help, --help-np**

show this help message. --help-np forcefully disables the use of a pager for displaying the help message

### **-d DATASET, --dataset DATASET**

specify the dataset to query. If no dataset is given, an attempt is made to identify the dataset based on the input and/or the current working directory. Constraints: Value must be a Dataset or a valid identifier of a Dataset (e.g. a path)

### **--fulfilled FULFILLED**

if given, must be a boolean flag indicating whether to report either only locally present or absent datasets. By default subdatasets are reported regardless of their status. Constraints: value must be convertible to type bool

### **-r, --recursive**

if set, recurse into potential subdataset.

### **-R LEVELS, --recursion-limit LEVELS**

limit recursion into subdataset to the given number of levels. Constraints: value must be convertible to type 'int'

### **--contains PATH**

limit report to the subdatasets containing the given path. If a root path of a subdataset is given the last reported dataset will be the subdataset itself. This option can be given multiple times, in which case datasets will be reported that contain any of the given paths. Constraints: value must be a string

### –bottomup

whether to report subdatasets in bottom-up order along each branch in the dataset tree, and not top-down.

### –set-property NAME VALUE

Name and value of one or more subdataset properties to be set in the parent dataset’s .gitmodules file. The property name is case-insensitive, must start with a letter, and consist only of alphanumeric characters. The value can be a Python format() template string wrapped in ‘<>’ (e.g. ‘<{gitmodule\_name}>’). Supported keywords are any item reported in the result properties of this command, plus ‘refds\_relpath’ and ‘refds\_relname’: the relative path of a subdataset with respect to the base dataset of the command call, and, in the latter case, the same string with all directory separators replaced by dashes. This option can be given multiple times. Constraints: value must be a string

### –delete-property NAME

Name of one or more subdataset properties to be removed from the parent dataset’s .gitmodules file. This option can be given multiple times. Constraints: value must be a string

## Authors

datalad is developed by The DataLad Team and Contributors <team@datalad.org>.

## 1.5.2 Python module reference

This module reference extends the manual with a comprehensive overview of the available functionality built into datalad. Each module in the package is documented by a general summary of its purpose and the list of classes and functions it provides.

### High-level user interface

#### Dataset operations

|  |  |
|--|--|
| <i>api.Dataset</i> (path)                                | Representation of a DataLad dataset/repository                 |
| <i>api.create</i> ([path, initopts, force, ...])         | Create a new dataset from scratch.                             |
| <i>api.create_sibling</i> (sshurl[, name, ...])          | Create a dataset sibling on a UNIX-like SSH-accessible machine |
| <i>api.create_sibling_github</i> (reponame[, ...])       | Create dataset sibling on Github.                              |
| <i>api.create_sibling_gitlab</i> ([path, site, ...])     | Create dataset sibling at a GitLab site                        |
| <i>api.drop</i> ([path, dataset, recursive, ...])        | Drop file content from datasets                                |
| <i>api.get</i> ([path, source, dataset, recursive, ...]) | Get any dataset content (files/directories/subdatasets).       |
| <i>api.install</i> ([path, source, dataset, ...])        | Install a dataset from a (remote) source.                      |
| <i>api.publish</i> ([path, dataset, to, since, ...])     | Publish a dataset to a known sibling.                          |
| <i>api.remove</i> ([path, dataset, recursive, ...])      | Remove components from datasets                                |
| <i>api.save</i> ([path, message, dataset, ...])          | Save the current state of a dataset                            |
| <i>api.update</i> ([path, sibling, merge, follow, ...])  | Update a dataset from a sibling.                               |
| <i>api.uninstall</i> ([path, dataset, recursive, ...])   | Uninstall subdatasets  |
| <i>api.unlock</i> ([path, dataset, recursive, ...])      | Unlock file(s) of a dataset                                    |

## datalad.api.Dataset

**class** `datalad.api.Dataset` (*path*)

Representation of a DataLad dataset/repository

This is the core data type of DataLad: a representation of a dataset. At its core, datasets are (git-annex enabled) Git repositories. This class provides all operations that can be performed on a dataset.

Creating a dataset instance is cheap, all actual operations are delayed until they are actually needed. Creating multiple *Dataset* class instances for the same Dataset location will automatically yield references to the same object.

A dataset instance comprises of two major components: a *repo* attribute, and a *config* attribute. The former offers access to low-level functionality of the Git or git-annex repository. The latter gives access to a dataset's configuration manager.

Most functionality is available via methods of this class, but also as stand-alone functions with the same name in *datalad.api*.

`__init__` (*path*)

**Parameters** *path* (*str* or *Path*) – Path to the dataset location. This location may or may not exist yet.

### Methods

---

|                                       |   |
|---------------------------------------|---|
| <code>__init__</code> ( <i>path</i> ) | <b>param path</b> Path to the dataset location.<br>This location may or may not exist |
|---------------------------------------|---|

---

|  |  |
|--|--|
| <code>Dataset.add</code>   |  |
| <code>add_readme</code> ( <i>filename</i> , <i>existing</i> )                      | Add basic information about DataLad datasets to a README file                |
| <code>addurls</code> ( <i>urlfile</i> , <i>urlformat</i> , <i>filenameformat</i> ) | Create and update a dataset from a list of URLs.                             |
| <code>aggregate_metadata</code> ( <i>[dataset, recursive, ...]</i> )               | Aggregate metadata of one or more datasets for later query.                  |
| <code>annotate_paths</code> ( <i>[dataset, recursive, ...]</i> )                   | Analyze and act upon input paths   |
| <code>clean</code> ( <i>[what, recursive, recursion_limit]</i> )                   | Clean up after DataLad (possible temporary files etc.)                       |
| <code>clone</code> ( <i>[path, dataset, description, reckless]</i> )               | Obtain a dataset (copy) from a URL or local directory                        |
| <code>close</code> ()  | Perform operations which would close any possible process using this Dataset |
| <code>create</code> ( <i>[initopts, force, description, ...]</i> )                 | Create a new dataset from scratch.   |
| <code>create_sibling</code> ( <i>[name, target_dir, ...]</i> )                     | Create a dataset sibling on a UNIX-like SSH-accessible machine               |
| <code>create_sibling_github</code> ( <i>[dataset, recursive, ...]</i> )            | Create dataset sibling on Github.  |
| <code>create_sibling_gitlab</code> ( <i>[site, project, ...]</i> )                 | Create dataset sibling at a GitLab site                                      |
| <code>diff</code> ( <i>[fr, to, dataset, annex, untracked, ...]</i> )              | Report differences between two states of a dataset (hierarchy)               |
| <code>download_url</code> ( <i>[dataset, path, overwrite, ...]</i> )               | Download content   |
| <code>drop</code> ( <i>[dataset, recursive, recursion_limit, ...]</i> )            | Drop file content from datasets  |
| <code>export_archive</code> ( <i>[filename, archivetype, ...]</i> )                | Export the content of a dataset as a TAR/ZIP archive.                        |

Continued on next page

Table 3 – continued from previous page

|  |   |
|--|---|
| <code>export_to_figshare([filename, ...])</code>               | Export the content of a dataset as a ZIP archive to figshare  |
| <code>extract_metadata([files, dataset])</code>                | Run one or more of DataLad’s metadata extractors on a dataset or file.  |
| <code>get([source, dataset, recursive, ...])</code>            | Get any dataset content (files/directories/subdatasets).  |
| <code>get_superdataset([datalad_only, topmost, ...])</code>    | Get the dataset’s superdataset  |
| <code>install([source, dataset, get_data, ...])</code>         | Install a dataset from a (remote) source.   |
| <code>is_installed()</code>                                    | Returns whether a dataset is installed.   |
| <code>metadata([dataset, get_aggregates, ...])</code>          | Metadata reporting for files and entire datasets  |
| <code>no_annex(pattern[, ref_dir, makedirs])</code>            | Configure a dataset to never put some content into the dataset’s annex  |
| <code>publish([dataset, to, since, missing, ...])</code>       | Publish a dataset to a known sibling.   |
| <code>recall_state(whereto)</code>                             | Something that can be used to checkout a particular state (tag, commit) to “undo” a change or switch to a otherwise desired previous state. |
| <code>remove([dataset, recursive, check, save, ...])</code>    | Remove components from datasets   |
| <code>rerun([since, dataset, branch, message, ...])</code>     | Re-execute previous <i>datalad run</i> commands.  |
| <code>run([dataset, inputs, outputs, expand, ...])</code>      | Run an arbitrary shell command and record its impact on a dataset.  |
| <code>run_procedure([dataset, discover, help_proc])</code>     | Run prepared procedures (DataLad scripts) on a dataset  |
| <code>save([message, dataset, version_tag, ...])</code>        | Save the current state of a dataset   |
| <code>search([dataset, force_reindex, ...])</code>             | Search dataset metadata   |
| <code>siblings([dataset, name, url, pushurl, ...])</code>      | Manage sibling configuration  |
| <code>status([dataset, annex, untracked, ...])</code>          | Report on the state of dataset content.   |
| <code>subdatasets([dataset, fulfilled, recursive, ...])</code> | Report subdatasets and their properties.  |
| <code>uninstall([dataset, recursive, check, if_dirty])</code>  | Uninstall subdatasets   |
| <code>unlock([dataset, recursive, recursion_limit])</code>     | Unlock file(s) of a dataset   |
| <code>update([sibling, merge, follow, dataset, ...])</code>    | Update a dataset from a sibling.  |
| <code>wtf([sensitive, sections, decor, clipboard])</code>      | Generate a report about the DataLad installation and configuration  |

### Attributes

|                      |  |
|----------------------|--|
| <code>config</code>  | Get an instance of the parser for the persistent dataset configuration.  |
| <code>id</code>      | Identifier of the dataset.   |
| <code>path</code>    | path to the dataset  |
| <code>pathobj</code> | pathobj for the dataset  |
| <code>repo</code>    | Get an instance of the version control system/repo for this dataset, or None if there is none yet (or none anymore). |

### datalad.api.create

`datalad.api.create` (*path=None, intopts=None, force=False, description=None, dataset=None, no\_annex=False, fake\_dates=False, cfg\_proc=None*)  
 Create a new dataset from scratch.

This command initializes a new dataset at a given location, or the current directory. The new dataset can optionally be registered in an existing superdataset (the new dataset's path needs to be located within the superdataset for that, and the superdataset needs to be given explicitly via *dataset*). It is recommended to provide a brief description to label the dataset's nature *and* location, e.g. "Michael's music on black laptop". This helps humans to identify data locations in distributed scenarios. By default an identifier comprised of user and machine name, plus path will be generated.

This command only creates a new dataset, it does not add existing content to it, even if the target directory already contains additional files or directories.

Plain Git repositories can be created via the *no\_annex* flag. However, the result will not be a full dataset, and, consequently, not all features are supported (e.g. a description).

To create a local version of a remote dataset use the `~datalad.api.install` command instead.

---

**Note:** Power-user info: This command uses `git init` and `git annex init` to prepare the new dataset. Registering to a superdataset is performed via a git submodule add operation in the discovered superdataset.

---

## Examples

Create a dataset 'mydataset' in the current directory:

```
> create(path='mydataset')
```

Apply the `text2git` procedure upon creation of a dataset:

```
> create(path='mydataset', cfg_proc='text2git')
```

Create a subdataset in the root of an existing dataset:

```
> create(dataset='.', path='mysubdataset')
```

Create a dataset in an existing, non-empty directory:

```
> create(force=True)
```

Create a plain Git repository:

```
> create(path='mydataset', no_annex=True)
```

## Parameters

- **path** (*str or Dataset or None, optional*) – path where the dataset shall be created, directories will be created as necessary. If no location is provided, a dataset will be created in the current working directory. Either way the command will error if the target directory is not empty. Use *force* to create a dataset in a non-empty directory. [Default: None]
- **initopts** – options to pass to `git init`. Options can be given as a list of command line arguments or as a GitPython-style option dictionary. Note that not all options will lead to viable results. For example '`-bare`' will not yield a repository where DataLad can adjust files in its working tree. [Default: None]
- **force** (*bool, optional*) – enforce creation of a dataset in a non-empty directory. [Default: False]



- **description** (*str or None, optional*) – short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the create operation on. If a dataset is given, a new subdataset will be created in it. [Default: None]
- **no\_annex** (*bool, optional*) – if set, a plain Git repository will be created without any annex. [Default: False]
- **fake\_dates** (*bool, optional*) – Configure the repository to use fake dates. The date for a new commit will be set to one second later than the latest commit in the repository. This can be used to anonymize dates. [Default: False]
- **cfg\_proc** – Run `cfg_PROC` procedure(s) (can be specified multiple times) on the created dataset. Use `run_procedure(discover=True)` to get a list of available procedures, such as `cfg_text2git`. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## datalad.api.create\_sibling

```
datalad.api.create_sibling(sshurl, name=None, target_dir=None, target_url=None, target_pushurl=None, dataset=None, recursive=False, recursion_limit=None, existing='error', shared=None, group=None, ui=False, as_common_datasrc=None, publish_by_default=None, publish_depends=None, annex_wanted=None, annex_group=None, annex_groupwanted=None, inherit=False, since=None)
```

Create a dataset sibling on a UNIX-like SSH-accessible machine

Given a local dataset, and SSH login information this command creates a remote dataset repository and configures it as a dataset sibling to be used as a publication target (see *publish* command).

Various properties of the remote sibling can be configured (e.g. name location on the server, read and write access URLs, and access permissions).

Optionally, a basic web-viewer for DataLad datasets can be installed at the remote location.

This command supports recursive processing of dataset hierarchies, creating a remote sibling for each dataset in the hierarchy. By default, remote siblings are created in hierarchical structure that reflects the organization on the local file system. However, a simple templating mechanism is provided to produce a flat list of datasets (see *-target-dir*).

### Parameters

- **sshurl** (*str*) – Login information for the target server. This can be given as a URL (*ssh://host/path*) or SSH-style (*user@host:path*). Unless overridden, this also serves the future dataset’s access URL and path on the server.
- **name** (*str or None, optional*) – sibling name to create for this publication target. If *recursive* is set, the same name will be used to label all the subdatasets’ siblings. When creating a target dataset fails, no sibling is added. [Default: None]
- **target\_dir** (*str or None, optional*) – path to the directory *on the server* where the dataset shall be created. By default the SSH access URL is used to identify this directory. If a relative path is provided here, it is interpreted as being relative to the user’s home directory on the server. Additional features are relevant for recursive processing of datasets with subdatasets. By default, the local dataset structure is replicated on the server. However, it is possible to provide a template for generating different target directory names for all (sub)datasets. Templates can contain certain placeholder that are substituted for each (sub)dataset. For example: *“/mydirectory/dataset%%RELNAME”*. Supported placeholders: *%%RELNAME* - the name of the datasets, with any slashes replaced by dashes. [Default: None]
- **target\_url** (*str or None, optional*) – “public” access URL of the to-be-created target dataset(s) (default: *sshurl*). Accessibility of this URL determines the access permissions of potential consumers of the dataset. As with *target\_dir*, templates (same set of placeholders) are supported. Also, if specified, it is provided as the annex description. [Default: None]
- **target\_pushurl** (*str or None, optional*) – In case the *target\_url* cannot be used to publish to the dataset, this option specifies an alternative URL for this purpose. As with *target\_url*, templates (same set of placeholders) are supported. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to create the publication target for. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]

- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **existing** (*{'skip', 'error', 'reconfigure', 'replace'}, optional*) – action to perform, if a sibling is already configured under the given name and/or a target (non-empty) directory already exists. In this case, a dataset can be skipped ('skip'), the sibling configuration be updated ('reconfigure'), or process interrupts with error ('error'). DANGER ZONE: If 'replace' is used, an existing target directory will be forcefully removed, re-initialized, and the sibling (re-)configured (thus implies 'reconfigure'). *replace* could lead to data loss, so use with care. To minimize possibility of data loss, in interactive mode DataLad will ask for confirmation, but it would just issue a warning and proceed in non-interactive mode. [Default: 'error']
- **shared** (*str or bool or None, optional*) – if given, configures the access permissions on the server for multi- users (this could include access by a webserver!). Possible values for this option are identical to those of *git init -shared* and are described in its documentation. [Default: None]
- **group** (*str or None, optional*) – Filesystem group for the repository. Specifying the group is particularly important when *shared="group"*. [Default: None]
- **ui** (*bool or str, optional*) – publish a web interface for the dataset with an optional user- specified name for the html at publication target. defaults to *index.html* at dataset root. [Default: False]
- **as\_common\_datasrc** – configure the created sibling as a common data source of the dataset that can be automatically used by all consumers of the dataset (technical: git-annex auto-enabled special remote). [Default: None]
- **publish\_by\_default** (*list of str or None, optional*) – add a refspec to be published to this sibling by default if nothing specified. [Default: None]
- **publish\_depends** (*list of str or None, optional*) – add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item 'remote.SIBLINGNAME.datalad-publish-depends'. Multiple dependencies can be given as a list of sibling names. [Default: None]
- **annex\_wanted** (*str or None, optional*) – expression to specify 'wanted' content for the repository/sibling. See <https://git-annex.branchable.com/git-annex-wanted/> for more information. [Default: None]
- **annex\_group** (*str or None, optional*) – expression to specify a group for the repository. See <https://git-annex.branchable.com/git-annex-group/> for more information. [Default: None]
- **annex\_groupwanted** (*str or None, optional*) – expression for the groupwanted. Makes sense only if *annex\_wanted="groupwanted"* and *annex-group* is given too. See <https://git-annex.branchable.com/git-annex-groupwanted/> for more information. [Default: None]
- **inherit** (*bool, optional*) – if sibling is missing, inherit settings (git config, git annex wanted/group/groupwanted) from its super-dataset. [Default: False]
- **since** (*str or None, optional*) – limit processing to datasets that have been changed since a given state (by tag, branch, commit, etc). This can be used to create siblings for recently added subdatasets. [Default: None]

## datalad.api.create\_sibling\_github

```
datalad.api.create_sibling_github(reponame, dataset=None, recursive=False, recursion_limit=None, name='github', existing='error', github_login=None, github_passwd=None, github_organization=None, access_protocol='https', publish_depends=None, dryrun=False)
```

Create dataset sibling on Github.

An existing GitHub project, or a project created via the GitHub website can be configured as a sibling with the siblings command. Alternatively, this command can create a repository under a user's Github account, or any organization a user is a member of (given appropriate permissions). This is particularly helpful for recursive sibling creation for subdatasets. In such a case, a dataset hierarchy is represented as a flat list of GitHub repositories.

Github cannot host dataset content. However, in combination with other data sources (and siblings), publishing a dataset to Github can facilitate distribution and exchange, while still allowing any dataset consumer to obtain actual data content from alternative sources.

For Github authentication user credentials can be given as arguments. Alternatively, they are obtained interactively or queried from the systems credential store. Lastly, an *oauth* token stored in the Git configuration under variable *hub.oauthtoken* will be used automatically. Such a token can be obtained, for example, using the commandline Github interface (<https://github.com/sociomantic/git-hub>) by running: `git hub setup` (if no 2FA is used).

### Parameters

- **reponame** (*str*) – Github repository name. When operating recursively, a suffix will be appended to this name for each subdataset.
- **dataset** (*Dataset or None, optional*) – specify the dataset to create the publication target for. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **name** (*str, optional*) – name to represent the Github repository in the local dataset installation. [Default: 'github']
- **existing** (*{'skip', 'error', 'reconfigure'}, optional*) – desired behavior when already existing or configured siblings are discovered. 'skip': ignore; 'error': fail immediately; 'reconfigure': use the existing repository and reconfigure the local dataset to use it as a sibling. [Default: 'error']
- **github\_login** (*str or None, optional*) – Github user name or access token. [Default: None]
- **github\_passwd** (*str or None, optional*) – Github user password. [Default: None]
- **github\_organization** (*str or None, optional*) – If provided, the repository will be created under this Github organization. The respective Github user needs appropriate permissions. [Default: None]
- **access\_protocol** (*{'https', 'ssh'}, optional*) – Which access protocol/URL to configure for the sibling. [Default: 'https']

- **publish\_depends** (*list of str or None, optional*) – add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item ‘remote.SIBLINGNAME.datalad-publish-depends’. Multiple dependencies can be given as a list of sibling names. [Default: None]
- **dryrun** (*bool, optional*) – If this flag is set, no communication with Github is performed, and no repositories will be created. Instead would-be repository names are reported for all relevant datasets. [Default: False]

### datalad.api.create\_sibling\_gitlab

`datalad.api.create_sibling_gitlab` (*path=None, site=None, project=None, layout=None, dataset=None, recursive=False, recursion\_limit=None, name=None, existing='error', access=None, publish\_depends=None, description=None, dryrun=False*)

Create dataset sibling at a GitLab site

An existing GitLab project, or a project created via the GitLab web interface can be configured as a sibling with the `siblings` command. Alternatively, this command can create a GitLab project at any location/path a given user has appropriate permissions for. This is particularly helpful for recursive sibling creation for subdatasets. API access and authentication are implemented via `python-gitlab`, and all its features are supported. A particular GitLab site must be configured in a named section of a `python-gitlab.cfg` file (see <https://python-gitlab.readthedocs.io/en/stable/cli.html#configuration> for details), such as:

```
[mygit]
url = https://git.example.com
api_version = 4
private_token = abcdefghijklmnopqrst
```

Subsequently, this site is identified by its name (‘mygit’ in the example above).

(Recursive) sibling creation for all, or a selected subset of subdatasets is supported with three different project layouts (see `-layout`):

“**hierarchy**” Each dataset is placed into its own group, and the actual GitLab project for a dataset is put in a project named “\_repo\_” inside this group. Using this layout, arbitrarily deep hierarchies of nested datasets can be represented, while the hierarchical structure is reflected in the project path. This is the default layout, if no project path is specified.

“**flat**” All datasets are placed in the same group. The name of a project is its relative path within the root dataset, with all path separator characters replaced by ‘-’.

“**collection**” This is a hybrid layout, where the root dataset is placed in a “\_repo\_” project inside a group, and all nested subdatasets are represented inside the group using a “flat” layout.

GitLab cannot host dataset content. However, in combination with other data sources (and siblings), publishing a dataset to GitLab can facilitate distribution and exchange, while still allowing any dataset consumer to obtain actual data content from alternative sources.

#### Configuration

All configuration switches and options for GitLab sibling creation can be provided arguments to the command. However, it is also possible to specify a particular setup in a dataset’s configuration. This is particularly important when managing large collections of datasets. Configuration options are:

“**datalad.gitlab-default-site**” Name of the default GitLab site (see `-site`)

“**datalad.gitlab-SITENAME-siblingname**” Name of the sibling configured for the local dataset that points to the GitLab instance SITENAME (see `-name`)

“**datalad.gitlab-SITENAME-layout**” Project layout used at the GitLab instance SITENAME (see `-layout`)

“**datalad.gitlab-SITENAME-access**” Access method used for the GitLab instance SITENAME (see `-access`)

“**datalad.gitlab-SITENAME-project**” Project location/path used for a datasets at GitLab instance SITENAME (see `-project`). Configuring this is useful for deriving project paths for subdatasets, relative to superdataset.

### Parameters

- **path** – selectively create siblings for any datasets underneath a given path. By default only the root dataset is considered. [Default: None]
- **site** (*None or str, optional*) – name of the GitLab site to create a sibling at. Must match an existing `python-gitlab` configuration section with location and authentication settings (see <https://python-gitlab.readthedocs.io/en/stable/cli.html#configuration>). By default the dataset configuration is consulted. [Default: None]
- **project** (*None or str, optional*) – project path at the GitLab site. If a subdataset of the reference dataset is processed, its project path is automatically determined by the `layout` configuration, by default. [Default: None]
- **layout** (*{None, 'hierarchy', 'collection', 'flat'}, optional*) – layout of projects at the GitLab site, if a collection, or a hierarchy of datasets and subdatasets is to be created. By default the dataset configuration is consulted. [Default: None]
- **dataset** (*Dataset or None, optional*) – reference or root dataset. If no path constraints are given, a sibling for this dataset will be created. In this and all other cases, the reference dataset is also consulted for the GitLab configuration, and desired project layout. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **name** (*str or None, optional*) – name to represent the GitLab sibling remote in the local dataset installation. If not specified a name is looked up in the dataset configuration, or defaults to the `site` name. [Default: None]
- **existing** (*{'skip', 'error', 'reconfigure'}, optional*) – desired behavior when already existing or configured siblings are discovered. ‘skip’: ignore; ‘error’: fail, if access URLs differ; ‘reconfigure’: use the existing repository and reconfigure the local dataset to use it as a sibling. [Default: ‘error’]
- **access** (*{None, 'http', 'ssh', 'ssh+http'}, optional*) – access method used for data transfer to and from the sibling. ‘ssh’: read and write access used the SSH protocol; ‘http’: read and write access use HTTP requests; ‘ssh+http’: read access is done via HTTP and write access performed with SSH. Dataset configuration is consulted for a default, ‘http’ is used otherwise. [Default: None]
- **publish\_depends** (*list of str or None, optional*) – add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item ‘remote.SIBLINGNAME.datalad-publish-depends’. Multiple dependencies can be given as a list of sibling names. [Default: None]
- **description** (*str or None, optional*) – brief description for the GitLab project (displayed on the site). [Default: None]

- **dryrun** (*bool, optional*) – If this flag is set, no communication with GitLab is performed, and no repositories will be created. Instead would-be repository names and configurations are reported for all relevant datasets. [Default: False]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## datalad.api.drop

`datalad.api.drop` (*path=None, dataset=None, recursive=False, recursion\_limit=None, check=True, if\_dirty='save-before'*)

Drop file content from datasets

This command takes any number of paths of files and/or directories. If a common (super)dataset is given explicitly, the given paths are interpreted relative to this dataset.

Recursion into subdatasets needs to be explicitly enabled, while recursion into subdirectories within a dataset is done automatically. An optional recursion limit is applied relative to each given input path.

By default, the availability of at least one remote copy is verified before file content is dropped. As these checks could lead to slow operation (network latencies, etc), they can be disabled.

### Examples

Drop single file content:

```
> drop('path/to/file')
```

Drop all file content in the current dataset:

```
> drop('.')
```

Drop all file content in a dataset and all its subdatasets:

```
> drop(dataset='.', recursive=True)
```

Disable check to ensure the configured minimum number of remote sources for dropped data:

```
> drop(path='path/to/content', check=False)
```

### Parameters

- **path** (*sequence of str or None, optional*) – path/name of the component to be dropped. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the operation on. If no dataset is given, an attempt is made to identify a dataset based on the *path* given. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **check** (*bool, optional*) – whether to perform checks to assure the configured minimum number (remote) source for data. [Default: True]
- **if\_dirty** – desired behavior if a dataset with unsaved changes is discovered: ‘fail’ will trigger an error and further processing is aborted; ‘save-before’ will save all changes prior any further action; ‘ignore’ let’s datalad proceed as if the dataset would not have unsaved changes. [Default: ‘save-before’]
- **on\_failure** (*{‘ignore’, ‘continue’, ‘stop’}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{‘default’, ‘json’, ‘json\_pp’, ‘tailored’} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{‘datasets’, ‘successdatasets-or-none’, ‘paths’, ‘relpaths’, ‘metadata’} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]



- **return\_type** (*{'generator', 'list', 'item-or-list'}*, *optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.get

`datalad.api.get` (*path=None, source=None, dataset=None, recursive=False, recursion\_limit=None, get\_data=True, description=None, reckless=None, jobs='auto'*)

Get any dataset content (files/directories/subdatasets).

This command only operates on dataset content. To obtain a new independent dataset from some source use the *install* command.

By default this command operates recursively within a dataset, but not across potential subdatasets, i.e. if a directory is provided, all files in the directory are obtained. Recursion into subdatasets is supported too. If enabled, relevant subdatasets are detected and installed in order to fulfill a request.

Known data locations for each requested file are evaluated and data are obtained from some available location (according to git-annex configuration and possibly assigned remote priorities), unless a specific source is specified.

---

**Note:** Power-user info: This command uses git annex get to fulfill file handles.

---

## Examples

Get a single file:

```
> get('path/to/file')
```

Get contents of a directory:

```
> get('path/to/dir/')
```

Get all contents of the current dataset and its subdatasets:

```
> get(dataset='.', recursive=True)
```

Get (clone) a registered subdataset, but don't retrieve data:

```
> get('path/to/subds', get_data=False)
```

## Parameters

- **path** (*sequence of str or None, optional*) – path/name of the requested dataset component. The component must already be known to a dataset. To add new components to a dataset use the *add* command. [Default: *None*]
- **source** (*str or None, optional*) – label of the data source to be used to fulfill requests. This can be the name of a dataset sibling or another known source. [Default: *None*]
- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the add operation on, in which case *path* arguments are interpreted as being relative to this dataset.

If no dataset is given, an attempt is made to identify a dataset for each input *path*. [Default: None]

- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or {'existing'} or None, optional*) – limit recursion into subdataset to the given number of levels. Alternatively, ‘existing’ will limit recursion to subdatasets that already existed on the filesystem at the start of processing, and prevent new subdatasets from being obtained recursively. [Default: None]
- **get\_data** (*bool, optional*) – whether to obtain data for all file handles. If disabled, *get* operations are limited to dataset handles. [Default: True]
- **description** (*str or None, optional*) – short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. [Default: None]
- **reckless** (*{None, True, False, 'auto', 'ephemeral'}, optional*) – Set up the dataset to be able to obtain content in the cheapest/fastest possible way, even if this poses a potential risk the data integrity (‘auto’: hardlink files from a local clone of the dataset, ‘ephemeral’: symlink annex to origin’s annex and discard local availability info via git-annex-dead ‘here’. Please note, that with a symlinked annex you share the annex with origin w/o git-annex knowing this. In case of a change in origin you need to update the clone before you’re able to save new content on your end.). Use with care, and limit to “read-only” use cases. With this flag the installed dataset will be marked as untrusted. The reckless mode is stored in a dataset’s local configuration under ‘datalad.clone.reckless’, and will be inherited to any of its subdatasets. [Default: None]
- **jobs** (*int or None or {'auto'}, optional*) – how many parallel jobs (where possible) to use. [Default: ‘auto’]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an *IncompleteResultsError* that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a *ValueError* exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) –

return value behavior switch. If 'item-or-list' a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: 'list']

## datalad.api.install

```
datalad.api.install (path=None, source=None, dataset=None, get_data=False, description=None, recursive=False, recursion_limit=None, save=True, reckless=None, jobs='auto')
```

Install a dataset from a (remote) source.

This command creates a local sibling of an existing dataset from a (remote) location identified via a URL or path. Optional recursion into potential subdatasets, and download of all referenced data is supported. The new dataset can be optionally registered in an existing superdataset by identifying it via the *dataset* argument (the new dataset's path needs to be located within the superdataset for that).

It is recommended to provide a brief description to label the dataset's nature *and* location, e.g. "Michael's music on black laptop". This helps humans to identify data locations in distributed scenarios. By default an identifier comprised of user and machine name, plus path will be generated.

When only partial dataset content shall be obtained, it is recommended to use this command without the *get-data* flag, followed by a `~datalad.api.get` operation to obtain the desired data.

---

**Note:** Power-user info: This command uses git clone, and git annex init to prepare the dataset. Registering to a superdataset is performed via a git submodule add operation in the discovered superdataset.

---

## Examples

Install a dataset from Github into the current directory:

```
> install(source='https://github.com/datalad-datasets/longnow-podcasts.git')
```

Install a dataset as a subdataset into the current dataset:

```
> install(dataset='.',
          source='https://github.com/datalad-datasets/longnow-podcasts.git')
```

Install a dataset, and get all content right away:

```
> install(source='https://github.com/datalad-datasets/longnow-podcasts.git',
          get_data=True)
```

Install a dataset with all its subdatasets:

```
> install(source='https://github.com/datalad-datasets/longnow-podcasts.git',
          recursive=True)
```

## Parameters

- **path** – path/name of the installation target. If no *path* is provided a destination path will be derived from a source URL similar to git clone. [Default: None]
- **source** (*str* or *None*, *optional*) – URL or local path of the installation source. [Default: None]

- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the install operation on. If no dataset is given, an attempt is made to identify the dataset in a parent directory of the current working directory and/or the *path* given. [Default: None]
- **get\_data** (*bool, optional*) – if given, obtain all data content too. [Default: False]
- **description** (*str or None, optional*) – short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **save** (*bool, optional*) – by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior. [Default: True]
- **reckless** (*{None, True, False, 'auto', 'ephemeral'}, optional*) – Set up the dataset to be able to obtain content in the cheapest/fastest possible way, even if this poses a potential risk the data integrity ('auto': hardlink files from a local clone of the dataset, 'ephemeral': symlink annex to origin’s annex and discard local availability info via git-annex-dead 'here'. Please note, that with a symlinked annex you share the annex with origin w/o git-annex knowing this. In case of a change in origin you need to update the clone before you’re able to save new content on your end.). Use with care, and limit to “read-only” use cases. With this flag the installed dataset will be marked as untrusted. The reckless mode is stored in a dataset’s local configuration under 'datalad.clone.reckless', and will be inherited to any of its subdatasets. [Default: None]
- **jobs** (*int or None or {'auto'}, optional*) – how many parallel jobs (where possible) to use. [Default: 'auto']
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: 'ignore' any failure is reported, but does not cause an exception; 'continue' if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; 'stop': processing will stop on first failure and an exception is raised. A failure is any result with status 'impossible' or 'error'. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: 'continue']
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]

- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.publish

`datalad.api.publish` (*path=None, dataset=None, to=None, since=None, missing='fail', force=False, transfer\_data='auto', recursive=False, recursion\_limit=None, git\_opts=None, annex\_opts=None, annex\_copy\_opts=None, jobs=None*)

Publish a dataset to a known sibling.

This makes the last saved state of a dataset available to a sibling or special remote data store of a dataset. Any target sibling must already exist and be known to the dataset.

Optionally, it is possible to limit publication to change sets relative to a particular point in the version history of a dataset (e.g. a release tag). By default, the state of the local dataset is evaluated against the last known state of the target sibling. Actual publication is only attempted if there was a change compared to the reference state, in order to speed up processing of large collections of datasets. Evaluation with respect to a particular “historic” state is only supported in conjunction with a specified reference dataset. Change sets are also evaluated recursively, i.e. only those subdatasets are published where a change was recorded that is reflected in to current state of the top-level reference dataset. See “since” option for more information.

Only publication of saved changes is supported. Any unsaved changes in a dataset (hierarchy) have to be saved before publication.

---

**Note:** Power-user info: This command uses git push, and git annex copy to publish a dataset. Publication targets are either configured remote Git repositories, or git-annex special remotes (if they support data upload).

---

### Parameters

- **path** (*sequence of str or None, optional*) – path(s), that may point to file handle(s) to publish including their actual content or to subdataset(s) to be published. If a file handle is published with its data, this implicitly means to also publish the (sub)dataset it belongs to. ‘.’ as a path is treated in a special way in the sense, that it is passed to subdatasets in case *recursive* is also given. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the (top-level) dataset to be published. If no dataset is given, the datasets are determined based on the input arguments. [Default: None]
- **to** (*str or None, optional*) – name of the target sibling. If no name is given an attempt is made to identify the target based on the dataset’s configuration (i.e. a configured tracking branch, or a single sibling that is configured for publication). [Default: None]
- **since** (*str or None, optional*) – When publishing dataset(s), specifies commit (treeish, tag, etc) from which to look for changes to decide whether updated publishing is necessary for this and which children. If empty argument is provided, then we would take from the previously published to that remote/sibling state (for the current branch). [Default: None]
- **missing** (*{'fail', 'inherit', 'skip'}, optional*) – action to perform, if a sibling does not exist in a given dataset. By default it would fail the run (‘fail’ setting). With ‘inherit’ a ‘create-sibling’ with ‘–inherit-settings’ will be used to create sibling on the remote. With ‘skip’ - it simply will be skipped. [Default: ‘fail’]

- **force** (*bool, optional*) – enforce doing publish activities (git push etc) regardless of the analysis if they seemed needed. [Default: False]
- **transfer\_data** (*{'auto', 'none', 'all'}, optional*) – ADDME. [Default: 'auto']
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **git\_opts** (*str or None, optional*) – option string to be passed to git calls. [Default: None]
- **annex\_opts** (*str or None, optional*) – option string to be passed to git annex calls. [Default: None]
- **annex\_copy\_opts** (*str or None, optional*) – option string to be passed to git annex copy calls. [Default: None]
- **jobs** (*int or None or {'auto'}, optional*) – how many parallel jobs (where possible) to use. [Default: None]

## datalad.api.remove

`datalad.api.remove` (*path=None, dataset=None, recursive=False, check=True, save=True, message=None, if\_dirty='save-before'*)

Remove components from datasets

This command can remove subdatasets and paths, including non-empty directories, from datasets. Removing a component implies dropping present content and uninstalling associated subdatasets. Subsequently, the component is “unregistered” from the respective dataset. This means that the component is no longer present on the file system.

By default, the availability of at least one remote copy is verified before file content is dropped. As these checks could lead to slow operation (network latencies, etc), they can be disabled.

## Examples

Permanently remove a subdataset from a dataset and wipe out the subdataset association too:

```
> remove(dataset='path/to/dataset', path='path/to/subds')
```

Permanently remove a dataset and all subdatasets:

```
> remove(dataset='path/to/dataset', recursive=True)
```

Permanently remove a dataset and all subdatasets even if there are fewer than the configured minimum number of (remote) sources for data:

```
> remove(dataset='path/to/dataset', recursive=True, check=False)
```

## Parameters

- **path** (*sequence of str or None, optional*) – path/name of the component to be removed. [Default: None]

- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the operation on. If no dataset is given, an attempt is made to identify a dataset based on the *path* given. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **check** (*bool, optional*) – whether to perform checks to assure the configured minimum number (remote) source for data. [Default: True]
- **save** (*bool, optional*) – by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior. [Default: True]
- **message** (*str or None, optional*) – a description of the state or the changes made to a dataset. [Default: None]
- **if\_dirty** – desired behavior if a dataset with unsaved changes is discovered: ‘fail’ will trigger an error and further processing is aborted; ‘save-before’ will save all changes prior any further action; ‘ignore’ let’s datalad proceed as if the dataset would not have unsaved changes. [Default: ‘save-before’]
- **on\_failure** (*{‘ignore’, ‘continue’, ‘stop’}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{‘default’, ‘json’, ‘json\_pp’, ‘tailored’} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{‘datasets’, ‘successdatasets-or-none’, ‘paths’, ‘relpaths’, ‘metadata’} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{‘generator’, ‘list’, ‘item-or-list’}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.save

`datalad.api.save` (*path=None, message=None, dataset=None, version\_tag=None, recursive=False, recursion\_limit=None, updated=False, message\_file=None, to\_git=None*)

Save the current state of a dataset

Saving the state of a dataset records changes that have been made to it. This change record is annotated with a user-provided description. Optionally, an additional tag, such as a version, can be assigned to the saved state. Such tag enables straightforward retrieval of past versions at a later point in time.

---

**Note:** Before Git v2.22, any Git repository without an initial commit located inside a Dataset is ignored, and content underneath it will be saved to the respective superdataset. DataLad datasets always have an initial commit, hence are not affected by this behavior.

---

## Examples

Save any content underneath the current directory, without altering any potential subdataset:

```
> save(path='.')
```

Save specific content in the dataset:

```
> save(path='myfile.txt')
```

Attach a commit message to save:

```
> save(path='myfile.txt', message='add file')
```

Save any content underneath the current directory, and recurse into any potential subdatasets:

```
> save(path='.', recursive=True)
```

Save any modification of known dataset content in the current directory, but leave untracked files (e.g. temporary files) untouched:

```
> save(path='.', updated=True)
```

Tag the most recent saved state of a dataset:

```
> save(version_tag='bestyet')
```

## Parameters

- **path** (*sequence of str or None, optional*) – path/name of the dataset component to save. If given, only changes made to those components are recorded in the new state. [Default: None]
- **message** (*str or None, optional*) – a description of the state or the changes made to a dataset. [Default: None]
- **dataset** (*Dataset or None, optional*) – “specify the dataset to save. [Default: None]
- **version\_tag** (*str or None, optional*) – an additional marker for that state. Every dataset that is touched will receive the tag. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]



- **updated** (*bool, optional*) – if given, only saves previously tracked paths. [Default: False]
- **message\_file** (*str or None, optional*) – take the commit message from this file. This flag is mutually exclusive with -m. [Default: None]
- **to\_git** (*bool, optional*) – flag whether to add data directly to Git, instead of tracking data identity only. Usually this is not desired, as it inflates dataset sizes and impacts flexibility of data transport. If not specified - it will be up to git-annex to decide, possibly on .gitattributes options. Use this flag with a simultaneous selection of paths to save. In general, it is better to pre-configure a dataset to track particular paths, file types, or file sizes with either Git or git-annex. See <https://git-annex.branchable.com/tips/largefiles/>. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an IncompleteResultsError that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a ValueError exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.update

`datalad.api.update` (*path=None, sibling=None, merge=False, follow='sibling', dataset=None, recursive=False, recursion\_limit=None, fetch\_all=None, reobtain\_data=False*)

Update a dataset from a sibling.

### Examples

Update from a particular sibling:

```
> update(sibling='siblingname')
```

Update from a particular sibling and merge the changes from a configured or matching branch from the sibling (see *follow* for details):

```
> update(sibling='siblingname', merge=True)
```

Update from the sibling 'origin', traversing into subdatasets. For subdatasets, merge the revision registered in the parent dataset into the current branch:

```
> update(sibling='origin', merge=True, follow='parentds', recursive=True)
```

### Parameters

- **path** (*sequence of str or None, optional*) – constrain to-be-updated sub-datasets to the given path for recursive operation. [Default: None]
- **sibling** (*str or None, optional*) – name of the sibling to update from. If no sibling is given, updates from all siblings are obtained. [Default: None]
- **merge** (*bool or {'any', 'ff-only'}, optional*) – merge obtained changes from the sibling. If a sibling is not explicitly given and there is only a single known sibling, that sibling is used. Otherwise, an unspecified sibling defaults to the configured remote for the current branch. By default, changes are fetched from the sibling but not merged into the current branch. With `merge=True` or `merge="any"`, the changes will be merged into the current branch. A value of 'ff-only' restricts the allowed merges to fast-forwards. [Default: False]
- **follow** (*{'sibling', 'parentds'}, optional*) – source of updates for sub-datasets. For 'sibling', the update will be done by merging in a branch from the (specified or inferred) sibling. The branch brought in will either be the current branch's configured branch, if it points to a branch that belongs to the sibling, or a sibling branch with a name that matches the current branch. For 'parentds', the revision registered in the parent dataset of the subdataset is merged in. Note that the current dataset is always updated according to 'sibling'. This option has no effect unless a merge is requested and `recursive=True` is specified. [Default: 'sibling']
- **dataset** (*Dataset or None, optional*) – specify the dataset to update. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **fetch\_all** (*bool, optional*) – this option has no effect and will be removed in a future version. When no siblings are given, an all-sibling update will be performed. [Default: None]
- **reobtain\_data** (*bool, optional*) – if enabled, file content that was present before an update will be re-obtained in case a file was changed by the update. [Default: False]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: 'ignore' any failure is reported, but does not cause an exception; 'continue' if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; 'stop': processing will stop on first failure and an exception is raised. A failure is any result with status 'impossible' or 'error'. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: 'continue']

- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a ValueError exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.uninstall

`datalad.api.uninstall` (*path=None, dataset=None, recursive=False, check=True, if\_dirty='save-before'*)

Uninstall subdatasets

This command can be used to uninstall any number of installed subdatasets. This command will error if individual files or non-dataset directories are given as input (use the drop or remove command depending on the desired goal), nor will it uninstall top-level datasets (i.e. datasets that are not a subdataset in another dataset; use the remove command for this purpose).

By default, the availability of at least one remote copy for each currently available file in any dataset is verified. As these checks could lead to slow operation (network latencies, etc), they can be disabled.

Any number of paths to process can be given as input. Recursion into subdatasets needs to be explicitly enabled, while recursion into subdirectories within a dataset is done automatically. An optional recursion limit is applied relative to each given input path.

## Examples

Uninstall a subdataset (undo installation):

```
> uninstall(path='path/to/subds')
```

Uninstall a subdataset and all potential subdatasets:

```
> uninstall(path='path/to/subds', recursive=True)
```

Skip checks that ensure a minimal number of (remote) sources:

```
> uninstall(path='path/to/subds', check=False)
```

## Parameters

- **path** (*sequence of str or None, optional*) – path/name of the component to be uninstalled. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the operation on. If no dataset is given, an attempt is made to identify a dataset based on the *path* given. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **check** (*bool, optional*) – whether to perform checks to assure the configured minimum number (remote) source for data. [Default: True]
- **if\_dirty** – desired behavior if a dataset with unsaved changes is discovered: ‘fail’ will trigger an error and further processing is aborted; ‘save-before’ will save all changes prior any further action; ‘ignore’ let’s datalad proceed as if the dataset would not have unsaved changes. [Default: ‘save-before’]
- **on\_failure** (*{‘ignore’, ‘continue’, ‘stop’}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{‘default’, ‘json’, ‘json\_pp’, ‘tailored’} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{‘datasets’, ‘successdatasets-or-none’, ‘paths’, ‘relpaths’, ‘metadata’} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{‘generator’, ‘list’, ‘item-or-list’}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.unlock

`datalad.api.unlock` (*path=None, dataset=None, recursive=False, recursion\_limit=None*)

Unlock file(s) of a dataset

Unlock files of a dataset in order to be able to edit the actual content

## Examples

Unlock a single file:

```
> unlock(path='path/to/file')
```

Unlock all contents in the dataset:

```
> unlock('.')
```

## Parameters

- **path** (*sequence of str or None, optional*) – file(s) to unlock. [Default: None]
- **dataset** (*Dataset or None, optional*) – “specify the dataset to unlock files in. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## Metadata handling

|   |  |
|---|--|
| <code>api.search([query, dataset, force_reindex, ...])</code> | Search dataset metadata  |
| <code>api.metadata([path, dataset, ...])</code>               | Metadata reporting for files and entire datasets                       |
| <code>api.aggregate_metadata([path, dataset, ...])</code>     | Aggregate metadata of one or more datasets for later query.            |
| <code>api.extract_metadata(types[, files, dataset])</code>    | Run one or more of DataLad's metadata extractors on a dataset or file. |

## datalad.api.search

`datalad.api.search` (*query=None, dataset=None, force\_reindex=False, max\_nresults=None, mode=None, full\_record=False, show\_keys=None, show\_query=False*)

Search dataset metadata

DataLad can search metadata extracted from a dataset and/or aggregated into a superdataset (see the *aggregate-metadata* command). This makes it possible to discover datasets, or individual files in a dataset even when they are not available locally.

Ultimately DataLad metadata are a graph of linked data structures. However, this command does not (yet) support queries that can exploit all information stored in the metadata. At the moment the following search modes are implemented that represent different trade-offs between the expressiveness of a query and the computational and storage resources required to execute a query.

- `egrep` (default)
- `egrepcs` [case-sensitive `egrep`]
- `textblob`
- `autofield`

An alternative default mode can be configured by tuning the configuration variable ‘`datalad.search.default-mode`’:

```
[datalad "search"]
  default-mode = egrepcs
```

Each search mode has its own default configuration for what kind of documents to query. The respective default can be changed via configuration variables:

```
[datalad "search"]
  index-<mode_name>-documenttype = (all|datasets|files)
```

*Mode: `egrep/egrepcs`*

These search modes are largely ignorant of the metadata structure, and simply perform matching of a search pattern against a flat string-representation of metadata. This is advantageous when the query is simple and the metadata structure is irrelevant, or precisely known. Moreover, it does not require a search index, hence results can be reported without an initial latency for building a search index when the underlying metadata has changed (e.g. due to a dataset update). By default, these search modes only consider datasets and do not investigate records for individual files for speed reasons. Search results are reported in the order in which they were discovered.

Queries can make use of Python regular expression syntax (<https://docs.python.org/3/library/re.html>). In *egrep* mode, matching is case-insensitive when the query does not contain upper case characters, but is case-sensitive when it does. In *egrepcs* mode, matching is always case-sensitive. Expressions will match anywhere in a metadata string, not only at the start.

When multiple queries are given, all queries have to match for a search hit (AND behavior).

It is possible to search individual metadata key/value items by prefixing the query with a metadata key name, separated by a colon (':'). The key name can also be a regular expression to match multiple keys. A query match happens when any value of an item with a matching key name matches the query (OR behavior). See examples for more information.

## Examples

Query for (what happens to be) an author:

```
% datalad search haxby
```

Queries are case-INsensitive when the query contains no upper case characters, and can be regular expressions. Use *egrepcs* mode when it is desired to perform a case-sensitive lowercase match:

```
% datalad search --mode egrepcs halchenko.*haxby
```

This search mode performs NO analysis of the metadata content. Therefore queries can easily fail to match. For example, the above query implicitly assumes that authors are listed in alphabetical order. If that is the case (which may or may not be true), the following query would yield NO hits:

```
% datalad search Haxby.*Halchenko
```

The `textblob` search mode represents an alternative that is more robust in such cases.

For more complex queries multiple query expressions can be provided that all have to match to be considered a hit (AND behavior). This query discovers all files (non-default behavior) that match 'bids.type=T1w' AND 'nifti1.qform\_code=scanner':

```
% datalad -c datalad.search.index-egrep-documenttype=all search bids.type:T1w_
↳nifti1.qform_code:scanner
```

Key name selectors can also be expressions, which can be used to select multiple keys or construct “fuzzy” queries. In such cases a query matches when any item with a matching key matches the query (OR behavior). However, multiple queries are always evaluated using an AND conjunction. The following query extends the example above to match any files that have either 'nifti1.qform\_code=scanner' or 'nifti1.sform\_code=scanner':

```
% datalad -c datalad.search.index-egrep-documenttype=all search bids.type:T1w_
↳nifti1.(q|s)form_code:scanner
```

### Mode: *textblob*

This search mode is very similar to the `egrep` mode, but with a few key differences. A search index is built from the string-representation of metadata records. By default, only datasets are included in this index, hence the indexing is usually completed within a few seconds, even for hundreds of datasets. This mode uses its own query language (not regular expressions) that is similar to other search engines. It supports logical conjunctions and fuzzy search terms. More information on this is available from the Whoosh project (search engine implementation):

- Description of the Whoosh query language: <http://whoosh.readthedocs.io/en/latest/querylang.html>
- Description of a number of query language customizations that are enabled in DataLad, such as, fuzzy term matching: <http://whoosh.readthedocs.io/en/latest/parsing.html#common-customizations>

Importantly, search hits are scored and reported in order of descending relevance, hence limiting the number of search results is more meaningful than in the 'egrep' mode and can also reduce the query duration.

## Examples

Search for (what happens to be) two authors, regardless of the order in which those names appear in the metadata:

```
% datalad search --mode textblob halchenko haxby
```

Fuzzy search when you only have an approximate idea what you are looking for or how it is spelled:

```
% datalad search --mode textblob haxbi~
```

Very fuzzy search, when you are basically only confident about the first two characters and how it sounds approximately (or more precisely: allow for three edits and require matching of the first two characters):

```
% datalad search --mode textblob haksbi~3/2
```

Combine fuzzy search with logical constructs:

```
% datalad search --mode textblob 'haxbi~ AND (hanke OR halchenko)'
```

### *Mode: autofield*

This mode is similar to the ‘textblob’ mode, but builds a vastly more detailed search index that represents individual metadata variables as individual fields. By default, this search index includes records for datasets and individual fields, hence it can grow very quickly into a huge structure that can easily take an hour or more to build and require more than a GB of storage. However, limiting it to documents on datasets (see above) retains the enhanced expressiveness of queries while dramatically reducing the resource demands.

## Examples

List names of search index fields (auto-discovered from the set of indexed datasets):

```
% datalad search --mode autofield --show-keys name
```

Fuzzy search for datasets with an author that is specified in a particular metadata field:

```
% datalad search --mode autofield bids.author:haxbi~ type:dataset
```

Search for individual files that carry a particular description prefix in their ‘nifti1’ metadata:

```
% datalad search --mode autofield nifti1.description:FSL* type:file
```

### *Reporting*

Search hits are returned as standard DataLad results. On the command line the ‘–output-format’ (or ‘-f’) option can be used to tweak results for further processing.

## Examples

Format search hits as a JSON stream (one hit per line):

```
% datalad -f json search haxby
```

Custom formatting: which terms matched the query of particular results. Useful for investigating fuzzy search results:



```
$ datalad -f '{path}: {query_matched}' search --mode autofield bids.author:haxbi~
```

## Parameters

- **query** – query string, supported syntax and features depends on the selected search mode (see documentation). [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to perform the query operation on. If no dataset is given, an attempt is made to identify the dataset based on the current working directory and/or the *path* given. [Default: None]
- **force\_reindex** (*bool, optional*) – force rebuilding the search index, even if no change in the dataset’s state has been detected, for example, when the index documenttype configuration has changed. [Default: False]
- **max\_nresults** (*int or None, optional*) – maximum number of search results to report. Setting this to 0 will report all search matches. Depending on the mode this can search substantially slower. If not specified, a mode-specific default setting will be used. [Default: None]
- **mode** – Mode of search index structure and content. See section SEARCH MODES for details. [Default: None]
- **full\_record** (*bool, optional*) – If set, return the full metadata record for each search hit. Depending on the search mode this might require additional queries. By default, only data that is available to the respective search modes is returned. This always includes essential information, such as the path and the type. [Default: False]
- **show\_keys** – if given, a list of known search keys is shown. If ‘name’ - only the name is printed one per line. If ‘short’ or ‘full’, statistics (in how many datasets, and how many unique values) are printed. ‘short’ truncates the listing of unique values. No other action is performed (except for reindexing), even if other arguments are given. Each key is accompanied by a term definition in parenthesis (TODO). In most cases a definition is given in the form of a URL. If an ontology definition for a term is known, this URL can resolve to a webpage that provides a comprehensive definition of the term. However, for speed reasons term resolution is solely done on information contained in a local dataset’s metadata, and definition URLs might be outdated or point to no longer existing resources. [Default: None]
- **show\_query** (*bool, optional*) – if given, the formal query that was generated from the given query string is shown, but not actually executed. This is mostly useful for debugging purposes. [Default: False]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]

- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'}* or callable or *None*, optional) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}*, optional) – return value behavior switch. If 'item-or-list' a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: 'list']

## datalad.api.metadata

`datalad.api.metadata` (*path=None, dataset=None, get\_aggregates=False, reporton='all', recursive=False*)

Metadata reporting for files and entire datasets

Two types of metadata are supported:

1. metadata describing a dataset as a whole (dataset-global metadata), and
2. metadata for files in a dataset (content metadata).

Both types can be accessed with this command.

## Examples

Report the metadata of a single file, as aggregated into the closest locally available dataset, containing the query path:

```
% datalad metadata somedir/subdir/thisfile.dat
```

Sometimes it is helpful to get metadata records formatted in a more accessible form, here as pretty-printed JSON:

```
% datalad -f json_pp metadata somedir/subdir/thisfile.dat
```

Same query as above, but specify which dataset to query (must be containing the query path):

```
% datalad metadata -d . somedir/subdir/thisfile.dat
```

Report any metadata record of any dataset known to the queried dataset:

```
% datalad metadata --recursive --reporton datasets
```

Get a JSON-formatted report of aggregated metadata in a dataset, incl. information on enabled metadata extractors, dataset versions, dataset IDs, and dataset paths:

```
% datalad -f json metadata --get-aggregates
```

## Parameters

- **path** (*sequence of str or None, optional*) – path(s) to query metadata for. [Default: None]

- **dataset** (*Dataset or None, optional*) – dataset to query. If given, metadata will be reported as stored in this dataset. Otherwise, the closest available dataset containing a query path will be consulted. [Default: None]
- **get\_aggregates** (*bool, optional*) – if set, yields all (sub)datasets for which aggregate metadata are available in the dataset. No other action is performed, even if other arguments are given. The reported results contain a datasets’s ID, the commit hash at which metadata aggregation was performed, and the location of the object file(s) containing the aggregated metadata. [Default: False]
- **reporton** (*{'all', 'datasets', 'files', 'none'}, optional*) – choose on what type result to report on: ‘datasets’, ‘files’, ‘all’ (both datasets and files), or ‘none’ (no report). [Default: ‘all’]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top- level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## datalad.api.aggregate\_metadata

`datalad.api.aggregate_metadata` (*path=None, dataset=None, recursive=False, recursion\_limit=None, update\_mode='target', incremental=False, force\_extraction=False, save=True*)

Aggregate metadata of one or more datasets for later query.

Metadata aggregation refers to a procedure that extracts metadata present in a dataset into a portable representation that is stored a single standardized format. Moreover, metadata aggregation can also extract metadata in this format from one dataset and store it in another (super)dataset. Based on such collections of aggregated

metadata it is possible to discover particular datasets and specific parts of their content, without having to obtain the target datasets first (see the DataLad ‘search’ command).

To enable aggregation of metadata that are contained in files of a dataset, one has to enable one or more metadata extractor for a dataset. DataLad supports a number of common metadata standards, such as the Exchangeable Image File Format (EXIF), Adobe’s Extensible Metadata Platform (XMP), and various audio file metadata systems like ID3. DataLad extension packages can provide metadata data extractors for additional metadata sources. For example, the neuroimaging extension provides extractors for scientific (meta)data standards like BIDS, DICOM, and NIFTI1. Some metadata extractors depend on particular 3rd-party software. The list of metadata extractors available to a particular DataLad installation is reported by the ‘wtf’ command (‘datalad wtf’).

Enabling a metadata extractor for a dataset is done by adding its name to the ‘datalad.metadata.nativetype’ configuration variable – typically in the dataset’s configuration file (.datalad/config), e.g.:

```
[datalad "metadata"]
  nativetype = exif
  nativetype = xmp
```

If an enabled metadata extractor is not available in a particular DataLad installation, metadata extraction will not succeed in order to avoid inconsistent aggregation results.

Enabling multiple extractors is supported. In this case, metadata are extracted by each extractor individually, and stored alongside each other. Metadata aggregation will also extract DataLad’s own metadata (extractors ‘datalad\_core’, and ‘annex’).

Metadata aggregation can be performed recursively, in order to aggregate all metadata across all subdatasets, for example, to be able to search across any content in any dataset of a collection. Aggregation can also be performed for subdatasets that are not available locally. In this case, pre-aggregated metadata from the closest available superdataset will be considered instead.

Depending on the versatility of the present metadata and the number of dataset or files, aggregated metadata can grow prohibitively large. A number of configuration switches are provided to mitigate such issues.

**datalad.metadata.aggregate-content-<extractor-name>** If set to false, content metadata aggregation will not be performed for the named metadata extractor (a potential underscore ‘\_’ in the extractor name must be replaced by a dash ‘-’). This can substantially reduce the runtime for metadata extraction, and also reduce the size of the generated metadata aggregate. Note, however, that some extractors may not produce any metadata when this is disabled, because their metadata might come from individual file headers only. ‘datalad.metadata.store-aggregate-content’ might be a more appropriate setting in such cases.

**datalad.metadata.aggregate-ignore-fields** Any metadata key matching any regular expression in this configuration setting is removed prior to generating the dataset-level metadata summary (keys and their unique values across all dataset content), and from the dataset metadata itself. This switch can also be used to filter out sensitive information prior aggregation.

**datalad.metadata.generate-unique-<extractor-name>** If set to false, DataLad will not auto-generate a summary of unique content metadata values for a particular extractor as part of the dataset-global metadata (a potential underscore ‘\_’ in the extractor name must be replaced by a dash ‘-’). This can be useful if such a summary is bloated due to minor uninformative (e.g. numerical) differences, or when a particular extractor already provides a carefully designed content metadata summary.

**datalad.metadata.maxfieldsize** Any metadata value that exceeds the size threshold given by this configuration setting (in bytes/characters) is removed.

**datalad.metadata.store-aggregate-content** If set, extracted content metadata are still used to generate a dataset-level summary of present metadata (all keys and their unique values across all files in a dataset are determined and stored as part of the dataset-level metadata aggregate, see `datalad.metadata.generate-unique-<extractor-name>`), but metadata on individual files are not stored. This switch can be used to avoid

prohibitively large metadata files. Discovery of datasets containing content matching particular metadata properties will still be possible, but such datasets would have to be obtained first in order to discover which particular files in them match these properties.

### Parameters

- **path** (*sequence of str or None, optional*) – path to datasets that shall be aggregated. When a given path is pointing into a dataset, the metadata of the containing dataset will be aggregated. If no paths given, current dataset metadata is aggregated. [Default: None]
- **dataset** (*Dataset or None, optional*) – topmost dataset metadata will be aggregated into. All dataset between this dataset and any given path will receive updated aggregated metadata from all given paths. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **update\_mode** (*{'all', 'target'}, optional*) – which datasets to update with newly aggregated metadata: all datasets from any leaf dataset to the top-level target dataset including all intermediate datasets (all), or just the top-level target dataset (target). [Default: 'target']
- **incremental** (*bool, optional*) – If set, all information on metadata records of sub-datasets that have not been (re-)aggregated in this run will be kept unchanged. This is useful when (re-)aggregation only a subset of a dataset hierarchy, for example, because not all subdatasets are locally available. [Default: False]
- **force\_extraction** (*bool, optional*) – If set, all enabled extractors will be engaged regardless of whether change detection indicates that metadata has already been extracted for a given dataset state. [Default: False]
- **save** (*bool, optional*) – by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior. [Default: True]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: 'ignore' any failure is reported, but does not cause an exception; 'continue' if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; 'stop': processing will stop on first failure and an exception is raised. A failure is any result with status 'impossible' or 'error'. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: 'continue']
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable's return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations

that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]

- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.extract\_metadata

`datalad.api.extract_metadata` (*types, files=None, dataset=None*)

Run one or more of DataLad’s metadata extractors on a dataset or file.

The result(s) are structured like the metadata DataLad would extract during metadata aggregation. There is one result per dataset/file.

### Examples

Extract metadata with two extractors from a dataset in the current directory and also from all its files:

```
$ datalad extract-metadata -d . --type frictionless_datapackage --type datalad_
↪core
```

Extract XMP metadata from a single PDF that is not part of any dataset:

```
$ datalad extract-metadata --type xmp Downloads/freshfromtheweb.pdf
```

### Parameters

- **types** – Name of a metadata extractor to be executed.
- **files** (*sequence of str or None, optional*) – Path of a file to extract metadata from. [Default: None]
- **dataset** (*Dataset or None, optional*) – “Dataset to extract metadata from. If no *file* is given, metadata is extracted from all files of the dataset. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value

becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]

- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## Reproducible execution

|  |  |
|--|--|
| <code>api.run(cmd, dataset, inputs, outputs, ...)</code>     | Run an arbitrary shell command and record its impact on a dataset. |
| <code>api.rerun([revision, since, dataset, ...])</code>      | Re-execute previous <i>datalad run</i> commands.                   |
| <code>api.run_procedure(spec, dataset, discover, ...)</code> | Run prepared procedures (DataLad scripts) on a dataset             |

## datalad.api.run

`datalad.api.run` (*cmd=None, dataset=None, inputs=None, outputs=None, expand=None, explicit=False, message=None, sidecar=None*)

Run an arbitrary shell command and record its impact on a dataset.

It is recommended to craft the command such that it can run in the root directory of the dataset that the command will be recorded in. However, as long as the command is executed somewhere underneath the dataset root, the exact location will be recorded relative to the dataset root.

If the executed command did not alter the dataset in any way, no record of the command execution is made.

If the given command errors, a *CommandError* exception with the same exit code will be raised, and no modifications will be saved.

### Command format

A few placeholders are supported in the command via Python format specification. “{pwd}” will be replaced with the full path of the current working directory. “{dspath}” will be replaced with the full path of the dataset that run is invoked on. “{tmpdir}” will be replaced with the full path of a temporary directory. “{inputs}” and “{outputs}” represent the values specified by *inputs* and *outputs*. If multiple values are specified, the values will be joined by a space. The order of the values will match that order from the command line, with any globs expanded in alphabetical order (like bash). Individual values can be accessed with an integer index (e.g., “{inputs[0]}”).

Note that the representation of the inputs or outputs in the formatted command string depends on whether the command is given as a list of arguments or as a string. The concatenated list of inputs or outputs will be surrounded by quotes when the command is given as a list but not when it is given as a string. This means that the string form is required if you need to pass each input as a separate argument to a preceding script (i.e., write the command as “./script {inputs}”, quotes included). The string form should also be used if the input or output paths contain spaces or other characters that need to be escaped.

To escape a brace character, double it (i.e., “{{” or “}}”).

Custom placeholders can be added as configuration variables under “datalad.run.substitutions”. As an example:

Add a placeholder “name” with the value “joe”:

```
% git config --file=.datalad/config datalad.run.substitutions.name joe
% datalad add -m "Configure name placeholder" .datalad/config
```

Access the new placeholder in a command:

```
% datalad run "echo my name is {name} >me"
```

## Examples

Run an executable script and record the impact on a dataset:

```
> run(message='run my script', cmd='code/script.sh')
```

Run a command and specify a directory as a dependency for the run. The contents of the dependency will be retrieved prior to running the script:

```
> run(cmd='code/script.sh', message='run my script',
      inputs=['data/*'])
```

Run an executable script and specify output files of the script to be unlocked prior to running the script:

```
> run(cmd='code/script.sh', message='run my script',
      inputs=['data/*'], outputs=['output_dir'])
```

Specify multiple inputs and outputs:

```
> run(cmd='code/script.sh',
      message='run my script',
      inputs=['data/*', 'datafile.txt'],
      outputs=['output_dir', 'outfile.txt'])
```

## Parameters

- **cmd** – command for execution. A leading ‘-’ can be used to disambiguate this command from the preceding options to DataLad. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to record the command results in. An attempt is made to identify the dataset based on the current working directory. If a dataset is given, the command will be executed in the root directory of this dataset. [Default: None]
- **inputs** – A dependency for the run. Before running the command, the content of this file will be retrieved. A value of “.” means “run datalad get .”. The value can also be a glob. [Default: None]
- **outputs** – Prepare this file to be an output file of the command. A value of “.” means “run datalad unlock .” (and will fail if some content isn’t present). For any other value, if the content of this file is present, unlock the file. Otherwise, remove it. The value can also be a glob. [Default: None]
- **expand** (*{None, 'inputs', 'outputs', 'both'}, optional*) – Expand globs when storing inputs and/or outputs in the commit message. [Default: None]
- **explicit** (*bool, optional*) – Consider the specification of inputs and outputs to be explicit. Don’t warn if the repository is dirty, and only save modifications to the listed outputs. [Default: False]
- **message** (*str or None, optional*) – a description of the state or the changes made to a dataset. [Default: None]



- **sidecar** (*None or bool, optional*) – By default, the configuration variable ‘datalad.run.record-sidecar’ determines whether a record with information on a command’s execution is placed into a separate record file instead of the commit message (default: off). This option can be used to override the configured behavior on a case-by-case basis. Sidecar files are placed into the dataset’s ‘.datalad/runinfo’ directory (customizable via the ‘datalad.run.record-directory’ configuration variable). [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## datalad.api.rerun

`datalad.api.rerun` (*revision='HEAD', since=None, dataset=None, branch=None, message=None, onto=None, script=None, report=False, explicit=False*)

Re-execute previous `datalad run` commands.

This will unlock any dataset content that is on record to have been modified by the command in the specified revision. It will then re-execute the command in the recorded path (if it was inside the dataset). Afterwards, all modifications will be saved.

### Report mode

When called with `report=True`, this command reports information about what would be re-executed as a series of records. There will be a record for each revision in the specified revision range. Each of these will have one of the following “`rerun_action`” values:

- `run`: the revision has a recorded command that would be re-executed
- `skip-or-pick`: the revision does not have a recorded command and would be either skipped or cherry picked
- `merge`: the revision is a merge commit and a corresponding merge would be made

The decision to skip rather than cherry pick a revision is based on whether the revision would be reachable from HEAD at the time of execution.

In addition, when a starting point other than HEAD is specified, there is a `rerun_action` value “checkout”, in which case the record includes information about the revision the would be checked out before rerunning any commands.

---

**Note:** Currently the “onto” feature only sets the working tree of the current dataset to a previous state. The working trees of any subdatasets remain unchanged.

---

## Examples

Re-execute the command from the previous commit:

```
> rerun()
```

Re-execute any commands in the last five commits:

```
> rerun(since='HEAD~5')
```

Do the same as above, but re-execute the commands on top of HEAD~5 in a detached state:

```
> rerun(onto='', since='HEAD~5')
```

## Parameters

- **revision** (*str, optional*) – rerun command(s) in *revision*. By default, the command from this commit will be executed, but *since* can be used to construct a revision range. [Default: ‘HEAD’]
- **since** (*str or None, optional*) – If *since* is a commit-ish, the commands from all commits that are reachable from *revision* but not *since* will be re-executed (in other words, the commands in `git log SINCE..REVISION`). If *SINCE* is an empty string, it is set to the parent of the first commit that contains a recorded command (i.e., all commands in `git log REVISION` will be re-executed). [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset from which to rerun a recorded command. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. If a dataset is given, the command will be executed in the root directory of this dataset. [Default: None]
- **branch** (*str or None, optional*) – create and checkout this branch before rerunning the commands. [Default: None]
- **message** (*str or None, optional*) – use `MESSAGE` for the reran commit rather than the recorded commit message. In the case of a multi-commit rerun, all the reran commits will have this message. [Default: None]
- **onto** (*str or None, optional*) – start point for rerunning the commands. If not specified, commands are executed at HEAD. This option can be used to specify an alternative start point, which will be checked out with the branch name specified by *branch* or in a detached state otherwise. As a special case, an empty value for this option means the parent of the first run commit in the specified revision list. [Default: None]
- **script** (*str or None, optional*) – extract the commands into this file rather than rerunning. Use `-` to write to stdout instead. [Default: None]

- **report** (*bool, optional*) – Don’t actually re-execute anything, just display what would be done. [Default: False]
- **explicit** (*bool, optional*) – Consider the specification of inputs and outputs in the run record to be explicit. Don’t warn if the repository is dirty, and only save modifications to the outputs from the original record. Note that when several run commits are specified, this applies to every one. Care should also be taken when using *onto* because checking out a new HEAD can easily fail when the working tree has modifications. [Default: False]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an *IncompleteResultsError* that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a *ValueError* exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.run\_procedure

`datalad.api.run_procedure` (*spec=None, dataset=None, discover=False, help\_proc=False*)

Run prepared procedures (DataLad scripts) on a dataset

### Concept

A “procedure” is an algorithm with the purpose to process a dataset in a particular way. Procedures can be useful in a wide range of scenarios, like adjusting dataset configuration in a uniform fashion, populating a dataset with particular content, or automating other routine tasks, such as synchronizing dataset content with certain siblings.

Implementations of some procedures are shipped together with DataLad, but additional procedures can be provided by 1) any DataLad extension, 2) any (sub-)dataset, 3) a local user, or 4) a local system administrator. DataLad will look for procedures in the following locations and order:

Directories identified by the configuration settings

- ‘datalad.locations.user-procedures’ (determined by `appdirs.user_config_dir`; defaults to ‘\$HOME/.config/datalad/procedures’ on GNU/Linux systems)

- ‘datalad.locations.system-procedures’ (determined by `appdirs.site_config_dir`; defaults to ‘/etc/xdg/datalad/procedures’ on GNU/Linux systems)
- ‘datalad.locations.dataset-procedures’

and subsequently in the ‘resources/procedures/’ directories of any installed extension, and, lastly, of the DataLad installation itself.

Please note that a dataset that defines ‘datalad.locations.dataset-procedures’ provides its procedures to any dataset it is a subdataset of. That way you can have a collection of such procedures in a dedicated dataset and install it as a subdataset into any dataset you want to use those procedures with. In case of a naming conflict with such a dataset hierarchy, the dataset you’re calling run-procedures on will take precedence over its subdatasets and so on.

Each configuration setting can occur multiple times to indicate multiple directories to be searched. If a procedure matching a given name is found (filename without a possible extension), the search is aborted and this implementation will be executed. This makes it possible for individual datasets, users, or machines to override externally provided procedures (enabling the implementation of customizable processing “hooks”).

#### *Procedure implementation*

A procedure can be any executable. Executables must have the appropriate permissions and, in the case of a script, must contain an appropriate “shebang” line. If a procedure is not executable, but its filename ends with ‘.py’, it is automatically executed by the ‘python’ interpreter (whichever version is available in the present environment). Likewise, procedure implementations ending on ‘.sh’ are executed via ‘bash’.

Procedures can implement any argument handling, but must be capable of taking at least one positional argument (the absolute path to the dataset they shall operate on).

For further customization there are two configuration settings per procedure available:

- ‘datalad.procedures.<NAME>.call-format’ fully customizable format string to determine how to execute procedure NAME (see also `datalad-run`). It currently requires to include the following placeholders:
  - ‘{script}’: will be replaced by the path to the procedure
  - ‘{ds}’: will be replaced by the absolute path to the dataset the procedure shall operate on
  - ‘{args}’: (not actually required) will be replaced by all but the first element of *spec* if *spec* is a list or tuple As an example the default format string for a call to a python script is: “python {script} {ds} {args}”
- ‘datalad.procedures.<NAME>.help’ will be shown on `datalad run-procedure -help-proc NAME` to provide a description and/or usage info for procedure NAME

## **Examples**

Find out which procedures are available on the current system:

```
> run_procedure(discover=True)
```

Run the ‘yoda’ procedure in the current dataset:

```
> run_procedure(spec='cfg_yoda', recursive=True)
```

### **Parameters**

- **spec** – Name and possibly additional arguments of the to-be-executed procedure. [Default: None]

- **dataset** (*Dataset or None, optional*) – specify the dataset to run the procedure on. An attempt is made to identify the dataset based on the current working directory. [Default: None]
- **discover** (*bool, optional*) – if given, all configured paths are searched for procedures and one result record per discovered procedure is yielded, but no procedure is executed. [Default: False]
- **help\_proc** (*bool, optional*) – if given, get a help message for procedure NAME from config setting `datalad.procedures.NAME.help`. [Default: False]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## Plumbing commands

|   |  |
|---|--|
| <code>api.annotate_paths([path, dataset, ...])</code>         | Analyze and act upon input paths                               |
| <code>api.clean([dataset, what, recursive, ...])</code>       | Clean up after DataLad (possible temporary files etc.)         |
| <code>api.clone(source[, path, dataset, ...])</code>          | Obtain a dataset (copy) from a URL or local directory          |
| <code>api.create_test_dataset([path, spec, seed])</code>      | Create test (meta-)dataset.                                    |
| <code>api.diff([path, fr, to, dataset, annex, ...])</code>    | Report differences between two states of a dataset (hierarchy) |
| <code>api.download_url(urls[, dataset, path, ...])</code>     | Download content   |
| <code>api.ls(loc[, recursive, fast, all_, long_, ...])</code> | List summary information about URLs and dataset(s)             |
| <code>api.sshrun(login, cmd[, port, ipv4, ipv6, ...])</code>  | Run command on remote machines via SSH.                        |
| <code>api.siblings([action, dataset, name, url, ...])</code>  | Manage sibling configuration                                   |
| <code>api.subdatasets([path, dataset, fulfilled, ...])</code> | Report subdatasets and their properties.                       |

## datalad.api.annotate\_paths

```
datalad.api.annotate_paths (path=None, dataset=None, recursive=False, recur-
    sion_limit=None, action=None, unavailable_path_status=",
    unavailable_path_msg=None, nondataset_path_status='error',
    force_parentds_discovery=True, force_subds_discovery=True,
    force_no_revision_change_discovery=True,
    force_untracked_discovery=True, modified=None)
```

Analyze and act upon input paths

Given paths (or more generally location requests) are inspected and annotated with a number of properties. A list of recognized properties is provided below.

Input *paths* for this command can either be un-annotated (raw) path strings, or already (partially) annotated paths. In the latter case, further annotation is limited to yet-unknown properties, and is potentially faster than initial annotation.

### Recognized path properties

“**action**” label of the action that triggered the path annotation

“**annexkey**” annex key for the content of a file

“**logger**” logger for reporting a message

“**message**” message (plus possible tsring expansion arguments)

“**orig\_request**” original input by which a path was determined

“**parentds**” path of dataset containing the annotated path (superdataset for subdatasets)

“**path**” absolute path that is annotated

“**process\_content**” flag that content underneath the path is to be processed

“**process\_updated\_only**” flag that only known dataset components are to be processed

“**raw\_input**” flag whether this path was given as raw (non-annotated) input

“**refds**” path of a reference/base dataset the annotated path is part of

“**registered\_subds**” flag whether a dataset is known to be a true subdataset of *parentds*

“**revision**” the recorded commit for a subdataset in a superdataset

“**revision\_descr**” a human-readable description of *revision*

“**source\_url**” URL a dataset was installed from

“**staged**” flag whether a path is known to be “staged” in its containing dataset

“**state**” state indicator for a path in its containing dataset (clean, modified, absent (also for files), conflict)

“**status**” action result status (ok, notneeded, impossible, error)

“**type**” nature of the path (file, directory, dataset)

“**url**” registered URL for a subdataset in a superdataset

In the case of enabled modification detection the results may contain additional properties regarding the nature of the modification. See the documentation of the *diff* command for details.

### Parameters

- **path** (*sequence of str or None, optional*) – path to be annotated. [Default: None]

- **dataset** (*Dataset or None, optional*) – an optional reference/base dataset for the paths. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **action** (*str or None, optional*) – an “action” property value to include in the path annotation. [Default: None]
- **unavailable\_path\_status** (*str or None, optional*) – a “status” property value to include in the annotation for paths that are underneath a dataset, but do not exist on the filesystem. [Default: ‘’]
- **unavailable\_path\_msg** (*str or None, optional*) – a “message” property value to include in the annotation for paths that are underneath a dataset, but do not exist on the filesystem. [Default: None]
- **nondataset\_path\_status** (*str or None, optional*) – a “status” property value to include in the annotation for paths that are not underneath any dataset. [Default: ‘error’]
- **force\_parentds\_discovery** (*bool, optional*) – Flag to disable reports of parent dataset information for any path, in particular dataset root paths. Disabling saves on command run time, if this information is not needed. [Default: True]
- **force\_subds\_discovery** (*bool, optional*) – Flag to disable reporting type=‘dataset’ for subdatasets, even when they are not installed, or their mount point directory doesn’t exist. Disabling saves on command run time, if this information is not needed. [Default: True]
- **force\_no\_revision\_change\_discovery** (*bool, optional*) – Flag to disable discovery of changes which were not yet committed. Disabling saves on command run time, if this information is not needed. [Default: True]
- **force\_untracked\_discovery** (*bool, optional*) – Flag to disable discovery of untracked changes. Disabling saves on command run time, if this information is not needed. [Default: True]
- **modified** (*str or bool or None, optional*) – comparison reference specification for modification detection. This can be (mostly) anything that *git diff* understands (commit, treeish, tag, etc). See the documentation of *datalad diff –revision* for details. Unmodified paths will not be annotated. If a requested path was not modified but some content underneath it was, then the request is replaced by the modified paths and those are annotated instead. This option can be used with *True* as an argument to test against changes that have been made, but have not yet been staged for a commit. [Default: None]
- **on\_failure** (*{‘ignore’, ‘continue’, ‘stop’}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an *IncompleteResultsError* that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a *ValueError* exception is raised. If the given callable supports

*\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]

- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'}* or *None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'}* or *callable* or *None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}*, *optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.clean

`datalad.api.clean` (*dataset=None, what=None, recursive=False, recursion\_limit=None*)  
Clean up after DataLad (possible temporary files etc.)

Removes extracted temporary archives, etc.

### Examples

\$ `datalad clean`

#### Parameters

- **dataset** (*Dataset* or *None, optional*) – specify the dataset to perform the clean operation on. If no dataset is given, an attempt is made to identify the dataset in current working directory. [Default: None]
- **what** – What to clean. If none specified – all known targets are cleaned. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int* or *None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}*, *optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an *IncompleteResultsError* that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable* or *None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a *ValueError* exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]



- **result\_renderer** (`{'default', 'json', 'json_pp', 'tailored'}` or `None, optional`) – format of return value rendering on stdout. [Default: `None`]
- **result\_xfm** (`{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'}` or `callable` or `None, optional`) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: `None`]
- **return\_type** (`{'generator', 'list', 'item-or-list'}`, `optional`) – return value behavior switch. If `'item-or-list'` a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: `'list'`]

## datalad.api.clone

`datalad.api.clone` (*source, path=None, dataset=None, description=None, reckless=None*)

Obtain a dataset (copy) from a URL or local directory

The purpose of this command is to obtain a new clone (copy) of a dataset and place it into a not-yet-existing or empty directory. As such `clone` provides a strict subset of the functionality offered by `install`. Only a single dataset can be obtained, and immediate recursive installation of subdatasets is not supported. However, once a (super)dataset is installed via `clone`, any content, including subdatasets can be obtained by a subsequent `get` command.

Primary differences over a direct `git clone` call are 1) the automatic initialization of a dataset annex (pure Git repositories are equally supported); 2) automatic registration of the newly obtained dataset as a subdataset (submodule), if a parent dataset is specified; and 3) support for additional resource identifiers (DataLad resource identifiers as used on `datasets.datalad.org`, and RIA store URLs as used for `store.datalad.org`; see examples); and 4) automatic configurable generation of alternative access URL for common cases (such as appending `'.git'` to the URL in case the accessing the base URL failed).

By default, the command returns a single Dataset instance for an installed dataset, regardless of whether it was newly installed (`'ok'` result), or found already installed from the specified source (`'notneeded'` result).

**See also:**

**handbook:3-001** (<http://handbook.datalad.org/symbols>) More information on Remote Indexed Archive (RIA) stores

## Examples

Install a dataset from Github into the current directory:

```
> clone(source='https://github.com/datalad-datasets/longnow-podcasts.git')
```

Install a dataset into a specific directory:

```
> clone(source='https://github.com/datalad-datasets/longnow-podcasts.git',
        path='myfavpodcasts')
```

Install a dataset as a subdataset into the current dataset:

```
> clone(dataset='.',
        source='https://github.com/datalad-datasets/longnow-podcasts.git')
```

Install the main superdataset from datasets.datalad.org:

```
> clone(source='///')
```

Install a dataset identified by its ID from store.datalad.org:

```
> clone(source='ria+http://store.datalad.org#76b6ca66-36b1-11ea-a2e6-f0d5bf7b5561
↪')
```

### Parameters

- **source** (*str or None*) – URL, DataLad resource identifier, local path or instance of dataset to be cloned.
- **path** – path to clone into. If no *path* is provided a destination path will be derived from a source URL similar to git clone. [Default: None]
- **dataset** (*Dataset or None, optional*) – (parent) dataset to clone into. If given, the newly cloned dataset is registered as a subdataset of the parent. Also, if given, relative paths are interpreted as being relative to the parent dataset, and not relative to the working directory. [Default: None]
- **description** (*str or None, optional*) – short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. [Default: None]
- **reckless** (*{None, True, False, 'auto', 'ephemeral'}, optional*) – Set up the dataset to be able to obtain content in the cheapest/fastest possible way, even if this poses a potential risk the data integrity (‘auto’: hardlink files from a local clone of the dataset, ‘ephemeral’: symlink annex to origin’s annex and discard local availability info via git-annex-dead ‘here’. Please note, that with a symlinked annex you share the annex with origin w/o git-annex knowing this. In case of a change in origin you need to update the clone before you’re able to save new content on your end.). Use with care, and limit to “read-only” use cases. With this flag the installed dataset will be marked as untrusted. The reckless mode is stored in a dataset’s local configuration under ‘datalad.clone.reckless’, and will be inherited to any of its subdatasets. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]

- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

### datalad.api.create\_test\_dataset

`datalad.api.create_test_dataset` (*path=None, spec=None, seed=None*)

Create test (meta-)dataset.

#### Parameters

- **path** (*str or None, optional*) – path/name where to create (if specified, must not exist). [Default: None]
- **spec** (*str or None, optional*) – spec for hierarchy, defined as a min-max (min could be omitted to assume 0) defining how many (random number from min to max) of sub-datasets to generate at any given level of the hierarchy. Each level separated from each other with /. Example: 1-3/-2 would generate from 1 to 3 subdatasets at the top level, and up to two within those at the 2nd level. [Default: None]
- **seed** (*int or None, optional*) – seed for rng. [Default: None]

### datalad.api.diff

`datalad.api.diff` (*path=None, fr='HEAD', to=None, dataset=None, annex=None, untracked='normal', recursive=False, recursion\_limit=None*)

Report differences between two states of a dataset (hierarchy)

The two to-be-compared states are given via the `–from` and `–to` options. These state identifiers are evaluated in the context of the (specified or detected) dataset. In the case of a recursive report on a dataset hierarchy, corresponding state pairs for any subdataset are determined from the subdataset record in the respective superdataset. Only changes recorded in a subdataset between these two states are reported, and so on.

Any paths given as additional arguments will be used to constrain the difference report. As with Git’s diff, it will not result in an error when a path is specified that does not exist on the filesystem.

Reports are very similar to those of the *status* command, with the distinguished content types and states being identical.

#### Examples

Show unsaved changes in a dataset:

```
> diff()
```

Compare a previous dataset state identified by shasum against current worktree:

```
> diff(fr='SHASUM')
```

Compare two branches against each other:

```
> diff(fr='branch1', to='branch2')
```

Show unsaved changes in the dataset and potential subdatasets:

```
> diff(recursive=True)
```

Show unsaved changes made to a particular file:

```
> diff(path='path/to/file')
```

### Parameters

- **path** (*sequence of str or None, optional*) – path to constrain the report to. [Default: None]
- **fr** (*str, optional*) – original state to compare to, as given by any identifier that Git understands. [Default: 'HEAD']
- **to** (*str or None, optional*) – state to compare against the original state, as given by any identifier that Git understands. If none is specified, the state of the working tree will be compared. [Default: None]
- **dataset** (*Dataset or None, optional*) – specify the dataset to query. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. [Default: None]
- **annex** (*{None, 'basic', 'availability', 'all'}, optional*) – Switch whether to include information on the annex content of individual files in the status report, such as recorded file size. By default no annex information is reported (faster). Three report modes are available: basic information like file size and key name ('basic'); additionally test whether file content is present in the local annex ('availability'; requires one or two additional file system stat calls, but does not call git-annex), this will add the result properties 'has\_content' (boolean flag) and 'objloc' (absolute path to an existing annex object file); or 'all' which will report all available information (presently identical to 'availability'). [Default: None]
- **untracked** (*{'no', 'normal', 'all'}, optional*) – If and how untracked content is reported when comparing a revision to the state of the working tree. 'no': no untracked content is reported; 'normal': untracked files and entire untracked directories are reported as such; 'all': report individual files even in fully untracked directories. [Default: 'normal']
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: 'ignore' any failure is reported, but does not cause an exception; 'continue' if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; 'stop': processing will stop on first failure and an exception is raised. A failure is any result with status 'impossible' or 'error'. Raised exception

is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: 'continue']

- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to `False` or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: `None`]
- **result\_renderer** (`{'default', 'json', 'json_pp', 'tailored'}` or *None, optional*) – format of return value rendering on stdout. [Default: `None`]
- **result\_xfm** (`{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'}` or *callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: `None`]
- **return\_type** (`{'generator', 'list', 'item-or-list'}`, *optional*) – return value behavior switch. If 'item-or-list' a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: 'list']

## datalad.api.download\_url

`datalad.api.download_url(urls, dataset=None, path=None, overwrite=False, archive=False, save=True, message=None)`

Download content

It allows for a uniform download interface to various supported URL schemes, re-using or asking for authentication details maintained by datalad.

### Examples

Download files from an http and S3 URL:

```
> download_url(urls=['http://example.com/file.dat', 's3://bucket/file2.dat'])
```

Download a file to a path and provide a commit message:

```
> download_url(urls='s3://bucket/file2.dat', message='added a file', path='myfile.
↳dat')
```

### Parameters

- **urls** (*non-empty sequence of str*) – URL(s) to be downloaded.
- **dataset** (`Dataset or None, optional`) – specify the dataset to add files to. If no dataset is given, an attempt is made to identify the dataset based on the current working directory. Use `save=False` to prevent adding files to the dataset. [Default: `None`]
- **path** (*str or None, optional*) – target for download. If the path has a trailing separator, it is treated as a directory, and each specified URL is downloaded under that directory to a base name taken from the URL. Without a trailing separator, the value specifies the name of the downloaded file (file name extensions inferred from the URL may be added

to it, if they are not yet present) and only a single URL should be given. In both cases, leading directories will be created if needed. This argument defaults to the current directory. [Default: None]

- **overwrite** (*bool, optional*) – flag to overwrite it if target file exists. [Default: False]
- **archive** (*bool, optional*) – pass the downloaded files to `add_archive_content(..., delete=True)`. [Default: False]
- **save** (*bool, optional*) – by default all modifications to a dataset are immediately saved. Giving this option will disable this behavior. [Default: True]
- **message** (*str or None, optional*) – a description of the state or the changes made to a dataset. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an `IncompleteResultsError` that carries the result dictionaries of the failures in its `failed` attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a `ValueError` exception is raised. If the given callable supports `**kwargs` it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from `result_filter`, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. `None` is return in case of an empty list. [Default: ‘list’]

## datalad.api.ls

`datalad.api.ls` (*loc, recursive=False, fast=False, all\_=False, long\_=False, config\_file=None, list\_content=False, json=None*)

List summary information about URLs and dataset(s)

ATM only s3:// URLs and datasets are supported

## Examples

\$ `datalad ls s3://openfmri/tarballs/ds202 #` to list S3 bucket  
 \$ `datalad ls #` to list current dataset

**Parameters**

- **loc** (*sequence of str or None*) – URL or path to list, e.g. s3://...
- **recursive** (*bool, optional*) – recurse into subdirectories. [Default: False]
- **fast** (*bool, optional*) – only perform fast operations. Would be overridden by `-all`. [Default: False]
- **all** (*bool, optional*) – list all (versions of) entries, not e.g. only latest entries in case of S3. [Default: False]
- **long** (*bool, optional*) – list more information on entries (e.g. acl, urls in s3, annex sizes etc). [Default: False]
- **config\_file** (*str or None, optional*) – path to config file which could help the ‘ls’. E.g. for s3:// URLs could be some `~/s3cfg` file which would provide credentials. [Default: None]
- **list\_content** – list also the content or only first 10 bytes (`first10`), or md5 checksum of an entry. Might require expensive transfer and dump binary output to your screen. Do not enable unless you know what you are after. [Default: False]
- **json** – metadata json of dataset for creating web user interface. `display`: prints jsons to stdout or file: writes each subdir metadata to json file in subdir of dataset or `delete`: deletes all metadata json files in dataset. [Default: None]

**datalad.api.sshrun**

`datalad.api.sshrun` (*login, cmd, port=None, ipv4=False, ipv6=False, options=None, no\_stdin=False*)  
Run command on remote machines via SSH.

This is a replacement for a small part of the functionality of SSH. In addition to SSH alone, this command can make use of datalad’s SSH connection management. Its primary use case is to be used with Git as ‘core.sshCommand’ or via “GIT\_SSH\_COMMAND”.

Configure `datalad.ssh.identityfile` to pass a file to the ssh’s `-i` option.

**Parameters**

- **login** – [user@]hostname.
- **cmd** – command for remote execution.
- **port** – port to connect to on the remote host. [Default: None]
- **ipv4** (*bool, optional*) – use IPv4 addresses only. [Default: False]
- **ipv6** (*bool, optional*) – use IPv6 addresses only. [Default: False]
- **options** – configuration option passed to SSH. [Default: None]
- **no\_stdin** (*bool, optional*) – Do not connect stdin to the process. [Default: False]

## datalad.api.siblings

`datalad.api.siblings` (*action='query', dataset=None, name=None, url=None, pushurl=None, description=None, fetch=False, as\_common\_datasrc=None, publish\_depends=None, publish\_by\_default=None, annex\_wanted=None, annex\_required=None, annex\_group=None, annex\_groupwanted=None, inherit=False, get\_annex\_info=True, recursive=False, recursion\_limit=None*)

Manage sibling configuration

This command offers four different actions: ‘query’, ‘add’, ‘remove’, ‘configure’, ‘enable’. ‘query’ is the default action and can be used to obtain information about (all) known siblings. ‘add’ and ‘configure’ are highly similar actions, the only difference being that adding a sibling with a name that is already registered will fail, whereas re-configuring a (different) sibling under a known name will not be considered an error. ‘enable’ can be used to complete access configuration for non-Git sibling (aka git-annex special remotes). Lastly, the ‘remove’ action allows for the removal (or de-configuration) of a registered sibling.

For each sibling (added, configured, or queried) all known sibling properties are reported. This includes:

“**name**” Name of the sibling

“**path**” Absolute path of the dataset

“**url**” For regular siblings at minimum a “fetch” URL, possibly also a “pushurl”

Additionally, any further configuration will also be reported using a key that matches that in the Git configuration.

By default, sibling information is rendered as one line per sibling following this scheme:

```
<dataset_path>: <sibling_name>(<+|->) [<access_specification>]
```

where the + and - labels indicate the presence or absence of a remote data annex at a particular remote, and *access\_specification* contains either a URL and/or a type label for the sibling.

### Parameters

- **action** (*{'query', 'add', 'remove', 'configure', 'enable'} or None, optional*) – command action selection (see general documentation). [Default: ‘query’]
- **dataset** (*Dataset or None, optional*) – specify the dataset to configure. If no dataset is given, an attempt is made to identify the dataset based on the input and/or the current working directory. [Default: None]
- **name** (*str or None, optional*) – name of the sibling. For addition with path “URLs” and sibling removal this option is mandatory, otherwise the hostname part of a given URL is used as a default. This option can be used to limit ‘query’ to a specific sibling. [Default: None]
- **url** (*str or None, optional*) – the URL of or path to the dataset sibling named by *name*. For recursive operation it is required that a template string for building subdataset sibling URLs is given. List of currently available placeholders: %%NAME the name of the dataset, where slashes are replaced by dashes. [Default: None]
- **pushurl** (*str or None, optional*) – in case the *url* cannot be used to publish to the dataset sibling, this option specifies a URL to be used instead. If no *url* is given, *pushurl* serves as *url* as well. [Default: None]
- **description** (*str or None, optional*) – short description to use for a dataset location. Its primary purpose is to help humans to identify a dataset copy (e.g., “mike’s



dataset on lab server”). Note that when a dataset is published, this information becomes available on the remote side. [Default: None]

- **fetch** (*bool, optional*) – fetch the sibling after configuration. [Default: False]
- **as\_common\_datasrc** – configure the created sibling as a common data source of the dataset that can be automatically used by all consumers of the dataset (technical: git-annex auto-enabled special remote). [Default: None]
- **publish\_depends** (*list of str or None, optional*) – add a dependency such that the given existing sibling is always published prior to the new sibling. This equals setting a configuration item ‘remote.SIBLINGNAME.datalad-publish-depends’. Multiple dependencies can be given as a list of sibling names. [Default: None]
- **publish\_by\_default** (*list of str or None, optional*) – add a refspec to be published to this sibling by default if nothing specified. [Default: None]
- **annex\_wanted** (*str or None, optional*) – expression to specify ‘wanted’ content for the repository/sibling. See <https://git-annex.branchable.com/git-annex-wanted/> for more information. [Default: None]
- **annex\_required** (*str or None, optional*) – expression to specify ‘required’ content for the repository/sibling. See <https://git-annex.branchable.com/git-annex-required/> for more information. [Default: None]
- **annex\_group** (*str or None, optional*) – expression to specify a group for the repository. See <https://git-annex.branchable.com/git-annex-group/> for more information. [Default: None]
- **annex\_groupwanted** (*str or None, optional*) – expression for the group-wanted. Makes sense only if annex\_wanted=”groupwanted” and annex-group is given too. See <https://git-annex.branchable.com/git-annex-groupwanted/> for more information. [Default: None]
- **inherit** (*bool, optional*) – if sibling is missing, inherit settings (git config, git annex wanted/group/groupwanted) from its super-dataset. [Default: False]
- **get\_annex\_info** (*bool, optional*) – Whether to query all information about the annex configurations of siblings. Can be disabled if speed is a concern. [Default: True]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an IncompleteResultsError that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a ValueError exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]

- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'}* or *None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'}* or *callable* or *None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top-level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}*, *optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: ‘list’]

## datalad.api.subdatasets

`datalad.api.subdatasets` (*path=None, dataset=None, fulfilled=None, recursive=False, recursion\_limit=None, contains=None, bottomup=False, set\_property=None, delete\_property=None*)

Report subdatasets and their properties.

The following properties are reported (if possible) for each matching subdataset record.

“**name**” Name of the subdataset in the parent (often identical with the relative path in the parent dataset)

“**path**” Absolute path to the subdataset

“**parentds**” Absolute path to the parent dataset

“**gitshasum**” SHA1 of the subdataset commit recorded in the parent dataset

“**state**” Condition of the subdataset: ‘clean’, ‘modified’, ‘absent’, ‘conflict’ as reported by *git submodule*

“**gitmodule\_url**” URL of the subdataset recorded in the parent

“**gitmodule\_name**” Name of the subdataset recorded in the parent

“**gitmodule\_<label>**” Any additional configuration property on record.

Performance note: Property modification, requesting *bottomup* reporting order, or a particular numerical *recursion\_limit* implies an internal switch to an alternative query implementation for recursive query that is more flexible, but also notably slower (performs one call to Git per dataset versus a single call for all combined).

The following properties for subdatasets are recognized by DataLad (without the ‘gitmodule\_’ prefix that is used in the query results):

“**datalad-recursiveinstall**” If set to ‘skip’, the respective subdataset is skipped when DataLad is recursively installing its superdataset. However, the subdataset remains installable when explicitly requested, and no other features are impaired.

### Parameters

- **path** (*sequence of str* or *None, optional*) – path/name to query for subdatasets. Defaults to the current directory, or the entire dataset if called as a dataset method. [Default: None]
- **dataset** (*Dataset* or *None, optional*) – specify the dataset to query. If no dataset is given, an attempt is made to identify the dataset based on the input and/or the current working directory. [Default: None]

- **fulfilled** (*bool or None, optional*) – if given, must be a boolean flag indicating whether to report either only locally present or absent datasets. By default subdatasets are reported regardless of their status. [Default: None]
- **recursive** (*bool, optional*) – if set, recurse into potential subdataset. [Default: False]
- **recursion\_limit** (*int or None, optional*) – limit recursion into subdataset to the given number of levels. [Default: None]
- **contains** (*list of str or None, optional*) – limit report to the subdatasets containing the given path. If a root path of a subdataset is given the last reported dataset will be the subdataset itself. Can be a list with multiple paths, in which case datasets will be reported that contain any of the given paths. [Default: None]
- **bottomup** (*bool, optional*) – whether to report subdatasets in bottom-up order along each branch in the dataset tree, and not top-down. [Default: False]
- **set\_property** (*list of 2-item sequence of str or None, optional*) – Name and value of one or more subdataset properties to be set in the parent dataset’s .gitmodules file. The property name is case- insensitive, must start with a letter, and consist only of alphanumeric characters. The value can be a Python format() template string wrapped in ‘<>’ (e.g. ‘<{gitmodule\_name}>’). Supported keywords are any item reported in the result properties of this command, plus ‘refds\_relpath’ and ‘refds\_relname’: the relative path of a subdataset with respect to the base dataset of the command call, and, in the latter case, the same string with all directory separators replaced by dashes. [Default: None]
- **delete\_property** (*list of str or None, optional*) – Name of one or more subdataset properties to be removed from the parent dataset’s .gitmodules file. [Default: None]
- **on\_failure** (*{'ignore', 'continue', 'stop'}, optional*) – behavior to perform on failure: ‘ignore’ any failure is reported, but does not cause an exception; ‘continue’ if any failure occurs an exception will be raised at the end, but processing other actions will continue for as long as possible; ‘stop’: processing will stop on first failure and an exception is raised. A failure is any result with status ‘impossible’ or ‘error’. Raised exception is an IncompleteResultsError that carries the result dictionaries of the failures in its *failed* attribute. [Default: ‘continue’]
- **result\_filter** (*callable or None, optional*) – if given, each to-be-returned status dictionary is passed to this callable, and is only returned if the callable’s return value does not evaluate to False or a ValueError exception is raised. If the given callable supports *\*\*kwargs* it will additionally be passed the keyword arguments of the original API call. [Default: None]
- **result\_renderer** (*{'default', 'json', 'json\_pp', 'tailored'} or None, optional*) – format of return value rendering on stdout. [Default: None]
- **result\_xfm** (*{'datasets', 'successdatasets-or-none', 'paths', 'relpaths', 'metadata'} or callable or None, optional*) – if given, each to-be-returned result status dictionary is passed to this callable, and its return value becomes the result instead. This is different from *result\_filter*, as it can perform arbitrary transformation of the result value. This is mostly useful for top- level command invocations that need to provide the results in a particular format. Instead of a callable, a label for a pre-crafted result transformation can be given. [Default: None]
- **return\_type** (*{'generator', 'list', 'item-or-list'}, optional*) – return value behavior switch. If ‘item-or-list’ a single value is returned instead of a one-item

return value list, or a list in case of multiple return values. *None* is return in case of an empty list. [Default: 'list']

## Miscellaneous commands

---

`api.add_archive_content`(archive[, annex, Add content of an archive under git annex control. ...])

---

`api.test`([module, verbose, nocapture, pdb, stop]) Run internal DataLad (unit)tests.

---

### datalad.api.add\_archive\_content

```
datalad.api.add_archive_content (archive, annex=None, add_archive_leading_dir=False,
                                strip_leading_dirs=False, leading_dirs_depth=None,
                                leading_dirs_consider=None, use_current_dir=False,
                                delete=False, key=False, exclude=None, rename=None,
                                existing='fail', annex_options=None, copy=False, com-
                                mit=True, allow_dirty=False, stats=None, drop_after=False,
                                delete_after=False)
```

Add content of an archive under git annex control.

This results in the files within archive (which must be already under annex control itself) added under annex referencing original archive via custom special remotes mechanism

### Example

```
annex-repo$ datalad add-archive-content my_big_tarball.tar.gz
```

#### Parameters

- **archive** (*str*) – archive file or a key (if *key=True* specified).
- **annex** – annex instance to use. [Default: None]
- **add\_archive\_leading\_dir** (*bool, optional*) – flag to place extracted content under a directory which would correspond to archive name with suffix stripped. E.g. for archive *example.zip* its content will be extracted under a directory *example/*. [Default: False]
- **strip\_leading\_dirs** (*bool, optional*) – flag to move all files directories up, from how they were stored in an archive, if that one contained a number (possibly more than 1 down) single leading directories. [Default: False]
- **leading\_dirs\_depth** – maximal depth to strip leading directories to. If not specified (None), no limit. [Default: None]
- **leading\_dirs\_consider** (*list of str or None, optional*) – regular expression(s) for directories to consider to strip away. [Default: None]
- **use\_current\_dir** (*bool, optional*) – flag to extract archive under the current directory, not the directory where archive is located. Note that it will be of no effect if *key=True* is given. [Default: False]
- **delete** (*bool, optional*) – flag to delete original archive from the filesystem/git in current tree. Note that it will be of no effect if *key=True* is given. [Default: False]
- **key** (*bool, optional*) – flag to signal if provided archive is not actually a filename on its own but an annex key. [Default: False]

- **exclude** (*list of str or None, optional*) – regular expressions for filenames which to exclude from being added to annex. Applied after `–rename` if that one is specified. For exact matching, use anchoring. [Default: None]
- **rename** (*list of str or None, optional*) – regular expressions to rename files before being added under git. First letter defines how to split provided string into two parts: Python regular expression (with groups), and replacement string. [Default: None]
- **existing** – what operation to perform a file from archive tries to overwrite an existing file with the same name. ‘fail’ (default) leads to RuntimeError exception. ‘overwrite’ silently replaces existing file. ‘archive-suffix’ instructs to add a suffix (prefixed with a ‘-’) matching archive name from which file gets extracted, and if that one present, ‘numeric-suffix’ is in effect in addition, when incremental numeric suffix (prefixed with a ‘.’) is added until no name collision is longer detected. [Default: ‘fail’]
- **annex\_options** (*str or None, optional*) – additional options to pass to git-annex. [Default: None]
- **copy** (*bool, optional*) – flag to copy the content of the archive instead of moving. [Default: False]
- **commit** (*bool, optional*) – flag to not commit upon completion. [Default: True]
- **allow\_dirty** (*bool, optional*) – flag that operating on a dirty repository (uncommitted or untracked content) is ok. [Default: False]
- **stats** – ActivityStats instance for global tracking. [Default: None]
- **drop\_after** (*bool, optional*) – drop extracted files after adding to annex. [Default: False]
- **delete\_after** (*bool, optional*) – extract under a temporary directory, git-annex add, and delete after. To be used to “index” files within annex without actually creating corresponding files under git. Note that *annex dropunused* would later remove that load. [Default: False]

### Returns

**Return type** annex

## datalad.api.test

`datalad.api.test` (*module=None, verbose=False, nocapture=False, pdb=False, stop=False*)

Run internal DataLad (unit)tests.

This can be used to verify correct operation on the system. It is just a thin wrapper around a call to nose, so number of exposed options is minimal

### Parameters

- **module** – test name(s), by default all tests of DataLad core and any installed extensions are executed. [Default: None]
- **verbose** (*bool, optional*) – be verbose - list test names. [Default: False]
- **nocapture** (*bool, optional*) – do not capture stdout. [Default: False]
- **pdb** (*bool, optional*) – drop into debugger on failures or errors. [Default: False]
- **stop** (*bool, optional*) – stop running tests after the first error or failure. [Default: False]

## Plugins

DataLad can be customized by plugins. The following plugins are shipped with DataLad.

|                                 |   |
|---------------------------------|---|
| <code>add_readme</code>         | add a README file to a dataset                          |
| <code>addurls</code>            | Create and update a dataset from a list of URLs.        |
| <code>check_dates</code>        | Extension for checking dates within repositories.       |
| <code>export_archive</code>     | export a dataset as a compressed TAR/ZIP archive        |
| <code>export_to_figshare</code> | export a dataset as a TAR/ZIP archive to figshare       |
| <code>no_annex</code>           | configure which dataset parts to never put in the annex |
| <code>wtf</code>                | provide information about this DataLad installation     |

## Support functionality

|                                      |   |
|--------------------------------------|---|
| <code>auto</code>                    | Proxy basic file operations (e.g.   |
| <code>cmd</code>                     | Wrapper for command and function calls, allowing for dry runs and output handling |
| <code>consts</code>                  | constants for datalad   |
| <code>log</code>                     |   |
| <code>utils</code>                   |   |
| <code>version</code>                 | Defines version to be imported in the module and obtained from setup.py           |
| <code>support.gitrepo</code>         | Interface to Git via GitPython  |
| <code>support.annexrepo</code>       | Interface to git-annex by Joey Hess.  |
| <code>support.archives</code>        | Various handlers/functionality for different types of files (e.g.                 |
| <code>support.configparserinc</code> |   |
| <code>customremotes.main</code>      |   |
| <code>customremotes.base</code>      | Base classes to custom git-annex remotes (e.g.                                    |
| <code>customremotes.archives</code>  | Custom remote to support getting the load from archives present under annex       |

### datalad.auto

Proxy basic file operations (e.g. open) to auto-obtain files upon I/O

**class** `datalad.auto.AutomagicIO` (*autoget=True, activate=False, check\_once=False*)

Bases: `object`

Class to proxy commonly used API for accessing files so they get automatically fetched

Currently supports builtin `open()` and `h5py.File` when those are read

**activate** ()

**active**

**autoget**

**deactivate** ()

**datalad.cmd**

Wrapper for command and function calls, allowing for dry runs and output handling

**class** `datalad.cmd.BatchedCommand` (*cmd*, *path=None*, *output\_proc=None*)

Bases: `datalad.cmd.SafeDelCloseMixin`

Container for a process which would allow for persistent communication

**close** (*return\_stderr=False*)

Close communication and wait for process to terminate

**Returns** stderr output if *return\_stderr* and stderr file was there. None otherwise

**Return type** str

**procl** (*arg*)

Same as `__call__`, but only takes a single command argument

and returns a single result.

**yield\_** (*cmds*)

Same as `__call__`, but requires *cmds* to be an iterable

and yields results for each item.

**class** `datalad.cmd.GitRunner` (*\*args*, *\*\*kwargs*)

Bases: `datalad.cmd.Runner`

Runner to be used to run git and git annex commands

Overloads the runner class to check & update `GIT_DIR` and `GIT_WORK_TREE` environment variables set to the absolute path if is defined and is relative path

**static** `get_git_envIRON_adjusted` (*env=None*)

Replaces `GIT_DIR` and `GIT_WORK_TREE` with absolute paths if relative path and defined

**run** (*cmd*, *env=None*, *\*args*, *\*\*kwargs*)

Runs the command *cmd* using shell.

In case of dry-mode *cmd* is just added to *commands* and it is actually executed otherwise. Allows for separately logging stdout and stderr or streaming it to system's stdout or stderr respectively.

**Note: Using a string as *cmd* and *shell=True* allows for piping**, multiple commands, etc., but that implies `split_cmdline()` is not used. This is considered to be a security hazard. So be careful with input.

**Parameters**

- **cmd** (*str*, *list*) – String (or list) defining the command call. No shell is used if *cmd* is specified as a list
- **log\_stdout** (*bool*, *optional*) – If True, stdout is logged. Goes to `sys.stdout` otherwise.
- **log\_stderr** (*bool*, *optional*) – If True, stderr is logged. Goes to `sys.stderr` otherwise.
- **log\_online** (*bool*, *optional*) – Whether to log as output comes in. Setting to True is preferable for running user-invoked actions to provide timely output
- **expect\_stderr** (*bool*, *optional*) – Normally, having stderr output is a signal of a problem and thus it gets logged at level 11. But some utilities, e.g. `wget`, use stderr for their progress output. Whenever such output is expected, set it to True and output will be

logged at level 9 unless exit status is non-0 (in non-online mode only, in online – would log at 9)

- **expect\_fail** (*bool, optional*) – Normally, if command exits with non-0 status, it is considered an error and logged at level 11 (above DEBUG). But if the call intended for checking routine, such messages are usually not needed, thus it will be logged at level 9.
- **cwd** (*string, optional*) – Directory under which run the command (passed to Popen)
- **env** (*string, optional*) – Custom environment to pass
- **shell** (*bool, optional*) – Run command in a shell. If not specified, then it runs in a shell only if command is specified as a string (not a list)
- **stdin** (*file descriptor*) – input stream to connect to stdin of the process.

### Returns

**Return type** (stdout, stderr) - bytes!

**Raises** `CommandError` – if command’s exitcode wasn’t 0 or None. `exitcode` is passed to `CommandError`’s `code`-field. Command’s stdout and stderr are stored in `CommandError`’s `stdout` and `stderr` fields respectively.

**class** `datalad.cmd.KillOutput` (*done\_future*)

Bases: `datalad.cmd.WitlessProtocol`

WitlessProtocol that swallows stdout/stderr of a subprocess

**pipe\_data\_received** (*fd, data*)

Called when the subprocess writes data into stdout/stderr pipe.

`fd` is int file descriptor. `data` is bytes object.

**proc\_err** = True

**proc\_out** = True

**class** `datalad.cmd.NoCapture` (*done\_future*)

Bases: `datalad.cmd.WitlessProtocol`

WitlessProtocol that captures no subprocess output

As this is identical with the behavior of the `WitlessProtocol` base class, this class is merely a more readable convenience alias.

**class** `datalad.cmd.Runner` (*cwd=None, env=None, protocol=None, log\_outputs=None*)

Bases: `object`

Provides a wrapper for calling functions and commands.

An object of this class provides a methods that calls shell commands or python functions, allowing for protocoling the calls and output handling.

Outputs (stdout and stderr) can be either logged or streamed to system’s stdout/stderr during execution. This can be enabled or disabled for both of them independently. Additionally, a protocol object can be a used with the `Runner`. Such a protocol has to implement `datalad.support.protocol.ProtocolInterface`, is able to record calls and allows for dry runs.

**call** (*f, \*args, \*\*kwargs*)

Helper to unify collection of logging all “dry” actions.

Calls `f` if `Runner`-object is not in dry-mode. Adds `f` along with its arguments to `commands` otherwise.



**Parameters** *f* (*callable*) –**commands****cwd****dry****env**

**log** (*msg*, *\*args*, *\*\*kwargs*)  
log helper

Logs at level 9 by default and adds “Protocol:”-prefix in order to log the used protocol.

**log\_cwd****log\_env****log\_outputs****log\_stdin****protocol**

**run** (*cmd*, *log\_stdout=True*, *log\_stderr=True*, *log\_online=False*, *expect\_stderr=False*, *expect\_fail=False*, *cwd=None*, *env=None*, *shell=None*, *stdin=None*)  
Runs the command *cmd* using shell.

In case of dry-mode *cmd* is just added to *commands* and it is actually executed otherwise. Allows for separately logging stdout and stderr or streaming it to system’s stdout or stderr respectively.

**Note: Using a string as *cmd* and *shell=True* allows for piping**, multiple commands, etc., but that implies *split\_cmdline()* is not used. This is considered to be a security hazard. So be careful with input.

**Parameters**

- **cmd** (*str*, *list*) – String (or list) defining the command call. No shell is used if *cmd* is specified as a list
- **log\_stdout** (*bool*, *optional*) – If True, stdout is logged. Goes to *sys.stdout* otherwise.
- **log\_stderr** (*bool*, *optional*) – If True, stderr is logged. Goes to *sys.stderr* otherwise.
- **log\_online** (*bool*, *optional*) – Whether to log as output comes in. Setting to True is preferable for running user-invoked actions to provide timely output
- **expect\_stderr** (*bool*, *optional*) – Normally, having stderr output is a signal of a problem and thus it gets logged at level 11. But some utilities, e.g. *wget*, use stderr for their progress output. Whenever such output is expected, set it to True and output will be logged at level 9 unless exit status is non-0 (in non-online mode only, in online – would log at 9)
- **expect\_fail** (*bool*, *optional*) – Normally, if command exits with non-0 status, it is considered an error and logged at level 11 (above DEBUG). But if the call intended for checking routine, such messages are usually not needed, thus it will be logged at level 9.
- **cwd** (*string*, *optional*) – Directory under which run the command (passed to *Popen*)
- **env** (*string*, *optional*) – Custom environment to pass

- **shell** (*bool, optional*) – Run command in a shell. If not specified, then it runs in a shell only if command is specified as a string (not a list)
- **stdin** (*file descriptor*) – input stream to connect to stdin of the process.

**Returns**

**Return type** (stdout, stderr) - bytes!

**Raises** `CommandError` – if command’s exitcode wasn’t 0 or None. `exitcode` is passed to `CommandError`’s `code`-field. Command’s stdout and stderr are stored in `CommandError`’s `stdout` and `stderr` fields respectively.

**class** `datalad.cmd.SafeDelCloseMixin`

Bases: `object`

A helper class to use where `__del__` would call `.close()` which might fail if “too late in GC game”

**class** `datalad.cmd.StdErrCapture` (*done\_future*)

Bases: `datalad.cmd.WitlessProtocol`

WitlessProtocol that only captures and returns stderr of a subprocess

**proc\_err = True**

**class** `datalad.cmd.StdOutCapture` (*done\_future*)

Bases: `datalad.cmd.WitlessProtocol`

WitlessProtocol that only captures and returns stdout of a subprocess

**proc\_out = True**

**class** `datalad.cmd.StdOutErrCapture` (*done\_future*)

Bases: `datalad.cmd.WitlessProtocol`

WitlessProtocol that captures and returns stdout/stderr of a subprocess

**proc\_err = True**

**proc\_out = True**

**class** `datalad.cmd.WitlessProtocol` (*done\_future*)

Bases: `asyncio.protocols.SubprocessProtocol`

Subprocess communication protocol base class for `run_async_cmd`

This class implements basic subprocess output handling. Derived classes like `StdOutCapture` should be used for subprocess communication that need to capture and return output. In particular, the `pipe_data_received()` method can be overwritten to implement “online” processing of process output.

This class defines a default return value setup that causes `run_async_cmd()` to return a 2-tuple with the subprocess’s exit code and a list with bytestrings of all captured output streams.

**FD\_NAMES = ['stdin', 'stdout', 'stderr']**

**connection\_made** (*transport*)

Called when a connection is made.

The argument is the transport representing the pipe connection. To receive data, wait for `data_received()` calls. When the connection is closed, `connection_lost()` is called.

**pipe\_data\_received** (*fd, data*)

Called when the subprocess writes data into stdout/stderr pipe.

`fd` is int file descriptor. `data` is bytes object.

`proc_err = None`

`proc_out = None`

`process_exited()`

Called when subprocess has exited.

**class** `datalad.cmd.WitlessRunner` (*cwd=None, env=None*)

Bases: `object`

Minimal Runner with support for online command output processing

It aims to be as simple as possible, providing only essential functionality.

**cwd**

**env**

**run** (*cmd, protocol=None, stdin=None*)

Execute a command and communicate with it.

#### Parameters

- **cmd** (*list*) – Sequence of program arguments. Passing a single string means that it is simply the name of the program, no complex shell commands are supported.
- **protocol** (`WitlessProtocol`, *optional*) – Protocol class handling interaction with the running process (e.g. output capture). A number of pre-crafted classes are provided (e.g. `KillOutput`, `NoCapture`, `GitProgress`).
- **stdin** (*byte stream, optional*) – File descriptor like, used as stdin for the process. Passed verbatim to `subprocess.Popen()`.

**Returns** Unicode string with the cumulative standard output and error of the process.

**Return type** `stdout, stderr`

#### Raises

- `CommandError` – On execution failure (non-zero exit code) this exception is raised which provides the command (`cmd`), `stdout`, `stderr`, exit code (`status`), and a message identifying the failed command, as properties.
- `FileNotFoundError` – When a given executable does not exist.

`datalad.cmd.readline_rstripped` (*stdout*)

`datalad.cmd.run_async_cmd` (*loop, cmd, protocol, stdin, \*\*kwargs*)

Run a command in a subprocess managed by `asyncio`

This implementation has been inspired by <https://pymotw.com/3/asyncio/subprocesses.html>

#### Parameters

- **loop** (`asyncio.AbstractEventLoop`) – `asyncio` event loop instance. Must support subprocesses on the target platform.
- **cmd** (*list*) – Command to be executed, passed to `subprocess_exec`.
- **protocol** (`WitlessProtocol`) – Protocol class to be instantiated for managing communication with the subprocess.
- **stdin** (*file-like or None*) – Passed to the subprocess as its standard input.
- **kwargs** (Pass to `subprocess_exec`, will typically be parameters) – supported by `subprocess.Popen`.

**Returns** The nature of the return value is determined by the given protocol class.

**Return type** undefined

`datalad.cmd.run_gitcommand_on_file_list_chunks` (*func*, *cmd*, *files*, *\*args*, *\*\*kwargs*)

Run a git command multiple times if *files* is too long

**Parameters**

- **func** (*callable*) – Typically a `Runner.run` variant. Assumed to return a 2-tuple with `stdout` and `stderr` as strings.
- **cmd** (*list*) – Base Git command argument list, to be amended with ‘-’, followed by a file list chunk.
- **files** (*list*) – List of files.
- **kwargs** (*args,*) – Passed to *func*

**Returns** Concatenated `stdout` and `stderr`.

**Return type** `str`, `str`

## datalad.consts

constants for datalad

## datalad.log

**class** `datalad.log.ColorFormatter` (*use\_color=None*, *log\_name=False*, *log\_pid=False*)

Bases: `logging.Formatter`

**format** (*record*)

Format the specified record as text.

The record’s attribute dictionary is used as the operand to a string formatting operation which yields the returned string. Before formatting the dictionary, a couple of preparatory steps are carried out. The message attribute of the record is computed using `LogRecord.getMessage()`. If the formatting string uses the time (as determined by a call to `usesTime()`, `formatTime()` is called to format the event time. If there is exception information, it is formatted using `formatException()` and appended to the message.

## datalad.utils

**class** `datalad.utils.ArgSpecFake` (*args*, *varargs*, *keywords*, *defaults*)

Bases: `tuple`

**args**

Alias for field number 0

**defaults**

Alias for field number 3

**keywords**

Alias for field number 2

**varargs**

Alias for field number 1

**class** `datalad.utils.File` (*name*, *executable=False*)

Bases: `object`

Helper for a file entry in the `create_tree/@with_tree`

It allows to define additional settings for entries

**class** `datalad.utils.SequenceFormatter` (*separator=' '*, *element\_formatter=<string.Formatter object>*, *\*args*, *\*\*kwargs*)

Bases: `string.Formatter`

`string.Formatter` subclass with special behavior for sequences.

This class delegates formatting of individual elements to another formatter object. Non-list objects are formatted by calling the delegate formatter's "format\_field" method. List-like objects (list, tuple, set, frozenset) are formatted by formatting each element of the list according to the specified format spec using the delegate formatter and then joining the resulting strings with a separator (space by default).

**format\_element** (*elem*, *format\_spec*)

Format a single element

For sequences, this is called once for each element in a sequence. For anything else, it is called on the entire object. It is intended to be overridden in subclasses.

**format\_field** (*value*, *format\_spec*)

`datalad.utils.all_same` (*items*)

Quick check if all items are the same.

Identical to a check like `len(set(items)) == 1` but should be more efficient while working on generators, since would return `False` as soon as any difference detected thus possibly avoiding unnecessary evaluations

`datalad.utils.any_re_search` (*regexes*, *value*)

Return if any of *regexes* (list or str) searches successfully for *value*

`datalad.utils.as_unicode` (*val*, *cast\_types=<class 'object'>*)

Given an arbitrary value, would try to obtain unicode value of it

For unicode it would return original value, for python2 str or python3 bytes it would use `ensure_unicode`, for `None` - an empty (unicode) string, and for any other type (see *cast\_types*) - would apply the unicode constructor. If *value* is not an instance of *cast\_types*, `TypeError` is thrown

**Parameters** `cast_types` (*type*) – Which types to cast to unicode by providing to constructor

`datalad.utils.assert_no_open_files` (*\*args*, *\*\*kwargs*)

`datalad.utils.assure_bool` (*s*)

Convert value into boolean following convention for strings

to recognize on, True, yes as `True`, off, False, no as `False`

`datalad.utils.assure_bytes` (*s*, *encoding='utf-8'*)

Convert/encode unicode string to bytes.

If *s* isn't a string, return it as is.

**Parameters** `encoding` (*str*, *optional*) – Encoding to use. "utf-8" is the default

`datalad.utils.assure_dict_from_str` (*s*, *\*\*kwargs*)

Given a multiline string with key=value items convert it to a dictionary

**Parameters**

- `s` (*str* or *dict*) –
- **None if input `s` is empty** (*Returns*) –

`datalad.utils.assure_dir(*args)`

Make sure directory exists.

Joins the list of arguments to an os-specific path to the desired directory and creates it, if it not exists yet.

`datalad.utils.assure_iter(s, cls, copy=False, iterate=True)`

Given not a list, would place it into a list. If None - empty list is returned

#### Parameters

- **s** (*list or anything*) –
- **cls** (*class*) – Which iterable class to ensure
- **copy** (*bool, optional*) – If correct iterable is passed, it would generate its shallow copy
- **iterate** (*bool, optional*) – If it is not a list, but something iterable (but not a str) iterate over it.

`datalad.utils.assure_list(s, copy=False, iterate=True)`

Given not a list, would place it into a list. If None - empty list is returned

#### Parameters

- **s** (*list or anything*) –
- **copy** (*bool, optional*) – If list is passed, it would generate a shallow copy of the list
- **iterate** (*bool, optional*) – If it is not a list, but something iterable (but not a str) iterate over it.

`datalad.utils.assure_list_from_str(s, sep='\n')`

Given a multiline string convert it to a list of return None if empty

**Parameters** **s** (*str or list*) –

`datalad.utils.assure_tuple_or_list(obj)`

Given an object, wrap into a tuple if not list or tuple

`datalad.utils.assure_unicode(s, encoding=None, confidence=None)`

Convert/decode bytestring to unicode.

If *s* isn't a bytestring, return it as is.

#### Parameters

- **encoding** (*str, optional*) – Encoding to use. If None, “utf-8” is tried, and then if not a valid UTF-8, encoding will be guessed
- **confidence** (*float, optional*) – A value between 0 and 1, so if guessing of encoding is of lower than specified confidence, ValueError is raised

`datalad.utils.auto_repr(cls)`

Decorator for a class to assign it an automagic quick and dirty `__repr__`

It uses public class attributes to prepare repr of a class

Original idea: <http://stackoverflow.com/a/27799004/1265472>

`datalad.utils.better_wraps(to_be_wrapped)`

Decorator to replace `functools.wraps`

This is based on `wrapt` instead of `functools` and in opposition to `wraps` preserves the correct signature of the decorated function. It is written with the intention to replace the use of `wraps` without any need to rewrite the actual decorators.

`datalad.utils.bytes2human` (*n*, *format*='%(value).1f%(symbol)sB')

Convert *n* bytes into a human readable string based on *format*. symbols can be either “customary”, “customary\_ext”, “iec” or “iec\_ext”, see: <http://goo.gl/kTQMs>

```
>>> from datalad.utils import bytes2human
>>> bytes2human(1)
'1.0 B'
>>> bytes2human(1024)
'1.0 KB'
>>> bytes2human(1048576)
'1.0 MB'
>>> bytes2human(1099511627776127398123789121)
'909.5 YB'
```

```
>>> bytes2human(10000, "%(value).1f %(symbol)s/sec")
'9.8 K/sec'
```

```
>>> # precision can be adjusted by playing with %f operator
>>> bytes2human(10000, format="%(value).5f %(symbol)s")
'9.76562 K'
```

Taken from: <http://goo.gl/kTQMs> and subsequently simplified Original Author: Giampaolo Rodola' <g.rodola [AT] gmail [DOT] com> License: MIT

`datalad.utils.check_symlink_capability` (*path*, *target*)

helper similar to `datalad.tests.utils.has_symlink_capability`

However, for use in a datalad command context, we shouldn't assume to be able to write to tmpfile and also not import a whole lot from datalad's test machinery. Finally, we want to know, whether we can create a symlink at a specific location, not just somewhere. Therefore use arbitrary path to test-build a symlink and delete afterwards. Suitable location can therefore be determined by high lever code.

#### Parameters

- **path** (*Path*) –
- **target** (*Path*) –

#### Returns

**Return type** bool

**class** `datalad.utils.chpwd` (*path*, *makedirs*=False, *logsuffix*='')

Bases: object

Wrapper around `os.chdir` which also adjusts `environ['PWD']`

The reason is that otherwise PWD is simply inherited from the shell and we have no ability to assess directory path without dereferencing symlinks.

If used as a context manager it allows to temporarily change directory to the given path

`datalad.utils.create_tree` (*path*, *tree*, *archives\_leading\_dir*=True, *remove\_existing*=False)

Given a list of tuples (name, load) create such a tree

if load is a tuple itself – that would create either a subtree or an archive with that content and place it into the tree if name ends with `.tar.gz`

`datalad.utils.create_tree_archive` (*path*, *name*, *load*, *overwrite*=False, *archives\_leading\_dir*=True)

Given an archive *name*, create under *path* with specified *load* tree

`datalad.utils.decode_input(s)`

Given input string/bytes, decode according to stdin codepage (or UTF-8) if not defined

If fails – issue warning and decode allowing for errors being replaced

`datalad.utils.disable_logger(logger=None)`

context manager to temporarily disable logging

This is to provide one of `swallow_logs`' purposes without unnecessarily creating temp files (see gh-1865)

**Parameters** `logger` (*Logger*) – Logger whose handlers will be ordered to not log anything.  
Default: datalad's topmost Logger ('datalad')

`datalad.utils.dlabspath(path, norm=False)`

Symlinks-in-the-cwd aware abspath

`os.path.abspath` relies on `os.getcwd()` which would not know about symlinks in the path

TODO: we might want to `norm=True` by default to match behavior of `os.path.abspath?`

`datalad.utils.encode_filename(filename)`

Encode unicode filename

`datalad.utils.ensure_bool(s)`

Convert value into boolean following convention for strings

to recognize on, True, yes as True, off, False, no as False

`datalad.utils.ensure_bytes(s, encoding='utf-8')`

Convert/encode unicode string to bytes.

If `s` isn't a string, return it as is.

**Parameters** `encoding` (*str, optional*) – Encoding to use. "utf-8" is the default

`datalad.utils.ensure_dict_from_str(s, **kwargs)`

Given a multiline string with key=value items convert it to a dictionary

**Parameters**

- `s` (*str or dict*) –
- **None if input s is empty** (*Returns*) –

`datalad.utils.ensure_dir(*args)`

Make sure directory exists.

Joins the list of arguments to an os-specific path to the desired directory and creates it, if it not exists yet.

`datalad.utils.ensure_iter(s, cls, copy=False, iterate=True)`

Given not a list, would place it into a list. If None - empty list is returned

**Parameters**

- `s` (*list or anything*) –
- `cls` (*class*) – Which iterable class to ensure
- `copy` (*bool, optional*) – If correct iterable is passed, it would generate its shallow copy
- `iterate` (*bool, optional*) – If it is not a list, but something iterable (but not a str) iterate over it.

`datalad.utils.ensure_list(s, copy=False, iterate=True)`

Given not a list, would place it into a list. If None - empty list is returned



### Parameters

- **s** (*list or anything*) –
- **copy** (*bool, optional*) – If list is passed, it would generate a shallow copy of the list
- **iterate** (*bool, optional*) – If it is not a list, but something iterable (but not a str) iterate over it.

`datalad.utils.ensure_list_from_str(s, sep='\n')`

Given a multiline string convert it to a list of return None if empty

**Parameters** **s** (*str or list*) –

`datalad.utils.ensure_tuple_or_list(obj)`

Given an object, wrap into a tuple if not list or tuple

`datalad.utils.ensure_unicode(s, encoding=None, confidence=None)`

Convert/decode bytestring to unicode.

If *s* isn't a bytestring, return it as is.

### Parameters

- **encoding** (*str, optional*) – Encoding to use. If None, “utf-8” is tried, and then if not a valid UTF-8, encoding will be guessed
- **confidence** (*float, optional*) – A value between 0 and 1, so if guessing of encoding is of lower than specified confidence, ValueError is raised

`datalad.utils.escape_filename(filename)`

Surround filename in “” and escape ” in the filename

`datalad.utils.expandpath(path, force_absolute=True)`

Expand all variables and user handles in a path.

By default return an absolute path

`datalad.utils.file_basename(name, return_ext=False)`

Strips up to 2 extensions of length up to 4 characters and starting with alpha not a digit, so we could get rid of .tar.gz etc

`datalad.utils.find_files(regex, topdir='.', exclude=None, exclude_vcs=True, exclude_datalad=False, dirs=False)`

Generator to find files matching regex

### Parameters

- **regex** (*basestring*) –
- **exclude** (*basestring, optional*) – Matches to exclude
- **exclude\_vcs** – If True, excludes commonly known VCS subdirectories. If string, used as regex to exclude those files (regex: `'/(?:(?:git|gitattributes|svn|bzr|hg)(?:/|$))'`)
- **exclude\_datalad** – If True, excludes files known to be datalad meta-data files (e.g. under .datalad/ subdirectory) (regex: `'/(?:(?:datalad)(?:/|$))'`)
- **topdir** (*basestring, optional*) – Directory where to search
- **dirs** (*bool, optional*) – Whether to match directories as well as files

`datalad.utils.generate_chunks(container, size)`

Given a container, generate chunks from it with size up to *size*

`datalad.utils.generate_file_chunks(files, cmd=None)`

Given a list of files, generate chunks of them to avoid exceeding cmdline length

### Parameters

- **files** (*list of str*) –
- **cmd** (*str or list of str, optional*) – Command to account for as well

`datalad.utils.get_dataset_root(path)`

Return the root of an existent dataset containing a given path

The root path is returned in the same absolute or relative form as the input argument. If no associated dataset exists, or the input path doesn't exist, None is returned.

If *path* is a symlink or something other than a directory, its the root dataset containing its parent directory will be reported. If none can be found, at a symlink at *path* is pointing to a dataset, *path* itself will be reported as the root.

`datalad.utils.get_encoding_info()`

Return a dictionary with various encoding/locale information

`datalad.utils.get_envvars_info()`

`datalad.utils.get_func_kwargs_doc(func)`

Provides args for a function

**Parameters** **func** (*str*) – name of the function from which args are being requested

**Returns** of the args that a function takes in

**Return type** list

`datalad.utils.get_ipython_shell()`

Detect if running within IPython and returns its *ip* (shell) object

Returns None if not under ipython (no *get\_ipython* function)

`datalad.utils.get_linux_distribution()`

Compatibility wrapper for `{platform,distro}.linux_distribution()`.

`datalad.utils.get_logfilename(dspath, cmd='datalad')`

Return a filename to use for logging under a dataset/repository

directory would be created if doesn't exist, but *dspath* must exist and be a directory

`datalad.utils.get_open_files(path, log_open=False)`

Get open files under a path

### Parameters

- **path** (*str*) – File or directory to check for open files under
- **log\_open** (*bool or int*) – If set - logger level to use

**Returns** path : pid

**Return type** dict

`datalad.utils.get_path_prefix(path, pwd=None)`

Get path prefix (for current directory)

Returns relative path to the topdir, if we are under topdir, and if not absolute path to topdir. If *pwd* is not specified - current directory assumed

`datalad.utils.get_suggestions_msg(values, known, sep='\n')`

Return a formatted string with suggestions for values given the known ones

`datalad.utils.get_tempfile_kwargs(tkargs=None, prefix='', wrapped=None)`

Updates kwargs to be passed to tempfile. calls depending on env vars

`datalad.utils.get_timestamp_suffix` (*time\_=None, prefix='-'*)

Return a time stamp (full date and time up to second)

primarily to be used for generation of log files names

`datalad.utils.get_trace` (*edges, start, end, trace=None*)

Return the trace/path to reach a node in a tree.

#### Parameters

- **edges** (*sequence (2-tuple)*) – The tree given by a sequence of edges (parent, child) tuples. The nodes can be identified by any value and data type that supports the ‘==’ operation.
- **start** – Identifier of the start node. Must be present as a value in the parent location of an edge tuple in order to be found.
- **end** – Identifier of the target/end node. Must be present as a value in the child location of an edge tuple in order to be found.
- **trace** (*list*) – Mostly useful for recursive calls, and used internally.

**Returns** Returns a list with the trace to the target (the starts and the target are not included in the trace, hence if start and end are directly connected an empty list is returned), or None when no trace to the target can be found, or start and end are identical.

**Return type** None or list

`datalad.utils.get_wrapped_class` (*wrapped*)

Determine the command class a wrapped `__call__` belongs to

`datalad.utils.getargspec` (*func*)

`datalad.utils.getpwd` ()

Try to return a CWD without dereferencing possible symlinks

This function will try to use PWD environment variable to provide a current working directory, possibly with some directories along the path being symlinks to other directories. Unfortunately, PWD is used/set only by the shell and such functions as `os.chdir` and `os.getcwd` nohow use or modify it, thus `os.getcwd()` returns path with links dereferenced.

While returning current working directory based on PWD env variable we verify that the directory is the same as `os.getcwd()` after resolving all symlinks. If that verification fails, we fall back to always use `os.getcwd()`.

Initial decision to either use PWD env variable or `os.getcwd()` is done upon the first call of this function.

`datalad.utils.import_module_from_file` (*modpath, pkg=None, log=<bound method Logger.debug of <Logger datalad.utils (INFO)>>*)

Import provided module given a path

TODO: - RF/make use of it in pipeline.py which has similar logic - join with `import_modules` above?

**Parameters** **pkg** (*module, optional*) – If provided, and `modpath` is under `pkg.__path__`, relative import will be used

`datalad.utils.import_modules` (*modnames, pkg, msg='Failed to import {module}', log=<bound method Logger.debug of <Logger datalad.utils (INFO)>>*)

Helper to import a list of modules without failing if N/A

#### Parameters

- **modnames** (*list of str*) – List of module names to import
- **pkg** (*str*) – Package under which to import

- **msg**(*str*, *optional*) – Message template for `.format()` to log at DEBUG level if import fails. Keys {`module`} and {`package`} will be provided and ‘: {`exception`}’ appended
- **log**(*callable*, *optional*) – Logger call to use for logging messages

`datalad.utils.is_explicit_path`(*path*)

Return whether a path explicitly points to a location

Any absolute path, or relative path starting with either ‘`./`’ or ‘`./`’ is assumed to indicate a location on the filesystem. Any other path format is not considered explicit.

`datalad.utils.is_interactive`()

Return True if all in/out are open and tty.

Note that in a somewhat abnormal case where e.g. `stdin` is explicitly closed, and any operation on it would raise a `ValueError` (“*I/O operation on closed file*”) exception, this function would just return False, since the session cannot be used interactively.

`datalad.utils.knows_annex`(*path*)

Returns whether at a given path there is information about an annex

It is just a thin wrapper around `GitRepo.is_with_annex()` classmethod which also checks for *path* to exist first.

This includes actually present annexes, but also uninitialized ones, or even the presence of a remote annex branch.

`datalad.utils.line_profile`(*func*)

`datalad.utils.lmtime`(*filepath*, *mtime*)

Set `mtime` for files, while not de-referencing symlinks.

To overcome absence of `os.lutime`

Works only on linux and OSX ATM

`datalad.utils.make_tempfile`(*content=None*, *wrapped=None*, *\*\*kwargs*)

Helper class to provide a temporary file name and remove it at the end (context manager)

#### Parameters

- **mkdir**(*bool*, *optional* (default: `False`)) – If True, temporary directory created using `tempfile.mkdtemp()`
- **content**(*str* or *bytes*, *optional*) – Content to be stored in the file created
- **wrapped**(*function*, *optional*) – If set, function name used to prefix temporary file name
- **\*\*kwargs** – All other arguments are passed into the call to `tempfile.mk{,d}temp()`, and resultant temporary filename is passed as the first argument into the function `t`. If no ‘`prefix`’ argument is provided, it will be constructed using module and function names (‘`.`’ replaced with ‘`_`’).
- **change the used directory without providing keyword argument ‘`dir`’ set**(*To*) –
- **DATALAD\_TESTS\_TEMP\_DIR.** –

#### Examples

```
>>> from os.path import exists
>>> from datalad.utils import make_tempfile
>>> with make_tempfile() as fname:
...     k = open(fname, 'w').write('silly test')
>>> assert not exists(fname) # was removed
```

```
>>> with make_tempfile(content="blah") as fname:
...     assert open(fname).read() == "blah"
```

`datalad.utils.map_items` (*func*, *v*)

A helper to apply *func* to all elements (keys and values) within dict

No type checking of values passed to *func* is done, so *func* should be resilient to values which it should not handle

Initial usecase - `apply_recursive(url_fragment, ensure_unicode)`

`datalad.utils.md5sum` (*filename*)

Compute an MD5 sum for the given file

`datalad.utils.never_fail` (*f*)

Assure that function never fails – all exceptions are caught

Returns *None* if function fails internally.

`datalad.utils.not_supported_on_windows` (*msg=None*)

A little helper to be invoked to consistently fail whenever functionality is not supported (yet) on Windows

`datalad.utils.nothing_cm` ()

Just a dummy cm to programmically switch context managers

`datalad.utils.open_r_encdetect` (*fname*, *readahead=1000*)

Return a file object in read mode with auto-detected encoding

This is helpful when dealing with files of unknown encoding.

**Parameters** `readahead` (*int*, *optional*) – How many bytes to read for guessing the encoding type. If negative - full file will be read

`datalad.utils.optional_args` (*decorator*)

allows a decorator to take optional positional and keyword arguments. Assumes that taking a single, callable, positional argument means that it is decorating a function, i.e. something like this:

```
@my_decorator
def function(): pass
```

Calls decorator with `decorator(f, *args, **kwargs)`

`datalad.utils.partition` (*items*, *predicate=<class 'bool'>*)

Partition *items* by *predicate*.

#### Parameters

- **items** (*iterable*) –
- **predicate** (*callable*) – A function that will be mapped over each element in *items*. The elements will be partitioned based on whether the return value is false or true.

#### Returns

- A tuple with two generators, the first for 'false' items and the second for
- 'true' ones.

## Notes

Taken from Peter Otten's snippet posted at [https://nedbatchelder.com/blog/201306/filter\\_a\\_list\\_into\\_two\\_parts.html](https://nedbatchelder.com/blog/201306/filter_a_list_into_two_parts.html)

`datalad.utils.path_is_subpath` (*path*, *prefix*)

Return True if path is a subpath of prefix

It will return False if path == prefix.

### Parameters

- **path** (*str*) –
- **prefix** (*str*) –

`datalad.utils.path_startswith` (*path*, *prefix*)

Return True if path starts with prefix path

### Parameters

- **path** (*str*) –
- **prefix** (*str*) –

`datalad.utils.posix_relpath` (*path*, *start=None*)

Behave like `os.path.relpath`, but always return POSIX paths...

on any platform.

`datalad.utils.quote_cmdlinearg` (*arg*)

Perform platform-appropriate argument quoting

`datalad.utils.read_csv_lines` (*fname*, *dialect=None*, *readahead=16384*, *\*\*kwargs*)

A generator of dict records from a CSV/TSV

Automatically guesses the encoding for each record to convert to UTF-8

### Parameters

- **fname** (*str*) – Filename
- **dialect** (*str*, *optional*) – Dialect to specify to `csv.reader`. If not specified – guessed from the file, if fails to guess, “excel-tab” is assumed
- **readahead** (*int*, *optional*) – How many bytes to read from the file to guess the type
- **\*\*kwargs** – Passed to `csv.reader`

`datalad.utils.rmdir` (*path*, *\*args*, *\*\*kwargs*)

`os.rmdir` with our optional checking for open files

`datalad.utils.rmtemp` (*f*, *\*args*, *\*\*kwargs*)

Wrapper to centralize removing of temp files so we could keep them around

It will not remove the temporary file/directory if `DATALAD_TESTS_TEMP_KEEP` environment variable is defined

`datalad.utils.rmtree` (*path*, *chmod\_files='auto'*, *children\_only=False*, *\*args*, *\*\*kwargs*)

To remove git-annex .git it is needed to make all files and directories writable again first

### Parameters

- **chmod\_files** (*string or bool*, *optional*) – Whether to make files writable also before removal. Usually it is just a matter of directories to have write permissions. If ‘auto’ it would `chmod` files on windows by default

- **children\_only** (*bool, optional*) – If set, all files and subdirectories would be removed while the path itself (must be a directory) would be preserved
- **\*args** –
- **\*\*kwargs** – Passed into `shutil.rmtree` call

`datalad.utils.rotree(path, ro=True, chmod_files=True)`

To make tree read-only or writable

#### Parameters

- **path** (*string*) – Path to the tree/directory to `chmod`
- **ro** (*bool, optional*) – Whether to make it R/O (default) or RW
- **chmod\_files** (*bool, optional*) – Whether to operate also on files (not just directories)

`datalad.utils.safe_print(s)`

Print with protection against UTF-8 encoding errors

`datalad.utils.saved_generator(gen)`

Given a generator returns two generators, where 2nd one just replays

So the first one would be going through the generated items and 2nd one would be yielding saved items

`datalad.utils.setup_exceptionhook(ipython=False)`

Overloads default `sys.excepthook` with our exceptionhook handler.

If interactive, our exceptionhook handler will invoke `pdb.post_mortem`; if not interactive, then invokes default handler.

`datalad.utils.shortened_repr(value, l=30)`

`datalad.utils.slash_join(base, extension)`

Join two strings with a `'/'`, avoiding duplicate slashes

If any of the strings is `None` the other is returned as is.

`datalad.utils.sorted_files(dout)`

Return a (sorted) list of files under `dout`

`datalad.utils.split_cmdline(s)`

Perform platform-appropriate command line splitting.

Identical to `shlex.split()` on non-windows platforms.

Modified from <https://stackoverflow.com/a/35900070>

`datalad.utils.swallow_logs(new_level=None, file_=None, name='datalad')`

Context manager to consume all logs.

`datalad.utils.swallow_outputs()`

Context manager to help consuming both `stdout` and `stderr`, and `print()`

`stdout` is available as `cm.out` and `stderr` as `cm.err` whenever `cm` is the yielded context manager. Internally uses temporary files to guarantee absent side-effects of swallowing into `StringIO` which lacks `.fileno`.

`print` mocking is necessary for some uses where `sys.stdout` was already bound to original `sys.stdout`, thus mocking it later had no effect. Overriding `print` function had desired effect

`datalad.utils.try_multiple(ntrials, exception, base, f, *args, **kwargs)`

Call `f` multiple times making exponentially growing delay between the calls

`datalad.utils.try_multiple_dec(f, ntrials=None, duration=0.1, exceptions=None, increment_type=None)`

`datalad.utils.unique` (*seq*, *key=None*, *reverse=False*)

Given a sequence return a list only with unique elements while maintaining order

This is the fastest solution. See <https://www.peterbe.com/plog/uniqifiers-benchmark> and <http://stackoverflow.com/a/480227/1265472> for more information. Enhancement – added ability to compare for uniqueness using a key function

### Parameters

- **seq** – Sequence to analyze
- **key** (*callable*, *optional*) – Function to call on each element so we could decide not on a full element, but on its member etc
- **reverse** (*bool*, *optional*) – If True, uniqueness checked in the reverse order, so that the later ones will take the order

`datalad.utils.unlink` (*f*)

‘Robust’ unlink. Would try multiple times

On windows boxes there is evidence for a latency of more than a second until a file is considered no longer “in-use”. `WindowsError` is not known on Linux, and if `IOError` or any other exception is thrown then if except statement has `WindowsError` in it – `NameError` also see gh-2533

`datalad.utils.updated` (*d*, *update*)

Return a copy of the input with the ‘update’

Primarily for updating dictionaries

`datalad.utils.with_pathsep` (*path*)

Little helper to guarantee that path ends with /

## datalad.version

Defines version to be imported in the module and obtained from `setup.py`

## datalad.support.gitrepo

Interface to Git via GitPython

For further information on GitPython see <http://gitpython.readthedocs.org/>

**class** `datalad.support.gitrepo.FetchInfo`

Bases: `dict`

dict that carries results of a fetch operation of a single head

Reduced variant of GitPython’s `RemoteProgress` class

**Original copyright:** Copyright (C) 2008, 2009 Michael Trier and contributors

**Original license:** BSD 3-Clause “New” or “Revised” License

**ERROR** = 128

**FAST\_FORWARD** = 64

**FORCED\_UPDATE** = 32

**HEAD\_UPTODATE** = 4

**NEW\_HEAD** = 2



**NEW\_TAG = 1**

**REJECTED = 16**

**TAG\_UPDATE = 8**

**class** datalad.support.gitrepo.**GitProgress**(\*args)

Bases: *datalad.cmd.WitlessProtocol*

Reduced variant of GitPython's RemoteProgress class

**Original copyright:** Copyright (C) 2008, 2009 Michael Trier and contributors

**Original license:** BSD 3-Clause "New" or "Revised" License

**BEGIN = 1**

**CHECKING\_OUT = 256**

**COMPRESSING = 8**

**COUNTING = 4**

**DONE\_TOKEN = 'done.'**

**END = 2**

**ENUMERATING = 512**

**FINDING\_SOURCES = 128**

**OP\_MASK = -4**

**RECEIVING = 32**

**RESOLVING = 64**

**STAGE\_MASK = 3**

**TOKEN\_SEPARATOR = ', '**

**WRITING = 16**

**connection\_made**(*transport*)

Called when a connection is made.

The argument is the transport representing the pipe connection. To receive data, wait for `data_received()` calls. When the connection is closed, `connection_lost()` is called.

**pipe\_data\_received**(*fd*, *bytes*)

Called when the subprocess writes data into stdout/stderr pipe.

`fd` is int file descriptor. `data` is bytes object.

**proc\_err = True**

**process\_exited**()

Called when subprocess has exited.

**re\_op\_absolute** = `re.compile(' (remote: )? ([\\w\\s]+) :\\s+ ( (\\d+) ) ( .* )')`

**re\_op\_relative** = `re.compile(' (remote: )? ([\\w\\s]+) :\\s+ (\\d+)% \\( (\\d+) / (\\d+) \\) ( .`

**class** datalad.support.gitrepo.**GitRepo**(*path*, *url=None*, *runner=None*, *create=True*,  
*git\_opts=None*, *repo=None*, *fake\_dates=False*,  
*create\_sanity\_checks=True*, *\*\*kwargs*)

Bases: *datalad.support.repo.RepoInterface*

Representation of a git repository

**add** (*files*, *git=True*, *git\_options=None*, *update=False*)  
 Adds file(s) to the repository.

**Parameters**

- **files** (*list*) – list of paths to add
- **git** (*bool*) – somewhat ugly construction to be compatible with AnnexRepo.add(); has to be always true.
- **update** (*bool*) –

–**update option for git-add. From git’s manpage:** Update the index just where it already has an entry matching <pathspec>. This removes as well as modifies index entries to match the working tree, but adds no new files.

If no <pathspec> is given when –update option is used, all tracked files in the entire working tree are updated (old versions of Git used to limit the update to the current directory and its subdirectories).

**Returns** Of status dicts.

**Return type** list

**add\_** (*files*, *git=True*, *git\_options=None*, *update=False*)  
 Like *add*, but returns a generator

**add\_fake\_dates** (*env*)  
 Add fake dates to *env*.

**Parameters** *env* (*dict* or *None*) – Environment variables.

**Returns**

- A *dict* (copied from *env*), with date-related environment
- variables for *git* and *git-annex* set.

**add\_remote** (*name*, *url*, *options=None*)  
 Register remote pointing to a url

**add\_submodule** (*path*, *name=None*, *url=None*, *branch=None*)  
 Add a new submodule to the repository.

This will alter the index as well as the .gitmodules file, but will not create a new commit. If the submodule already exists, no matter if the configuration differs from the one provided, the existing submodule is considered as already added and no further action is performed.

**Parameters**

- **path** (*str*) – repository-relative path at which the submodule should be located, and which will be created as required during the repository initialization.
- **name** (*str* or *None*) – name/identifier for the submodule. If *None*, the *path* will be used as name.
- **url** (*str* or *None*) – git-clone compatible URL. If *None*, the repository is assumed to exist, and the url of the first remote is taken instead. This is useful if you want to make an existing repository a submodule of another one.
- **branch** (*str* or *None*) – name of branch to be checked out in the submodule. The given branch must exist in the remote repository, and will be checked out locally as a tracking branch. If *None*, remote HEAD will be checked out.

**call\_git** (*args*, *files=None*, *expect\_stderr=False*, *expect\_fail=False*)

Call git and return standard output.

#### Parameters

- **args** (*list of str*) – Arguments to pass to *git*.
- **files** (*list of str, optional*) – File arguments to pass to *git*. The advantage of passing these here rather than as part of *args* is that the call will be split into multiple calls to avoid exceeding the maximum command line length.
- **expect\_stderr** (*bool, optional*) – Standard error is expected and should not be elevated above the DEBUG level.
- **expect\_fail** (*bool, optional*) – A non-zero exit is expected and should not be elevated above the DEBUG level.

#### Returns

**Return type** standard output (str)

**Raises** CommandError if the call exits with a non-zero status.

**call\_git\_items\_** (*args*, *files=None*, *expect\_stderr=False*, *sep=None*)

Call git, splitting output on *sep*.

#### Parameters

- **args** (*list of str*) – Arguments to pass to *git*.
- **files** (*list of str, optional*) – File arguments to pass to *git*. The advantage of passing these here rather than as part of *args* is that the call will be split into multiple calls to avoid exceeding the maximum command line length.
- **expect\_stderr** (*bool, optional*) – Standard error is expected and should not be elevated above the DEBUG level.
- **sep** (*str, optional*) – Split the output by *str.split(sep)* rather than *str.splitlines*.

#### Returns

**Return type** Generator that yields output items.

**Raises** CommandError if the call exits with a non-zero status.

**call\_git\_online** (*args*, *files=None*, *expect\_stderr=False*)

Call git for a single line of output.

#### Parameters

- **args** (*list of str*) – Arguments to pass to *git*.
- **files** (*list of str, optional*) – File arguments to pass to *git*. The advantage of passing these here rather than as part of *args* is that the call will be split into multiple calls to avoid exceeding the maximum command line length.
- **expect\_stderr** (*bool, optional*) – Standard error is expected and should not be elevated above the DEBUG level.
- **sep** (*str, optional*) – Split the output by *str.split(sep)* rather than *str.splitlines*.

#### Raises

- CommandError if the call exits with a non-zero status.
- AssertionError if there is more than one line of output.

**call\_git\_success** (*args*, *files=None*, *expect\_stderr=False*)

Call git and return true if the call exit code of 0.

**Parameters**

- **args** (*list of str*) – Arguments to pass to *git*.
- **files** (*list of str, optional*) – File arguments to pass to *git*. The advantage of passing these here rather than as part of *args* is that the call will be split into multiple calls to avoid exceeding the maximum command line length.
- **expect\_stderr** (*bool, optional*) – Standard error is expected and should not be elevated above the DEBUG level.

**Returns**

**Return type** bool

**checkout** (*name*, *options=None*)

**cherry\_pick** (*commit*)

Cherry pick *commit* to the current branch.

**Parameters** **commit** (*str*) – A single commit.

**classmethod clone** (*url*, *path*, *\*args*, *clone\_options=None*, *\*\*kwargs*)

Clone url into path

Provides workarounds for known issues (e.g. <https://github.com/datalad/datalad/issues/785>)

**Parameters**

- **url** (*str*) –
- **path** (*str*) –
- **clone\_options** (*dict*) – Key/value pairs of arbitrary options that will be passed on to the underlying call to *git-clone*.
- **expect\_fail** (*bool*) – Whether expect that command might fail, so error should be logged then at DEBUG level instead of ERROR

**commit** (*msg=None*, *options=None*, *\_datalad\_msg=False*, *careless=True*, *files=None*, *date=None*, *index\_file=None*)

Commit changes to git.

**Parameters**

- **msg** (*str, optional*) – commit-message
- **options** (*list of str, optional*) – cmdline options for git-commit
- **\_datalad\_msg** (*bool, optional*) – To signal that commit is automated commit by datalad, so it would carry the [DATALAD] prefix
- **careless** (*bool, optional*) – if False, raise when there’s nothing actually committed; if True, don’t care
- **files** (*list of str, optional*) – path(s) to commit
- **date** (*str, optional*) – Date in one of the formats git understands
- **index\_file** (*str, optional*) – An alternative index to use

**commit\_exists** (*commitish*)

Does *commitish* exist in the repo?

**Parameters** **commitish** (*str*) – A commit or an object that can be dereferenced to one.

**Returns**

**Return type** bool

**config**

Get an instance of the parser for the persistent repository configuration.

Note: This allows to also read/write `.datalad/config`, not just `.git/config`

**Returns**

**Return type** *ConfigManager*

**configure\_fake\_dates** ()

Configure repository to use fake dates.

**count\_objects**

return dictionary with count, size(in KiB) information of git objects

**deinit\_submodule** (*path*, *\*\*kwargs*)

Deinit a submodule

**Parameters**

- **path** (*str*) – path to the submodule; relative to *self.path*
- **kwargs** – see `__init__`

**describe** (*commitish=None*, *\*\*kwargs*)

Quick and dirty implementation to call git-describe

**Parameters** **kwargs** – transformed to cmdline options for git-describe; see `__init__` for description of the transformation

**diff** (*fr*, *to*, *paths=None*, *untracked='all'*, *eval\_submodule\_state='full'*)

Like status(), but reports changes between to arbitrary revisions

**Parameters**

- **fr** (*str or None*) – Revision specification (anything that Git understands). Passing *None* considers anything in the target state as new.
- **to** (*str or None*) – Revision specification (anything that Git understands), or *None* to compare to the state of the work tree.
- **paths** (*list or None*) – If given, limits the query to the specified paths. To query all paths specify *None*, not an empty list.
- **untracked** (*{'no', 'normal', 'all'}*) – If and how untracked content is reported when *to* is *None*: ‘no’: no untracked files are reported; ‘normal’: untracked files and entire untracked directories are reported as such; ‘all’: report individual files even in fully untracked directories.
- **eval\_submodule\_state** (*{'full', 'commit', 'no'}*) – If ‘full’ (the default), the state of a submodule is evaluated by considering all modifications, with the treatment of untracked files determined by *untracked*. If ‘commit’, the modification check is restricted to comparing the submodule’s HEAD commit to the one recorded in the superdataset. If ‘no’, the state of the subdataset is not evaluated.

**Returns**

Each content item has an entry under a pathlib *Path* object instance pointing to its absolute path inside the repository (this path is guaranteed to be underneath *Repo.path*). Each value is a dictionary with properties:

**type** Can be ‘file’, ‘symlink’, ‘dataset’, ‘directory’

**state** Can be ‘added’, ‘untracked’, ‘clean’, ‘deleted’, ‘modified’.

**Return type** dict

**diffstatus** (*fr, to, paths=None, untracked='all', eval\_submodule\_state='full', eval\_file\_type=True, \_cache=None*)

Like `diff()`, but reports the status of ‘clean’ content too

**dirty**

Is the repository dirty?

Note: This provides a quick answer when you simply want to know if there are any untracked changes or modifications in this repository or its submodules. For finer-grained control and more detailed reporting, use `status()` instead.

**fake\_dates\_enabled**

Is the repository configured to use fake dates?

**fetch** (*remote=None, refspec=None, all\_=False, git\_options=None, \*\*kwargs*)

Fetches changes from a remote (or all remotes).

**Parameters**

- **remote** (*str, optional*) – name of the remote to fetch from. If no remote is given and *all\_* is not set, the tracking branch is fetched.
- **refspec** (*str, optional*) – refspec to fetch.
- **all** (*bool, optional*) – fetch all remotes (and all of their branches). Fails if *remote* was given.
- **git\_options** (*list, optional*) – Additional command line options for git-fetch.
- **kwargs** – Deprecated. GitPython-style keyword argument for git-fetch. Will be appended to any *git\_options*.

**fetch\_** (*remote=None, refspec=None, all\_=False, git\_options=None*)

Like `fetch`, but returns a generator

**for\_each\_ref\_** (*fields=('objectname', 'objecttype', 'refname'), pattern=None, points\_at=None, sort=None, count=None, contains=None*)

Wrapper for `git for-each-ref`

Please see manual page `git-for-each-ref(1)` for a complete overview of its functionality. Only a subset of it is supported by this wrapper.

**Parameters**

- **fields** (*iterable or str*) – Used to compose a NULL-delimited specification for `for-each-ref's --format` option. The default field list reflects the standard behavior of `for-each-ref` when the `--format` option is not given.
- **pattern** (*list or str, optional*) – If provided, report only refs that match at least one of the given patterns.
- **points\_at** (*str, optional*) – Only list refs which points at the given object.
- **sort** (*list or str, optional*) – Field name(s) to sort-by. If multiple fields are given, the last one becomes the primary key. Prefix any field name with ‘-’ to sort in descending order.
- **count** (*int, optional*) – Stop iteration after the given number of matches.
- **contains** (*str, optional*) – Only list refs which contain the specified commit.

**Yields** dict with items matching the given *fields*

**Raises**

- `ValueError` – if no *fields* are given
- `RuntimeError` – if *git for-each-ref* returns a record where the number of properties does not match the number of *fields*

**format\_commit** (*fmt*, *commitish=None*)

Return *git show* output for *commitish*.

**Parameters**

- **fmt** (*str*) – A format string accepted by *git show*.
- **commitish** (*str*, *optional*) – Any commit identifier (defaults to “HEAD”).

**Returns**

**Return type** *str* or, if there are not commits yet, `None`.

**gc** (*allow\_background=False*, *auto=False*)

Perform house keeping (garbage collection, repacking)

**get\_active\_branch** ()

Get the name of the active branch

**Returns** Returns `None` if there is no active branch, i.e. detached HEAD, and the branch name otherwise.

**Return type** *str* or `None`

**get\_branch\_commits\_** (*branch=None*, *limit=None*, *stop=None*)

Return commit hexshas for a branch

**Parameters**

- **branch** (*str*, *optional*) – If not provided, assumes current branch
- **limit** (*None* | *'left-only'*, *optional*) – Limit which commits to report. If `None` – all commits (merged or not), if *'left-only'* – only the commits from the left side of the tree upon merges
- **stop** (*str*, *optional*) – hexsha of the commit at which stop reporting (matched one is not reported either)

**Yields** *str*

**get\_branches** ()

Get all branches of the repo.

**Returns** Names of all branches of this repository.

**Return type** [*str*]

**get\_commit\_date** (*branch=None*, *date='authored'*)

Get the date stamp of the last commit (in a branch or head otherwise)

**Parameters** **date** (*{'authored', 'committed'}*) – Which date to return. “authored” will be the date shown by “git show” and the one possibly specified via *-date* to *git commit*

**Returns** `None` if no commit

**Return type** *int* or `None`

**get\_content\_info** (*paths=None, ref=None, untracked='all', eval\_file\_type=True*)

Get identifier and type information from repository content.

This is simplified front-end for *git ls-files/tree*.

Both commands differ in their behavior when queried about subdataset paths. *ls-files* will not report anything, *ls-tree* will report on the subdataset record. This function uniformly follows the behavior of *ls-tree* (report on the respective subdataset mount).

#### Parameters

- **paths** (*list (pathlib.PurePath)*) – Specific paths, relative to the resolved repository root, to query info for. Paths must be normed to match the reporting done by Git, i.e. no parent dir components (ala “some/../this”). If none are given, info is reported for all content.
- **ref** (*gitref or None*) – If given, content information is retrieved for this Git reference (via *ls-tree*), otherwise content information is produced for the present work tree (via *ls-files*). With a given reference, the reported content properties also contain a ‘bytesize’ record, stating the size of a file in bytes.
- **untracked** (*{'no', 'normal', 'all'}*) – If and how untracked content is reported when no *ref* was given: ‘no’: no untracked files are reported; ‘normal’: untracked files and entire untracked directories are reported as such; ‘all’: report individual files even in fully untracked directories.
- **eval\_file\_type** (*bool*) – If True, inspect file type of untracked files, and report annex symlink pointers as type ‘file’. This convenience comes with a cost; disable to get faster performance if this information is not needed.

#### Returns

Each content item has an entry under a *pathlib Path* object instance pointing to its absolute path inside the repository (this path is guaranteed to be underneath *Repo.path*). Each value is a dictionary with properties:

**type** Can be ‘file’, ‘symlink’, ‘dataset’, ‘directory’

Note that the reported type will not always match the type of content committed to Git, rather it will reflect the nature of the content minus platform/mode-specifics. For example, a symlink to a locked annexed file on Unix will have a type ‘file’, reported, while a symlink to a file in Git or directory will be of type ‘symlink’.

**gitshasum** SHASUM of the item as tracked by Git, or None, if not tracked. This could be different from the SHASUM of the file in the worktree, if it was modified.

**Return type** dict

**Raises** *ValueError* – In case of an invalid Git reference (e.g. ‘HEAD’ in an empty repository)

**get\_files** (*branch=None*)

Get a list of files in git.

Lists the files in the (remote) branch.

**Parameters** **branch** (*str*) – Name of the branch to query. Default: active branch.

**Returns** list of files.

**Return type** [str]

**get\_git\_attributes** ()

Query gitattributes which apply to top level directory



It is a thin compatibility/shortcut wrapper around more versatile `get_gitattributes` which operates on a list of paths and returns a dictionary per each path

**Returns** a dictionary with attribute name and value items relevant for the top (‘.’) directory of the repository, and thus most likely the default ones (if not overwritten with more rules) for all files within repo.

**Return type** dict

**static** `get_git_dir(repo)`

figure out a repo’s gitdir

‘.git’ might be a directory, a symlink or a file

---

**Note:** Please try using `GitRepo.dot_git` instead! That one’s not static, but it’s cheaper and you should avoid not having an instance of a repo you’re working on anyway. Note, that the property in opposition to this method returns an absolute path.

---

**Parameters** `repo` (*path or Repo instance*) – currently expected to be the repos base dir

**Returns** relative path to the repo’s git dir; So, default would be “.git”

**Return type** str

`get_gitattributes(path, index_only=False)`

Query gitattributes for one or more paths

**Parameters**

- **path** (*path or list*) – Path(s) to query. Paths may be relative or absolute.
- **index\_only** (*bool*) – Flag whether to consider only gitattribute setting that are reflected in the repository index, not just in the work tree content.

**Returns** Each key is a queried path (always relative to the repository root), each value is a dictionary with attribute name and value items. Attribute values are either True or False, for set and unset attributes, or are the literal attribute value.

**Return type** dict

`get_hexsha(commitish=None, short=False)`

Return a hexsha for a given commitish.

**Parameters**

- **commitish** (*str, optional*) – Any identifier that refers to a commit (defaults to “HEAD”).
- **short** (*bool, optional*) – Return the abbreviated form of the hexsha.

**Returns**

**Return type** str or, if there are not commits yet, None.

`get_indexed_files()`

Get a list of files in git’s index

**Returns** list of paths rooting in git’s base dir

**Return type** list

**get\_last\_commit\_hexsha** (*files*)

Return the hash of the last commit the modified any of the given paths

**get\_merge\_base** (*commitishes*)

Get a merge base hexsha

**Parameters** **commitishes** (*str or list of str*) – List of commitishes (branches, hexshas, etc) to determine the merge base of. If a single value provided, returns `merge_base` with the current branch.

**Returns** If no merge-base for given commits, or specified treeish doesn't exist, `None` returned

**Return type** `str` or `None`

**get\_remote\_branches** ()

Get all branches of all remotes of the repo.

**Returns** Names of all remote branches.

**Return type** [`str`]

**get\_remote\_url** (*name, push=False*)

Get the url of a remote.

Reads the configuration of remote *name* and returns its url or `None`, if there is no url configured.

**Parameters**

- **name** (*str*) – name of the remote
- **push** (*bool*) – if `True`, get the pushurl instead of the fetch url.

**get\_remotes** (*with\_urls\_only=False*)

Get known remotes of the repository

**Parameters** **with\_urls\_only** (*bool, optional*) – return only remotes which have urls

**Returns** **remotes** – List of names of the remotes

**Return type** list of `str`

**get\_revisions** (*revrange=None, fmt='%H', options=None*)

Return list of revisions in *revrange*.

**Parameters**

- **revrange** (*str or list of str or None, optional*) – Revisions or revision ranges to walk. If `None`, revision defaults to `HEAD` unless a revision-modifying option like `-all` or `-branches` is included in *options*.
- **fmt** (*string, optional*) – Format accepted by `-format` option of `git log`. This should not contain new lines because the output is split on new lines.
- **options** (*list of str, optional*) – Options to pass to `git log`. This should not include `-format`.

**Returns**

**Return type** List of revisions (`str`), formatted according to *fmt*.

**get\_staged\_paths** ()

Returns a list of any stage repository path(s)

This is a rather fast call, as it will not depend on what is going on in the worktree.

**get\_submodules** (*sorted\_=True, paths=None, compat=True*)

Return list of submodules.

**Parameters**

- **sorted** (*bool, optional*) – Sort submodules by path name.
- **paths** (*list (pathlib.PurePath), optional*) – Restrict submodules to those under *paths*.
- **compat** (*bool, optional*) – If true, return a namedtuple that incompletely mimics the attributes of GitPython’s Submodule object in hope of backwards compatibility with previous callers. Note that this form should be considered temporary and callers should be updated; this flag will be removed in a future release.

**Returns**

- List of submodule namedtuples if *compat* is true or otherwise a list
- of dictionaries as returned by *get\_submodules\_*.

**get\_submodules\_** (*paths=None*)

Yield submodules in this repository.

**Parameters** **paths** (*list (pathlib.PurePath), optional*) – Restrict submodules to those under *paths*.

**Returns**

- *A generator that yields a dictionary with information for each*
- *submodule.*

**get\_tags** (*output=None*)

Get list of tags

**Parameters** **output** (*str, optional*) – If given, limit the return value to a list of values matching that particular key of the tag properties.

**Returns** Each item is a dictionary with information on a tag. At present this includes ‘hexsha’, and ‘name’, where the latter is the string label of the tag, and the former the hexsha of the object the tag is attached to. The list is sorted by the creator date (committer date for lightweight tags and tagger date for annotated tags), with the most recent commit being the last element.

**Return type** list

**classmethod** **get\_toppath** (*path, follow\_up=True, git\_options=None*)

Return top-level of a repository given the path.

**Parameters**

- **follow\_up** (*bool*) – If path has symlinks – they get resolved by git. If *follow\_up* is True, we will follow original path up until we hit the same resolved path. If no such path found, resolved one would be returned.
- **git\_options** (*list of str*) – options to be passed to the git rev-parse call
- **None if no parent directory contains a git repository.**  
(*Return*) –

**get\_tracking\_branch** (*branch=None, remote\_only=False*)

Get the tracking branch for *branch* if there is any.

**Parameters**

- **branch** (*str*) – local branch to look up. If none is given, active branch is used.

- **remote\_only** (*bool*) – Don’t return a value if the upstream remote is set to “.” (meaning this repository).

**Returns** (remote or None, refspec or None) of the tracking branch

**Return type** tuple

**is\_ancestor** (*reva, revb*)

Is *reva* an ancestor of *revb*?

**Parameters** **revb** (*reva,*) – Revisions.

**Returns**

**Return type** bool

**is\_valid\_git** ()

Returns whether the underlying repository appears to be still valid

Note, that this almost identical to the classmethod `is_valid_repo()`. However, if we are testing an existing instance, we can save Path object creations. Since this testing is done a lot, this is relevant. Creation of the Path objects in `is_valid_repo()` takes nearly half the time of the entire function.

Also note, that this method is bound to an instance but still class-dependent, meaning that a subclass cannot simply overwrite it. This is particularly important for the call from within `__init__()`, which in turn is called by the subclasses’ `__init__`. Using an overwrite would lead to the wrong thing being called.

**classmethod is\_valid\_repo** (*path*)

Returns if a given path points to a git repository

**is\_with\_annex** ()

Report if GitRepo (assumed) has (remotes with) a git-annex branch

**merge** (*name, options=None, msg=None, allow\_unrelated=False, \*\*kwargs*)

**precommit** ()

Perform pre-commit maintenance tasks

**pull** (*remote=None, refspec=None, git\_options=None, \*\*kwargs*)

Pulls changes from a remote.

**Parameters**

- **remote** (*str, optional*) – name of the remote to pull from. If no remote is given, the remote tracking branch is used.
- **refspec** (*str, optional*) – refspec to fetch.
- **git\_options** (*list, optional*) – Additional command line options for git-pull.
- **kwargs** – Deprecated. GitPython-style keyword argument for git-pull. Will be appended to any `git_options`.

**push** (*remote=None, refspec=None, all\_remotes=False, all\_=False, git\_options=None, \*\*kwargs*)

Push changes to a remote (or all remotes).

**Parameters**

- **remote** (*str, optional*) – name of the remote to push to. If no remote is given and `all_` is not set, the tracking branch is pushed.
- **refspec** (*str, optional*) – refspec to push.
- **all** (*bool, optional*) – push to all remotes. Fails if *remote* was given.
- **git\_options** (*list, optional*) – Additional command line options for git-push.

- **kwargs** – Deprecated. GitPython-style keyword argument for git-push. Will be appended to any `git_options`.

**push\_** (*remote=None, refspec=None, all\_=False, git\_options=None*)

Like *push*, but returns a generator

**remove** (*files, recursive=False, \*\*kwargs*)

Remove files.

Calls `git-rm`.

#### Parameters

- **files** (*str*) – list of paths to remove
- **recursive** (*False*) – whether to allow recursive removal from subdirectories
- **kwargs** – see `__init__`

**Returns** list of successfully removed files.

**Return type** [str]

**remove\_branch** (*branch*)

**remove\_remote** (*name*)

Remove existing remote

**repo**

**save** (*message=None, paths=None, \_status=None, \*\*kwargs*)

Save dataset content.

#### Parameters

- **message** (*str or None*) – A message to accompany the changeset in the log. If `None`, a default message is used.
- **paths** (*list or None*) – Any content with path matching any of the paths given in this list will be saved. Matching will be performed against the dataset status (`GitRepo.status()`), or a custom status provided via `_status`. If no paths are provided, ALL non-clean paths present in the repo status or `_status` will be saved.
- **\_status** (*dict or None*) – If `None`, `Repo.status()` will be queried for the given *ds*. If a dict is given, its content will be used as a constraint. For example, to save only modified content, but no untracked content, set *paths* to `None` and provide a `_status` that has no entries for untracked content.
- **\*\*kwargs** – Additional arguments that are passed to underlying Repo methods. Supported:
  - `git` : bool (passed to `Repo.add()`)
  - `eval_submodule_state` : { 'full', 'commit', 'no' } passed to `Repo.status()`
  - `untracked` : { 'no', 'normal', 'all' } - passed to `Repo.satus()`

**save\_** (*message=None, paths=None, \_status=None, \*\*kwargs*)

Like *save()* but working as a generator.

**set\_gitattributes** (*attrs, attrfile='.gitattributes', mode='a'*)

Set gitattributes

By default appends additional lines to *attrfile*. Note, that later lines in *attrfile* overrule earlier ones, which may or may not be what you want. Set *mode* to 'w' to replace the entire file by what you provided in *attrs*.

**Parameters**

- **attrs** (*list*) – Each item is a 2-tuple, where the first element is a path pattern, and the second element is a dictionary with attribute key/value pairs. The attribute dictionary must use the same semantics as those returned by `get_gitattributes()`. Path patterns can use absolute paths, in which case they will be normalized relative to the directory that contains the target `.gitattributes` file (see *attrfile*).
- **attrfile** (*path*) – Path relative to the repository root of the `.gitattributes` file the attributes shall be set in.
- **mode** (*str*) – ‘a’ to append `.gitattributes`, ‘w’ to replace it

**set\_remote\_url** (*name, url, push=False*)

Set the URL a remote is pointing to

Sets the URL of the remote *name*. Requires the remote to already exist.

**Parameters**

- **name** (*str*) – name of the remote
- **url** (*str*) –
- **push** (*bool*) – if True, set the push URL, otherwise the fetch URL

**status** (*paths=None, untracked='all', eval\_submodule\_state='full'*)

Simplified `git status` equivalent.

**Parameters**

- **paths** (*list or None*) – If given, limits the query to the specified paths. To query all paths specify *None*, not an empty list. If a query path points into a subdataset, a report is made on the subdataset record within the queried dataset only (no recursion).
- **untracked** (*{'no', 'normal', 'all'}*) – If and how untracked content is reported: ‘no’: no untracked files are reported; ‘normal’: untracked files and entire untracked directories are reported as such; ‘all’: report individual files even in fully untracked directories.
- **eval\_submodule\_state** (*{'full', 'commit', 'no'}*) – If ‘full’ (the default), the state of a submodule is evaluated by considering all modifications, with the treatment of untracked files determined by *untracked*. If ‘commit’, the modification check is restricted to comparing the submodule’s HEAD commit to the one recorded in the superdataset. If ‘no’, the state of the subdataset is not evaluated.

**Returns**

Each content item has an entry under a `pathlib Path` object instance pointing to its absolute path inside the repository (this path is guaranteed to be underneath *Repo.path*). Each value is a dictionary with properties:

**type** Can be ‘file’, ‘symlink’, ‘dataset’, ‘directory’

**state** Can be ‘added’, ‘untracked’, ‘clean’, ‘deleted’, ‘modified’.

**Return type** dict

**tag** (*tag, message=None, commit=None, options=None*)

Tag a commit

**Parameters**

- **tag** (*str*) – Custom tag label. Must be a valid tag name.

- **message** (*str*, *optional*) – If provided, adds [`-m`, `<message>`] to the list of *git tag* arguments.
- **commit** (*str*, *optional*) – If provided, will be appended as last argument to the *git tag* call, and can be used to identify the commit that shall be tagged, if not HEAD.
- **options** (*list*, *optional*) – Additional command options, inserted prior a potential *commit* argument.

#### **untracked\_files**

Legacy interface, do not use! Use the `status()` method instead.

Despite its name, it also reports on untracked datasets, and yields their names with trailing path separators.

#### **update\_ref** (*ref*, *value*, *symbolic=False*)

Update the object name stored in a ref “safely”.

Just a shim for *git update-ref* call if not symbolic, and *git symbolic-ref* if symbolic

##### **Parameters**

- **ref** (*str*) – Reference, such as *ref/heads/BRANCHNAME* or HEAD.
- **value** (*str*) – Value to update to, e.g. hexsha of a commit when updating for a branch ref, or branch ref if updating HEAD
- **symbolic** (*None*) – To instruct if ref is symbolic, e.g. should be used in case of `ref=HEAD`

#### **update\_remote** (*name=None*, *verbose=False*)

#### **update\_submodule** (*path*, *mode='checkout'*, *init=False*)

Update a registered submodule.

This will make the submodule match what the superproject expects by cloning missing submodules and updating the working tree of the submodules. The “updating” can be done in several ways depending on the value of `submodule.<name>.update` configuration variable, or the *mode* argument.

##### **Parameters**

- **path** (*str*) – Identifies which submodule to operate on by its repository-relative path.
- **mode** (*{checkout, rebase, merge}*) – Update procedure to perform. ‘checkout’: the commit recorded in the superproject will be checked out in the submodule on a detached HEAD; ‘rebase’: the current branch of the submodule will be rebased onto the commit recorded in the superproject; ‘merge’: the commit recorded in the superproject will be merged into the current branch in the submodule.
- **init** (*bool*) – If True, initialize all submodules for which “git submodule init” has not been called so far before updating. Primarily provided for internal purposes and should not be used directly since would result in not so annex-friendly `.git` symlinks/references instead of full featured `.git/` directories in the submodules

#### **class** `datalad.support.gitrepo.PushInfo`

Bases: `dict`

dict that carries results of a push operation of a single head

Reduced variant of `GitPython's RemoteProgress` class

**Original copyright:** Copyright (C) 2008, 2009 Michael Trier and contributors

**Original license:** BSD 3-Clause “New” or “Revised” License

**DELETED = 64**

```
ERROR = 1024
FAST_FORWARD = 256
FORCED_UPDATE = 128
NEW_HEAD = 2
NEW_TAG = 1
NO_MATCH = 4
REJECTED = 8
REMOTE_FAILURE = 32
REMOTE_REJECTED = 16
UP_TO_DATE = 512
```

`datalad.support.gitrepo.Repo(*args, **kwargs)`  
Factory method around `gitpy.Repo` to consistently initiate with different backend

```
class datalad.support.gitrepo.StdoutCaptureWithGitProgress(*args)
    Bases: datalad.support.gitrepo.GitProgress

    proc_out = True
```

```
class datalad.support.gitrepo.Submodule(name, path, url)
    Bases: tuple

    name
        Alias for field number 0

    path
        Alias for field number 1

    url
        Alias for field number 2
```

`datalad.support.gitrepo.guard_BadName(func)`  
A helper to guard against `BadName` exception

Workaround for <https://github.com/gitpython-developers/GitPython/issues/768> also see <https://github.com/datalad/datalad/issues/2550> Let's try to precommit (to flush anything flushable) and do it again

`datalad.support.gitrepo.to_options(split_single_char_options=True, **kwargs)`  
Transform keyword arguments into a list of cmdline options  
Imported from `GitPython`.

**Original copyright:** Copyright (C) 2008, 2009 Michael Trier and contributors

**Original license:** BSD 3-Clause "New" or "Revised" License

#### Parameters

- `split_single_char_options` (*bool*) –
- `kwargs` –

#### Returns

Return type list



## datalad.support.annexrepo

Interface to git-annex by Joey Hess.

For further information on git-annex see <https://git-annex.branchable.com/>.

```
class datalad.support.annexrepo.AnnexRepo(path, url=None, runner=None, back-
end=None, always_commit=True, cre-
ate=True, create_sanity_checks=True,
init=False, batch_size=None, version=None,
description=None, git_opts=None, an-
nex_opts=None, annex_init_opts=None,
repo=None, fake_dates=False)
```

Bases: `datalad.support.gitrepo.GitRepo`, `datalad.support.repo.RepoInterface`

Representation of an git-annex repository.

Paths given to any of the class methods will be interpreted as relative to PWD, in case this is currently beneath AnnexRepo's base dir (*self.path*). If PWD is outside of the repository, relative paths will be interpreted as relative to *self.path*. Absolute paths will be accepted either way.

```
GIT_ANNEX_MIN_VERSION = '7.20190503'
```

```
WEB_UUID = '00000000-0000-0000-0000-000000000001'
```

```
add(files, git=None, backend=None, options=None, jobs=None, git_options=None, an-
nex_options=None, update=False)
Add file(s) to the repository.
```

### Parameters

- **files** (*list of str*) – list of paths to add to the annex
- **git** (*bool*) – if True, add to git instead of annex.
- **backend** –
- **options** –
- **update** (*bool*) –

–**update option for git-add. From git's manpage:** Update the index just where it already has an entry matching <pathspec>. This removes as well as modifies index entries to match the working tree, but adds no new files.

If no <pathspec> is given when –update option is used, all tracked files in the entire working tree are updated (old versions of Git used to limit the update to the current directory and its subdirectories).

Note: Used only, if a call to git-add instead of git-annex-add is performed

### Returns

**Return type** list of dict or dict

```
add_(files, git=None, backend=None, options=None, jobs=None, git_options=None, an-
nex_options=None, update=False)
Like add, but returns a generator
```

```
add_remote(name, url, options=None)
```

Overrides method from GitRepo in order to set remote.<name>.annex-ssh-options in case of a SSH remote.

```
add_url_to_file(file_, url, options=None, backend=None, batch=False, git_options=None, an-
nex_options=None, unlink_existing=False)
```

Add file from url to the annex.

Downloads *file* from *url* and add it to the annex. If annex knows *file* already, records that it can be downloaded from *url*.

Note: Consider using the higher-level *download\_url* instead.

**Parameters**

- **file** (*str*) –
- **url** (*str*) –
- **options** (*list*) – options to the annex command
- **batch** (*bool*, *optional*) – initiate or continue with a batched run of annex addurl, instead of just calling a single git annex addurl command
- **unlink\_existing** (*bool*, *optional*) – by default crashes if file already exists and is under git. With this flag set to True would first remove it.

**Returns** In batch mode only ATM returns dict representation of json output returned by annex

**Return type** dict

**add\_urls** (*urls*, *options=None*, *backend=None*, *cwd=None*, *jobs=None*, *git\_options=None*, *annex\_options=None*)

Downloads each url to its own file, which is added to the annex.

**Parameters**

- **urls** (*list of str*) –
- **options** (*list*, *optional*) – options to the annex command
- **cwd** (*string*, *optional*) – working directory from within which to invoke git-annex

**adjust** (*options=None*)

enter an adjusted branch

This command is only available in a v6+ git-annex repository.

**Parameters** **options** (*list of str*) – currently requires ‘–unlock’ or ‘–fix’; default: –unlock

**annexstatus** (*paths=None*, *untracked='all'*)

**classmethod check\_direct\_mode\_support** ()

Does git-annex version support direct mode?

The result is cached at *cls.supports\_direct\_mode*.

**Returns**

**Return type** bool

**classmethod check\_repository\_versions** ()

Get information on supported and upgradable repository versions.

The result is cached at *cls.repository\_versions*.

**Returns** supported -> list of supported versions (int) upgradable -> list of upgradable versions (int)

**Return type** dict

**copy\_to** (*files*, *remote*, *options=None*, *jobs=None*)

Copy the actual content of *files* to *remote*

**Parameters**

- **files** (*str or list of str*) – path(s) to copy
- **remote** (*str*) – name of remote to copy *files* to

**Returns** files successfully copied

**Return type** list of str

#### **default\_backends**

**drop** (*files, options=None, key=False, jobs=None*)

Drops the content of annexed files from this repository.

Drops only if possible with respect to required minimal number of available copies.

#### **Parameters**

- **files** (*list of str*) – paths to drop
- **options** (*list of str, optional*) – commandline options for the git annex drop command
- **jobs** (*int, optional*) – how many jobs to run in parallel (passed to git-annex call)

**Returns** ‘success’ item in each object indicates failure/success per file path.

**Return type** list(JSON objects)

**drop\_key** (*keys, options=None, batch=False*)

Drops the content of annexed files from this repository referenced by keys

Dangerous: it drops without checking for required minimal number of available copies.

#### **Parameters**

- **keys** (*list of str, str*) –
- **batch** (*bool, optional*) – initiate or continue with a batched run of annex dropkey, instead of just calling a single git annex dropkey command

**enable\_remote** (*name, env=None*)

Enables use of an existing special remote

**Parameters** **name** (*str*) – name, the special remote was created with

**file\_has\_content** (*files, allow\_quick=True, batch=False*)

Check whether files have their content present under annex.

#### **Parameters**

- **files** (*list of str*) – file(s) to check for being actually present.
- **allow\_quick** (*bool, optional*) – allow quick check, based on having a symlink into `.git/annex/objects`. Works only in non-direct mode (TODO: thin mode)

**Returns** For each input file states whether file has content locally

**Return type** list of bool

**find** (*files, batch=False*)

Run `git annex find` on file(s).

#### **Parameters**

- **files** (*list of str*) – files to find under annex

- **batch** (*bool, optional*) – initiate or continue with a batched run of annex find, instead of just calling a single git annex find command. If any items in *files* are directories, this value is treated as False.

#### Returns

- A dictionary the maps each item in *files* to its *git annex find*
- *result*. Items without a successful result will be an empty string, and
- multi-item results (which can occur for if *files* includes a
- *directory*) will be returned as a list.

**fsck** (*paths=None, remote=None, fast=False, annex\_options=None, git\_options=None*)

Front-end for git-annex fsck

#### Parameters

- **paths** (*list*) – Limit operation to specific paths.
- **remote** (*str*) – If given, the identified remote will be fsck'ed instead of the local repository.
- **fast** (*bool*) – If True, typically means that no actual content is being verified, but tests are limited to the presence of files.

**get** (*files, remote=None, options=None, jobs=None, key=False*)

Get the actual content of files

#### Parameters

- **files** (*list of str*) – paths to get
- **remote** (*str, optional*) – from which remote to fetch content
- **options** (*list of str, optional*) – commandline options for the git annex get command
- **jobs** (*int or None, optional*) – how many jobs to run in parallel (passed to git-annex call). If not specified (None), then
- **key** (*bool, optional*) – If provided file value is actually a key

#### Returns files

**Return type** list of dict

**get\_annexed\_files** (*with\_content\_only=False, patterns=None*)

Get a list of files in annex

#### Parameters

- **with\_content\_only** (*bool, optional*) – Only list files whose content is present.
- **patterns** (*list, optional*) – Globs to pass to annex's *-include=*. Files that match any of these will be returned (i.e., they'll be separated by *-or*).

#### Returns

**Return type** A list of file names

**get\_content\_annexinfo** (*paths=None, init='git', ref=None, eval\_availability=False, key\_prefix="", \*\*kwargs*)

#### Parameters

- **paths** (*list*) – Specific paths to query info for. In none are given, info is reported for all content.
- **init** (*'git' or dict-like or None*) – If set to ‘git’ annex content info will amend the output of `GitRepo.get_content_info()`, otherwise the dict-like object supplied will receive this information and the present keys will limit the report of annex properties. Alternatively, if *None* is given, no initialization is done, and no limit is in effect.
- **ref** (*gitref or None*) – If not *None*, annex content info for this Git reference will be produced, otherwise for the content of the present worktree.
- **eval\_availability** (*bool*) – If this flag is given, evaluate whether the content of any annex’ed file is present in the local annex.
- **\*\*kwargs** – Additional arguments for `GitRepo.get_content_info()`, if *init* is set to ‘git’.

### Returns

Each content item has an entry under its relative path within the repository. Each value is a dictionary with properties:

**type** Can be ‘file’, ‘symlink’, ‘dataset’, ‘directory’

**revision** SHASUM is last commit affecting the item, or *None*, if not tracked.

**key** Annex key of a file (if an annex’ed file)

**bytesize** Size of an annexed file in bytes.

**has\_content** Bool whether a content object for this key exists in the local annex (with *eval\_availability*)

**objloc** `pathlib.Path` of the content object in the local annex, if one is available (with *eval\_availability*)

**Return type** dict

**get\_contentlocation** (*key, batch=False*)

Get location of the key content

Normally under `.git/annex` objects in indirect mode and within file tree in direct mode.

Unfortunately there is no (easy) way to discriminate situations when given key is simply incorrect (not known to annex) or its content not currently present – in both cases annex just silently exits with -1

### Parameters

- **key** (*str*) – key
- **batch** (*bool, optional*) – initiate or continue with a batched run of annex content-location

**Returns** path relative to the top directory of the repository. If no content is present, empty string is returned

**Return type** str

**get\_corresponding\_branch** (*branch=None*)

In case of a managed branch, get the corresponding one.

If *branch* is not a managed branch, return that branch without any changes.

Note: Since default for *branch* is the active branch, `get_corresponding_branch()` is equivalent to `get_active_branch()` if the active branch is not a managed branch.

**Parameters** **branch** (*str*) – name of the branch; defaults to active branch

**Returns** name of the corresponding branch if there is any, name of the queried branch otherwise.

**Return type** str

**get\_description** (*uuid=None*)

Get annex repository description

**Parameters** **uuid** (*str, optional*) – For which remote (based on uuid) to report description for

**Returns** None returned if not found

**Return type** str or None

**get\_file\_backend** (*files*)

Get the backend currently used for file(s).

**Parameters** **files** (*list of str*) –

**Returns** For each file in input list indicates the used backend by a str like “SHA256E” or “MD5”.

**Return type** list of str

**get\_file\_key** (*files, batch=None*)

Get key of an annexed file.

**Parameters**

- **files** (*str or list*) – file(s) to look up
- **batch** (*None or bool, optional*) – If True, *lookupkey -batch* process will be used, which would not crash even if provided file is not under annex (but directly under git), but rather just return an empty string. If False, invokes without *-batch*. If None, use batch mode if more than a single file is provided.

**Returns** keys used by git-annex for each of the files; in case of a list an empty string is returned if there was no key for that file

**Return type** str or list

**Raises**

- `FileInGitError` – If running in non-batch mode and a file is under git, not annex
- `FileNotInAnnexError` – If running in non-batch mode and a file is not under git at all

**get\_file\_size** (*path*)

**get\_groupwanted** (*name*)

Get *groupwanted* expression for a group *name*

**Parameters** **name** (*str*) – Name of the groupwanted group

**classmethod get\_key\_backend** (*key*)

Get the backend from a given key

**get\_metadata** (*files, timestamps=False, batch=False*)

Query git-annex file metadata

**Parameters**

- **files** (*str or iterable(str)*) – One or more paths for which metadata is to be queried. If one or more paths could be directories, *batch=False* must be given to prevent git-annex given an error. Due to technical limitations, such error will lead to a hanging process.

- **timestamps** (*bool, optional*) – If True, the output contains a ‘<metadatakey>-lastchanged’ key for every metadata item, reflecting the modification time, as well as a ‘lastchanged’ key with the most recent modification time of any metadata item.
- **batch** (*bool, optional*) – If True, a *metadata -batch* process will be used, and only confirmed annex’ed files can be queried (else query will hang indefinitely). If False, invokes without *-batch*, and gives all files as arguments (this can be problematic with a large number of files).

**Returns** One tuple per file (could be more items than input arguments when directories are given). First tuple item is the filename, second item is a dictionary with metadata key/value pairs. Note that annex metadata tags are stored under the key ‘tag’, which is a regular metadata item that can be manipulated like any other.

**Return type** generator

**get\_preferred\_content** (*property, remote=None*)

Get preferred content configuration of a repository or remote

**Parameters**

- **property** (*{'wanted', 'required', 'group'}*) – Type of property to query
- **remote** (*str, optional*) – If not specified (None), returns the property for the local repository.

**Returns** Whether the setting is returned, or an empty string if there is none.

**Return type** str

**Raises**

- `ValueError` – If an unknown property label is given.
- `CommandError` – If the annex call errors.

**get\_remotes** (*with\_urls\_only=False, exclude\_special\_remotes=False*)

Get known (special-) remotes of the repository

**Parameters**

- **exclude\_special\_remotes** (*bool, optional*) – if True, don’t return annex special remotes
- **with\_urls\_only** (*bool, optional*) – return only remotes which have urls

**Returns** `remotes` – List of names of the remotes

**Return type** list of str

**static get\_size\_from\_key** (*key*)

A little helper to obtain size encoded in a key

**Returns** size of the file or None if either no size is encoded in the key or key was None itself

**Return type** int or None

**Raises** `ValueError` – if key is considered invalid (at least its size-related part)

**get\_special\_remotes** ()

Get info about all known (not just enabled) special remotes.

**Returns**

Keys are special remote UUIDs. Each value is a dictionary with configuration information git-annex has for the remote. This should include the ‘type’ and ‘name’ as well as any *initremote* parameters that git-annex stores.

Note: This is a faithful translation of git-annex:remote.log with one exception. For a special remote initialized with the `-sameas` flag, git-annex stores the special remote name under the “sameas-name” key, we copy this value under the “name” key so that callers don’t have to check two places for the name. If you need to detect whether you’re working with a sameas remote, the presence of either “sameas-name” or “sameas-uuid” is a reliable indicator.

**Return type** dict

**classmethod** `get_toppath` (*path*, *follow\_up=True*, *git\_options=None*)

Return top-level of a repository given the path.

**Parameters**

- **follow\_up** (*bool*) – If path has symlinks – they get resolved by git. If follow\_up is True, we will follow original path up until we hit the same resolved path. If no such path found, resolved one would be returned.
- **git\_options** (*list of str*) – options to be passed to the git rev-parse call
- **None if no parent directory contains a git repository.**  
(Return) –

**get\_tracking\_branch** (*branch=None*, *remote\_only=False*, *corresponding=True*)

Get the tracking branch for *branch* if there is any.

By default returns the tracking branch of the corresponding branch if *branch* is a managed branch.

**Parameters**

- **branch** (*str*) – local branch to look up. If none is given, active branch is used.
- **remote\_only** (*bool*) – Don’t return a value if the upstream remote is set to “.” (meaning this repository).
- **corresponding** (*bool*) – If True actually look up the corresponding branch of *branch* (also if *branch* isn’t explicitly given)

**Returns** (remote or None, refspec or None) of the tracking branch

**Return type** tuple

**get\_urls** (*file\_*, *key=False*, *batch=False*)

Get URLs for a file/key

**Parameters**

- **file** (*str*) –
- **key** (*bool*, *optional*) – Whether provided files are actually annex keys

**Returns**

**Return type** A list of URLs

**git\_annex\_version = None**

**info** (*files*, *batch=False*, *fast=False*)

Provide annex info for file(s).

**Parameters** **files** (*list of str*) – files to look for

**Returns** Info for each file



**Return type** dict

**init\_remote** (*name, options*)

Creates a new special remote

**Parameters** **name** (*str*) – name of the special remote

**is\_available** (*file\_, remote=None, key=False, batch=False*)

Check if file or key is available (from a remote)

In case if key or remote is misspecified, it wouldn't fail but just keep returning False, although possibly also complaining out loud ;)

**Parameters**

- **file** (*str*) – Filename or a key
- **remote** (*str, optional*) – Remote which to check. If None, possibly multiple remotes are checked before positive result is reported
- **key** (*bool, optional*) – Whether provided files are actually annex keys
- **batch** (*bool, optional*) – Initiate or continue with a batched run of annex checkpresentkey

**Returns** with True indicating that file/key is available from (the) remote

**Return type** bool

**is\_crippled\_fs** ()

Return True if git-annex considers current filesystem 'crippled'.

**Returns**

**Return type** True if on crippled filesystem, False otherwise

**is\_direct\_mode** ()

Return True if annex is in direct mode

**Returns**

**Return type** True if in direct mode, False otherwise.

**is\_initialized** ()

quick check whether this appears to be an annex-init'ed repo

**is\_managed\_branch** (*branch=None*)

Whether *branch* is managed by git-annex.

ATM this returns true in direct mode (branch 'annex/direct/my\_branch') and if on an adjusted branch (annex v6+ repository: either 'adjusted/my\_branch(unlocked)' or 'adjusted/my\_branch(fixed)')

Note: The term 'managed branch' is used to make clear it's meant to be more general than the v6+ 'adjusted branch'.

**Parameters** **branch** (*str*) – name of the branch; default: active branch

**Returns** True if on a managed branch, False otherwise

**Return type** bool

**is\_remote\_annex\_ignored** (*remote*)

Return True if remote is explicitly ignored

**is\_special\_annex\_remote** (*remote, check\_if\_known=True*)

Return whether remote is a special annex remote

Decides based on the presence of an annex- option and lack of a configured URL for the remote.

**is\_under\_annex** (*files*, *allow\_quick=True*, *batch=False*)

Check whether files are under annex control

**Parameters**

- **files** (*list of str*) – file(s) to check for being under annex
- **allow\_quick** (*bool, optional*) – allow quick check, based on having a symlink into .git/annex/objects. Works only in non-direct mode (TODO: thin mode)

**Returns** For each input file states whether file is under annex

**Return type** list of bool

**is\_valid\_annex** (*allow\_noninitialized=False*, *check\_git=True*)

Returns whether the underlying repository appears to be still valid

Note, that this almost identical to the classmethod `is_valid_repo()`. However, if we are testing an existing instance, we can save Path object creations. Since this testing is done a lot, this is relevant. Creation of the Path objects in `is_valid_repo()` takes nearly half the time of the entire function.

Also note, that this method is bound to an instance but still class-dependent, meaning that a subclass cannot simply overwrite it. This is particularly important for the call from within `__init__()`, which in turn is called by the subclasses' `__init__`. Using an overwrite would lead to the wrong thing being called.

**classmethod is\_valid\_repo** (*path*, *allow\_noninitialized=False*)

Return True if given path points to an annex repository

**merge\_annex** (*remote=None*)

Merge git-annex branch

Merely calls `sync` with the appropriate arguments.

**Parameters remote** (*str, optional*) – Name of a remote to be “merged”.

**migrate\_backend** (*files*, *backend=None*)

Changes the backend used for *file*.

The backend used for the key-value of *files*. Only files currently present are migrated. Note: There will be no notification if migrating fails due to the absence of a file’s content!

**Parameters**

- **files** (*list*) – files to migrate.
- **backend** (*str*) – specify the backend to migrate to. If none is given, the default backend of this instance will be used.

**precommit** ()

Perform pre-commit maintenance tasks, such as closing all batched annexes since they might still need to flush their changes into index

**remove** (*files*, *force=False*, *\*\*kwargs*)

Remove files from git/annex

**Parameters**

- **files** –
- **force** (*bool, optional*) –

**repo\_info** (*fast=False*, *merge\_annex\_branches=True*)

Provide annex info for the entire repository.

**Parameters**

- **fast** (*bool, optional*) – Pass `-fast` to `git annex info`.
- **merge\_annex\_branches** (*bool, optional*) – Whether to allow `git-annex` if needed to merge annex branches, e.g. to make sure up to date descriptions for `git annex remotes`

**Returns** Info for the repository, with keys matching the ones returned by `annex`

**Return type** dict

**repository\_versions** = None

**rm\_url** (*file\_, url*)

Record that the file is no longer available at the url.

**Parameters**

- **file** (*str*) –
- **url** (*str*) –

**set\_default\_backend** (*backend, persistent=True, commit=True*)

Set default backend

**Parameters**

- **backend** (*str*) –
- **persistent** (*bool, optional*) – If persistent, would add/commit to `.gitattributes`. If not – would set within `.git/config`

**set\_groupwanted** (*name, expr*)

Set *expr* for the *name* groupwanted

**set\_metadata** (*files, reset=None, add=None, init=None, remove=None, purge=None, recursive=False*)

Manipulate `git-annex` file-metadata

**Parameters**

- **files** (*str or list(str)*) – One or more paths for which metadata is to be manipulated. The changes applied to each file item are uniform. However, the result may not be uniform across files, depending on the actual operation.
- **reset** (*dict, optional*) – Metadata items matching keys in the given dict are (re)set to the respective values.
- **add** (*dict, optional*) – The values of matching keys in the given dict appended to any possibly existing values. The metadata keys need not necessarily exist before.
- **init** (*dict, optional*) – Metadata items for the keys in the given dict are set to the respective values, if the key is not yet present in a file’s metadata.
- **remove** (*dict, optional*) – Values in the given dict are removed from the metadata items matching the respective key, if they exist in a file’s metadata. Non-existing values, or keys do not lead to failure.
- **purge** (*list, optional*) – Any metadata item with a key matching an entry in the given list is removed from the metadata.
- **recursive** (*bool, optional*) – If False, fail (with `CommandError`) when directory paths are given as *files*.

**Returns** JSON obj per modified file

**Return type** list

**set\_metadata\_** (*files*, *reset=None*, *add=None*, *init=None*, *remove=None*, *purge=None*, *recursive=False*)

Like `set_metadata()` but returns a generator

**set\_preferred\_content** (*property*, *expr*, *remote=None*)

Set preferred content configuration of a repository or remote

**Parameters**

- **property** (*{'wanted', 'required', 'group'}*) – Type of property to query
- **expr** (*str*) – Any expression or label supported by git-annex for the given property.
- **remote** (*str*, *optional*) – If not specified (`None`), sets the property for the local repository.

**Returns** Raw git-annex output in response to the set command.

**Return type** str

**Raises**

- `ValueError` – If an unknown property label is given.
- `CommandError` – If the annex call errors.

**set\_remote\_dead** (*name*)

Announce to annex that remote is “dead”

**set\_remote\_url** (*name*, *url*, *push=False*)

Set the URL a remote is pointing to

Sets the URL of the remote *name*. Requires the remote to already exist.

**Parameters**

- **name** (*str*) – name of the remote
- **url** (*str*) –
- **push** (*bool*) – if `True`, set the push URL, otherwise the fetch URL; if `True`, additionally set `annexurl` to *url*, to make sure annex uses it to talk to the remote, since access via fetch URL might be restricted.

**supports\_direct\_mode** = `None`

**supports\_unlocked\_pointers**

Return `True` if repository version supports unlocked pointers.

**sync** (*remotes=None*, *push=True*, *pull=True*, *commit=True*, *content=False*, *all=False*, *fast=False*)

Synchronize local repository with remotes

Use this command when you want to synchronize the local repository with one or more of its remotes. You can specify the remotes (or remote groups) to sync with by name; the default if none are specified is to sync with all remotes.

**Parameters**

- **remotes** (*str*, *list(str)*, *optional*) – Name of one or more remotes to be sync'ed.
- **push** (*bool*) – By default, git pushes to remotes.
- **pull** (*bool*) – By default, git pulls from remotes

- **commit** (*bool*) – A commit is done by default. Disable to avoid committing local changes.
- **content** (*bool*) – Normally, syncing does not transfer the contents of annexed files. This option causes the content of files in the work tree to also be uploaded and downloaded as necessary.
- **all** (*bool*) – This option, when combined with *content*, makes all available versions of all files be synced, when preferred content settings allow
- **fast** (*bool*) – Only sync with the remotes with the lowest annex-cost value configured

**unannex** (*files, options=None*)  
undo accidental add command

Use this to undo an accidental git annex add command. Note that for safety, the content of the file remains in the annex, until you use git annex unused and git annex dropunused.

**Parameters**

- **files** (*list of str*) –
- **options** (*list of str*) –

**Returns** successfully unannexed files

**Return type** list of str

**unlock** (*files*)  
unlock files for modification

Note: This method is silent about errors in unlocking a file (e.g, the file has not content). Use the higher-level interface `unlock` to get more informative reporting.

**Parameters** **files** (*list of str*) –

**Returns** successfully unlocked files

**Return type** list of str

**uuid**  
Annex UUID

**Returns** Returns a the annex UUID, if there is any, or *None* otherwise.

**Return type** str

**whereis** (*files, output='uuids', key=False, options=None, batch=False*)  
Lists repositories that have actual content of file(s).

**Parameters**

- **files** (*list of str*) – files to look for
- **output** (*{'descriptions', 'uuids', 'full'}, optional*) – If ‘descriptions’, a list of remotes descriptions returned is per each file. If ‘full’, for each file a dictionary of all fields is returned as returned by annex
- **key** (*bool, optional*) – Whether provided files are actually annex keys
- **options** (*list, optional*) – Options to pass into git-annex call

**Returns**

if output == ‘descriptions’, contains a list of descriptions of remotes for each input file, describing the remote for each remote, which was found by git-annex whereis, like:

```
u'me@mycomputer:~/where/my/repo/is [origin]' or
u'web' or
u'me@mycomputer:~/some/other/clone'
```

if output == 'uuids', returns a list of uuids. if output == 'full', returns a dictionary with filenames as keys and values a detailed record, e.g.:

```
{'00000000-0000-0000-0000-000000000001': {
  'description': 'web',
  'here': False,
  'urls': ['http://127.0.0.1:43442/about.txt', 'http://example.com/
↪someurl']
}}
```

**Return type** list of list of unicode or dict

**class** `datalad.support.annexrepo.BatchedAnnex` (*annex\_cmd*, *git\_options=None*, *annex\_options=None*, *path=None*, *json=False*, *output\_proc=None*)

Bases: `datalad.cmd.BatchedCommand`

Container for an annex process which would allow for persistent communication

**class** `datalad.support.annexrepo.BatchedAnnexes` (*batch\_size=0*, *git\_options=None*)

Bases: `datalad.cmd.SafeDelCloseMixin`, dict

Class to contain the registry of active batch'ed instances of annex for a repository

**clear** ()

Override just to make sure we don't rely on `__del__` to close all the pipes

**close** ()

Close communication to all the batched annexes

It does not remove them from the dictionary though

**get** (*codename*, *annex\_cmd=None*, *\*\*kwargs*)

Return the value for key if key is in the dictionary, else default.

**class** `datalad.support.annexrepo.ProcessAnnexProgressIndicators` (*expected=None*)

Bases: object

'Filter' for annex `-json` output to react to progress indicators

Instance of this beast should be passed into `log_stdout` option for git-annex commands runner

**finish** (*partial=False*)

**start** ()

`datalad.support.annexrepo.readline_json` (*stdout*)

`datalad.support.annexrepo.readlines_until_ok_or_failed` (*stdout*, *maxlines=100*)

Read stdout until line ends with ok or failed

## datalad.support.archives

Various handlers/functionality for different types of files (e.g. for archives)

**class** `datalad.support.archives.ArchivesCache` (*toppath=None*, *persistent=False*)

Bases: object

Cache to maintain extracted archives

#### Parameters

- **toppath** (*str*) – Top directory under `.git/` of which temp directory would be created. If not provided – random tempdir is used
- **persistent** (*bool, optional*) – Passed over into generated `ExtractedArchives`

**clean** (*force=False*)

**get\_archive** (*archive*)

**path**

```
class datalad.support.archives.ExtractedArchive (archive, path=None, persis-
tent=False)
```

Bases: object

Container for the extracted archive

**STAMP\_SUFFIX** = `' .stamp'`

**assure\_extracted** ()

Return path to the extracted *archive*. Extract archive if necessary

**clean** (*force=False*)

**get\_extracted\_file** (*afile*)

**get\_extracted\_filename** (*afile*)

Return full path to the *afile* within extracted *archive*

It does not actually extract any archive

**get\_extracted\_files** ()

Generator to provide filenames which are available under extracted archive

**get\_leading\_directory** (*depth=None, consider=None, exclude=None*)

Return leading directory of the content within archive

#### Parameters

- **depth** (*int or None, optional*) – Maximal depth of leading directories to consider. If None - no upper limit
- **consider** (*list of str, optional*) – Regular expressions for file/directory names to be considered (before exclude). Applied to the entire relative path to the file as in the archive
- **exclude** (*list of str, optional*) – Regular expressions for file/directory names to be excluded from consideration. Applied to the entire relative path to the file as in the archive

**Returns** If there is no single leading directory – None returned

**Return type** `str` or `None`

**is\_extracted**

**path**

Given an archive – return full path to it within cache (extracted)

**stamp\_path**

```
datalad.support.archives.decompress_file (archive, dir_, leading_directories='strip')
```

Decompress *archive* into a directory *dir\_*

### Parameters

- **archive** (*str*) –
- **dir** (*str*) –
- **leading\_directories** (*{'strip', None}*) – If *strip*, and archive contains a single leading directory under which all content is stored, all the content will be moved one directory up and that leading directory will be removed.

### datalad.support.configparserinc

**class** datalad.support.configparserinc.**SafeConfigParserWithIncludes** (*\*args*, *\*\*kwargs*)

Bases: configparser.ConfigParser

Class adds functionality to SafeConfigParser to handle included other configuration files (or may be urls, whatever in the future)

File should have section [includes] and only 2 options implemented are 'files\_before' and 'files\_after' where files are listed 1 per line.

Example:

```
[INCLUDES]
before = 1.conf
        3.conf

after = 1.conf
```

It is a simple implementation, so just basic care is taken about recursion. Includes preserve right order, ie new files are inserted to the list of read configs before original, and their includes correspondingly so the list should follow the leaves of the tree.

I wasn't sure what would be the right way to implement generic (aka c++ template) so we could base at any *\*configparser* class... so I will leave it for the future

**SECTION\_NAME = 'INCLUDES'**

**static** **getIncludes** (*resource*, *seen=[]*)

Given 1 config resource returns list of included files (recursively) with the original one as well Simple loops are taken care about

**read** (*filenames*)

Read and parse a filename or an iterable of filenames.

Files that cannot be opened are silently ignored; this is designed so that you can specify an iterable of potential configuration file locations (e.g. current directory, user's home directory, systemwide directory), and all existing configuration files in the iterable will be read. A single filename may also be given.

Return list of successfully read files.

### datalad.customremotes.main

datalad.customremotes.main.**main** (*args=None*, *backend=None*)

datalad.customremotes.main.**setup\_parser** (*backend*)



## datalad.customremotes.base

Base classes to custom git-annex remotes (e.g. extraction from archives)

**class** `datalad.customremotes.base.AnnexCustomRemote` (*path=None, cost=None, fin=None, fout=None*)

Bases: `object`

Base class to provide custom special remotes for git-annex

Implements git-annex special custom remotes protocol described at [http://git-annex.branchable.com/design/external\\_special\\_remote\\_protocol/](http://git-annex.branchable.com/design/external_special_remote_protocol/)

**AVAILABILITY** = 'LOCAL'

**COST** = 100

**CUSTOM\_REMOTE\_NAME** = None

**SUPPORTED\_SCHEMES** = ()

**debug** (*msg*)

**error** (*msg, annex\_err='ERROR'*)

**get\_DIRHASH** (*key, full=False*)

Gets a two level hash associated with a Key.

### Parameters

- **full** (*bool, optional*) – If True, would spit out full DIRHASH path, i.e. with a KEY/ directory
- **like "abc/def". This is always the same for any given Key, so (Something)** –
- **be used for eg, creating hash directory structures to store Keys in. (can)** –

**get\_URLS** (*key*)

Gets URL(s) associated with a Key.

**get\_contentlocation** (*key, absolute=False, verify\_exists=True*)

Return (relative to top or absolute) path to the file containing the key

This is a wrapper around `AnnexRepo.get_contentlocation` which provides caching of the result (we are asking the location for the same archive key often)

**heavydebug** (*msg, \*args, \*\*kwargs*)

**info** (*msg*)

**main** ()

Interface to the command line tool

**progress** (*bytes*)

**read** (*req=None, n=1*)

Read a message from git-annex

### Parameters

- **req** (*string, optional*) – Expected request - first msg of the response
- **n** (*int*) – Number of response elements after first msg

**req\_CHECKPRESENT** (*key*)

**CHECKPRESENT-SUCCESS Key** Indicates that a key has been positively verified to be present in the remote.

**CHECKPRESENT-FAILURE Key** Indicates that a key has been positively verified to not be present in the remote.

**CHECKPRESENT-UNKNOWN Key ErrorMessage** Indicates that it is not currently possible to verify if the key is present in the remote. (Perhaps the remote cannot be contacted.)

**req\_CHECKURL** (*url*)

The remote replies with one of CHECKURL-FAILURE, CHECKURL-CONTENTS, or CHECKURL-MULTI.

**CHECKURL-CONTENTS Size|UNKNOWN Filename** Indicates that the requested url has been verified to exist. The Size is the size in bytes, or use “UNKNOWN” if the size could not be determined. The Filename can be empty (in which case a default is used), or can specify a filename that is suggested to be used for this url.

**CHECKURL-MULTI Url Size|UNKNOWN Filename ...** Indicates that the requested url has been verified to exist, and contains multiple files, which can each be accessed using their own url. Note that since a list is returned, neither the Url nor the Filename can contain spaces.

**CHECKURL-FAILURE** Indicates that the requested url could not be accessed.

**req\_CLAIMURL** (*url*)

**req\_EXPORTSUPPORTED** ()

**req\_GETAVAILABILITY** ()

**req\_GETCOST** ()

**req\_INITREMOTE** (*\*args*)

Initialize this remote. Provides high level abstraction.

Specific implementation should go to `_initialize`

**req\_PREPARE** (*\*args*)

Prepare “to deliver”. Provides high level abstraction

Specific implementation should go to `_prepare`

**req\_REMOVE** (*key*)

**REMOVE-SUCCESS Key** Indicates the key has been removed from the remote. May be returned if the remote didn’t have the key at the point removal was requested.

**REMOVE-FAILURE Key ErrorMessage** Indicates that the key was unable to be removed from the remote.

**req\_TRANSFER** (*cmd, key, file*)

**req\_WHEREIS** (*key*)

Added in 5.20150812-17-g6bc46e3

provide any information about ways to access the content of a key stored in it, such as eg, public urls. This will be displayed to the user by eg, `git annex whereis`. The remote replies with WHEREIS-SUCCESS or WHEREIS-FAILURE. Note that users expect `git annex whereis` to run fast, without eg, network access. This is not needed when SETURIPRESENT is used, since such uris are automatically displayed by `git annex whereis`.

**WHEREIS-SUCCESS String** Indicates a location of a key. Typically an url, the string can be anything that it makes sense to display to the user about content stored in the special remote.

**WHEREIS-FAILURE** Indicates that no location is known for a key.

**send** (\*args)

Send a message to git-annex

**Parameters** \*args (*list of strings*) – arguments to be joined by a space and passed to git-annex

**send\_unsupported** (msg=None)

Send UNSUPPORTED-REQUEST to annex and log optional message in our log

**stop** (msg=None)

**class** datalad.customremotes.base.**AnnexExchangeProtocol** (repopath, custom\_remote\_name=None)

Bases: datalad.support.protocol.ProtocolInterface

A little helper to protocol interactions of custom remote with annex

**HEADER** = '#!/bin/bash\n\nset -e\n\n# Gets a VALUE response and stores it in \$RET\nrepor

**add\_section** (cmd, exception)

Adds a section to the protocol.

This is an alternative to the use of start\_section() and end\_section(). In opposition to start\_section, this one can be called anytime.

**Parameters**

- **cmd** (*list*) – The actual command and its options/arguments as a list
- **exception** (*Exception*) – The exception raised by the command if any or None otherwise.

**do\_execute\_callables**

Indicates whether or not the callables are supposed to actually be executed.

**Returns**

**Return type** bool

**do\_execute\_ext\_commands**

Indicates whether or not the called commands are supposed to actually be executed.

**Returns**

**Return type** bool

**end\_section** (id\_, exception)

Ends the section *id*.

To call after the command call to be recorded. This ends the section defined by *id* as returned by start\_section().

**Parameters**

- **id** (*int*) –
- **exception** (*Exception*) – The exception raised by the command if any or None otherwise.
- **Raises** –
- ----- –
- **IndexError** – in case *id* is invalid.

**initiate** ()

**records\_callables**

Indicates whether or not the protocol is supposed to include calls of python callables.

**Returns**

**Return type** bool

**records\_ext\_commands**

Indicates whether or not the protocol is supposed to include external command calls.

**Returns**

**Return type** bool

**start\_section** (*cmd*)

Starts a new section of the protocol.

To call before the command call to be recorded. To be used with a corresponding call of `end_section()`.

**Parameters** *cmd* (*list*) – The actual command and its options/arguments as a list

**Returns** An id of the started section to be used as argument of the corresponding call of `end_section()`.

**Return type** int

**write\_entries** (*entries*)

**write\_section** (*cmd*)

**exception** `datalad.customremotes.base.AnnexRemoteQuit`

Bases: `Exception`

`datalad.customremotes.base.generate_uuids()`

Generate UUIDs for our remotes. Even though quick, for consistency pre-generated and recorded in `consts.py`

`datalad.customremotes.base.get_function_nargs(f)`

`datalad.customremotes.base.init_datalad_remote(repo, remote, encryption=None, autoenable=False, opts=[])`

Initialize datalad special remote

**datalad.customremotes.archives**

Custom remote to support getting the load from archives present under annex

**class** `datalad.customremotes.archives.ArchiveAnnexCustomRemote` (*persistent\_cache=True, \*\*kwargs*)

Bases: `datalad.customremotes.base.AnnexCustomRemote`

Special custom remote allowing to obtain files from archives

Archives should also be under annex control.

**AVAILABILITY** = 'local'

**COST** = 500

**CUSTOM\_REMOTE\_NAME** = 'archive'

**SUPPORTED\_SCHEMES** = ('dl+archive',)

**URL\_PREFIX** = 'dl+archive:'

**URL\_SCHEME** = 'dl+archive'

**cache**

**get\_file\_url** (*archive\_file=None, archive\_key=None, file=None, size=None*)  
 Given archive (file or a key) and a file – compose URL for access

### Examples

**dl+archive:SHA256E-s176-69...3e.tar.gz#path=1/d2/2d&size=123** when size of file within archive was known to be 123

**dl+archive:SHA256E-s176-69...3e.tar.gz#path=1/d2/2d** when size of file within archive was not provided

**Parameters** **size** (*int, optional*) – Size of the file. If not provided, will simply be empty

**req\_CHECKPRESENT** (*key*)

Check if copy is available

TODO: just proxy the call to annex for underlying tarball

Replies

**CHECKPRESENT-SUCCESS Key** Indicates that a key has been positively verified to be present in the remote.

**CHECKPRESENT-FAILURE Key** Indicates that a key has been positively verified to not be present in the remote.

**CHECKPRESENT-UNKNOWN Key ErrorMessage** Indicates that it is not currently possible to verify if the key is present in the remote. (Perhaps the remote cannot be contacted.)

**req\_CHECKURL** (*url*)

Replies

**CHECKURL-CONTENTS Size|UNKNOWN Filename** Indicates that the requested url has been verified to exist. The Size is the size in bytes, or use “UNKNOWN” if the size could not be determined. The Filename can be empty (in which case a default is used), or can specify a filename that is suggested to be used for this url.

**CHECKURL-MULTI Url Size|UNKNOWN Filename ...** Indicates that the requested url has been verified to exist, and contains multiple files, which can each be accessed using their own url. Note that since a list is returned, neither the Url nor the Filename can contain spaces.

**CHECKURL-FAILURE** Indicates that the requested url could not be accessed.

**req\_REMOVE** (*key*)

**REMOVE-SUCCESS Key** Indicates the key has been removed from the remote. May be returned if the remote didn't have the key at the point removal was requested

**REMOVE-FAILURE Key ErrorMessage** Indicates that the key was unable to be removed from the remote.

**req\_WHEREIS** (*key*)

**WHEREIS-SUCCESS String** Indicates a location of a key. Typically an url, the string can be anything that it makes sense to display to the user about content stored in the special remote.

**WHEREIS-FAILURE** Indicates that no location is known for a key.

**stop** (*\*args*)

Stop communication with annex

```
datalad.customremotes.archives.link_file_load(src, dst, dry_run=False)
```

Just a little helper to hardlink files's load

```
datalad.customremotes.archives.main()
```

cmdline entry point

## Configuration management

---

*config*

---

### datalad.config

```
class datalad.config.ConfigManager(dataset=None, dataset_only=False, overrides=None,  
                                     source='any')
```

Bases: object

Thin wrapper around *git-config* with support for a dataset configuration.

The general idea is to have an object that is primarily used to read/query configuration option. Upon creation, current configuration is read via one (or max two, in the case of the presence of dataset-specific configuration) calls to *git config*. If this class is initialized with a Dataset instance, it supports reading and writing configuration from `.datalad/config` inside a dataset too. This file is committed to Git and hence useful to ship certain configuration items with a dataset.

The API aims to provide the most significant read-access API of a dictionary, the Python ConfigParser, and GitPython's config parser implementations.

This class is presently not capable of efficiently writing multiple configurations items at once. Instead, each modification results in a dedicated call to *git config*. This author thinks this is OK, as he cannot think of a situation where a large number of items need to be written during normal operation. If such need arises, various solutions are possible (via GitPython, or an independent writer).

Each instance carries a public *overrides* attribute. This dictionary contains variables that override any setting read from a file. The overrides are persistent across reloads, and are not modified by any of the manipulation methods, such as *set* or *unset*.

Any DATALAD\_\* environment variable is also presented as a configuration item. Settings read from environment variables are not stored in any of the configuration file, but are read dynamically from the environment at each *reload()* call. Their values take precedence over any specification in configuration files, and even overrides.

#### Parameters

- **dataset** (*Dataset*, *optional*) – If provided, all *git config* calls are executed in this dataset's directory. Moreover, any modifications are, by default, directed to this dataset's configuration file (which will be created on demand)
- **dataset\_only** (*bool*) – Legacy option, do not use.
- **overrides** (*dict*, *optional*) – Variable overrides, see general class documentation for details.
- **source** (`{'any', 'local', 'dataset', 'dataset-local'}`, *optional*) – Which sources of configuration setting to consider. If 'dataset', configuration items are only read from a dataset's persistent configuration file, if any is present (the one in `.datalad/config`, not `.git/config`); if 'local', any non-committed source is considered (local and global configuration in Git config's terminology); if 'dataset-local', persistent dataset configuration and local, but not global or system configuration are considered; if 'any' all possible sources of configuration are considered.

**add** (*var*, *value*, *where*='dataset', *reload*=True)  
 Add a configuration variable and value

#### Parameters

- **var** (*str*) – Variable name including any section like *git config* expects them, e.g. 'core.editor'
- **value** (*str*) – Variable value
- **where** ({'dataset', 'local', 'global', 'override'}, *optional*) – Indicator which configuration file to modify. 'dataset' indicates the persistent configuration in .datalad/config of a dataset; 'local' the configuration of a dataset's Git repository in .git/config; 'global' refers to the general configuration that is not specific to a single repository (usually in \$USER/.gitconfig); 'override' limits the modification to the Config-Manager instance, and the assigned value overrides any setting from any other source.
- **reload** (*bool*) – Flag whether to reload the configuration from file(s) after modification. This can be disable to make multiple sequential modifications slightly more efficient.

**get** (*k*, *d*) → D[k] if k in D, else d. d defaults to None.

**get\_value** (*section*, *option*, *default*=None)  
 Like *get()*, but with an optional default value

If the default is not None, the given default value will be returned in case the option did not exist. This behavior imitates GitPython's config parser.

**getbool** (*section*, *option*, *default*=None)  
 A convenience method which coerces the option value to a bool

Values "on", "yes", "true" and any int!=0 are considered True Values which evaluate to bool False, "off", "no", "false" are considered False TypeError is raised for other values.

**getfloat** (*section*, *option*)  
 A convenience method which coerces the option value to a float

**getint** (*section*, *option*)  
 A convenience method which coerces the option value to an integer

**has\_option** (*section*, *option*)  
 If the given section exists, and contains the given option

**has\_section** (*section*)  
 Indicates whether a section is present in the configuration

**items** (*section*=None)  
 Return a list of (name, value) pairs for each option  
 Optionally limited to a given section.

**keys** ()  
 Returns list of configuration item names

**obtain** (*var*, *default*=None, *dialog\_type*=None, *valtype*=None, *store*=False, *where*=None, *reload*=True, *\*\*kwargs*)  
 Convenience method to obtain settings interactively, if needed

A UI will be used to ask for user input in interactive sessions. Questions to ask, and additional explanations can be passed directly as arguments, or retrieved from a list of pre-configured items.

Additionally, this method allows for type conversion and storage of obtained settings. Both aspects can also be pre-configured.

#### Parameters

- **var** (*str*) – Variable name including any section like *git config* expects them, e.g. ‘core.editor’
- **default** (*any type*) – In interactive sessions and if *store* is True, this default value will be presented to the user for confirmation (or modification). In all other cases, this value will be silently assigned unless there is an existing configuration setting.
- **dialog\_type** (*{'question', 'yesno', None}*) – Which dialog type to use in interactive sessions. If *None*, pre-configured UI options are used.
- **store** (*bool*) – Whether to store the obtained value (or default)
- **where** (*{'dataset', 'local', 'global', 'override'}, optional*) – Indicator which configuration file to modify. ‘dataset’ indicates the persistent configuration in *.datalad/config* of a dataset; ‘local’ the configuration of a dataset’s Git repository in *.git/config*; ‘global’ refers to the general configuration that is not specific to a single repository (usually in *\$USER/.gitconfig*); ‘override’ limits the modification to the Config-Manager instance, and the assigned value overrides any setting from any other source.
- **reload** (*bool*) – Flag whether to reload the configuration from file(s) after modification. This can be disable to make multiple sequential modifications slightly more efficient.
- **\*\*kwargs** – Additional arguments for the UI function call, such as a question *text*.

**options** (*section*)

Returns a list of options available in the specified section.

**reload** (*force=False*)

Reload all configuration items from the configured sources

If *force* is False, all files configuration was previously read from are checked for differences in the modification times. If no difference is found for any file no reload is performed. This mechanism will not detect newly created global configuration files, use *force* in this case.

**remove\_section** (*sec, where='dataset', reload=True*)

Rename a configuration section

**Parameters**

- **sec** (*str*) – Name of the section to remove.
- **where** (*{'dataset', 'local', 'global', 'override'}, optional*) – Indicator which configuration file to modify. ‘dataset’ indicates the persistent configuration in *.datalad/config* of a dataset; ‘local’ the configuration of a dataset’s Git repository in *.git/config*; ‘global’ refers to the general configuration that is not specific to a single repository (usually in *\$USER/.gitconfig*); ‘override’ limits the modification to the Config-Manager instance, and the assigned value overrides any setting from any other source.
- **reload** (*bool*) – Flag whether to reload the configuration from file(s) after modification. This can be disable to make multiple sequential modifications slightly more efficient.

**rename\_section** (*old, new, where='dataset', reload=True*)

Rename a configuration section

**Parameters**

- **old** (*str*) – Name of the section to rename.
- **new** (*str*) – Name of the section to rename to.
- **where** (*{'dataset', 'local', 'global', 'override'}, optional*) – Indicator which configuration file to modify. ‘dataset’ indicates the persistent configuration in *.datalad/config* of a dataset; ‘local’ the configuration of a dataset’s Git repository



in `.git/config`; ‘global’ refers to the general configuration that is not specific to a single repository (usually in `$USER/.gitconfig`); ‘override’ limits the modification to the Config-Manager instance, and the assigned value overrides any setting from any other source.

- **reload** (*bool*) – Flag whether to reload the configuration from file(s) after modification. This can be disabled to make multiple sequential modifications slightly more efficient.

**rewrite\_url** (*url*)

Any matching ‘url.<base>.insteadOf’ configuration is applied

Any URL that starts with such a configuration will be rewritten to start, instead, with <base>. When more than one insteadOf strings match a given URL, the longest match is used.

#### Parameters

- **cfg** (*ConfigManager* or *dict*) – dict-like with configuration variable name/value-pairs.
- **url** (*str*) – URL to be rewritten, if matching configuration is found.

**Returns** Rewritten or unmodified URL.

**Return type** *str*

**sections** ()

Returns a list of the sections available

**set** (*var*, *value*, *where*='dataset', *reload*=True, *force*=False)

Set a variable to a value.

In opposition to *add*, this replaces the value of *var* if there is one already.

#### Parameters

- **var** (*str*) – Variable name including any section like *git config* expects them, e.g. ‘core.editor’
- **value** (*str*) – Variable value
- **force** (*bool*) – if set, replaces all occurrences of *var* by a single one with the given *value*. Otherwise raise if multiple entries for *var* exist already
- **where** (*{'dataset', 'local', 'global', 'override'}, optional*) – Indicator which configuration file to modify. ‘dataset’ indicates the persistent configuration in `.datalad/config` of a dataset; ‘local’ the configuration of a dataset’s Git repository in `.git/config`; ‘global’ refers to the general configuration that is not specific to a single repository (usually in `$USER/.gitconfig`); ‘override’ limits the modification to the Config-Manager instance, and the assigned value overrides any setting from any other source.
- **reload** (*bool*) – Flag whether to reload the configuration from file(s) after modification. This can be disabled to make multiple sequential modifications slightly more efficient.

**unset** (*var*, *where*='dataset', *reload*=True)

Remove all occurrences of a variable

#### Parameters

- **var** (*str*) – Name of the variable to remove
- **where** (*{'dataset', 'local', 'global', 'override'}, optional*) – Indicator which configuration file to modify. ‘dataset’ indicates the persistent configuration in `.datalad/config` of a dataset; ‘local’ the configuration of a dataset’s Git repository in `.git/config`; ‘global’ refers to the general configuration that is not specific to a single

repository (usually in \$USER/.gitconfig); ‘override’ limits the modification to the Config-Manager instance, and the assigned value overrides any setting from any other source.

- **reload** (*bool*) – Flag whether to reload the configuration from file(s) after modification. This can be disable to make multiple sequential modifications slightly more efficient.

`datalad.config.anything2bool` (*val*)

`datalad.config.get_git_version` (*runner*)

Return version of available git

`datalad.config.rewrite_url` (*cfg, url*)

Any matching ‘url.<base>.insteadOf’ configuration is applied

Any URL that starts with such a configuration will be rewritten to start, instead, with <base>. When more than one insteadOf strings match a given URL, the longest match is used.

**Parameters**

- **cfg** (*ConfigManager* or *dict*) – dict-like with configuration variable name/value-pairs.
- **url** (*str*) – URL to be rewritten, if matching configuration is found.

**Returns** Rewritten or unmodified URL.

**Return type** *str*

**Test infrastructure**

---

|                              |   |
|------------------------------|---|
| <i>tests.utils</i>           | Miscellaneous utilities to assist with testing    |
| <i>tests.utils.testrepos</i> |   |
| <i>tests.heavyoutput</i>     | Helper to provide heavy load on stdout and stderr |

---

**datalad.tests.utils**

Miscellaneous utilities to assist with testing

**class** `datalad.tests.utils.HTTPPath` (*path*)

Bases: *object*

Serve the content of a path via an HTTP URL.

This class can be used as a context manager, in which case it returns the URL.

Alternatively, the *start* and *stop* methods can be called directly.

**Parameters** **path** (*str*) – Directory with content to serve.

**start** ()

Start serving *path* via HTTP.

**stop** ()

Stop serving *path*.

**class** `datalad.tests.utils.SilentHTTPHandler` (*\*args, \*\*kwargs*)

Bases: *http.server.SimpleHTTPRequestHandler*

A little adapter to silence the handler

**log\_message** (*format, \*args*)

Log an arbitrary message.

This is used by all other logging functions. Override it if you have specific logging wishes.

The first argument, `FORMAT`, is a format string for the message to be logged. If the format string contains any `%` escapes requiring parameters, they should be specified as subsequent arguments (it's just like `printf!`).

The client ip and current date/time are prefixed to every message.

`datalad.tests.utils.assert_dict_equal` (*d1*, *d2*)

`datalad.tests.utils.assert_in_results` (*results*, *\*\*kwargs*)

Verify that the particular combination of keys and values is found in one of the results

`datalad.tests.utils.assert_is_generator` (*gen*)

`datalad.tests.utils.assert_message` (*message*, *results*)

Verify that each status dict in the results has a message

This only tests the message template string, and not a formatted message with args expanded.

`datalad.tests.utils.assert_no_errors_logged` (*func*, *skip\_re=None*)

Decorator around function to assert that no errors logged during its execution

`datalad.tests.utils.assert_not_in_results` (*results*, *\*\*kwargs*)

Verify that the particular combination of keys and values is not in any of the results

`datalad.tests.utils.assert_re_in` (*regex*, *c*, *flags=0*, *match=True*, *msg=None*)

Assert that container (list, str, etc) contains entry matching the regex

`datalad.tests.utils.assert_repo_status` (*path*, *annex=None*, *untracked\_mode='normal'*, *\*\*kwargs*)

Compare a repo status against (optional) exceptions.

Anything file/directory that is not explicitly indicated must have state 'clean', i.e. no modifications and recorded in Git.

This is an alternative to the traditional `ok_clean_git` helper.

#### Parameters

- **path** (*str* or *Repo*) – in case of a `str`: path to the repository's base dir; Note, that passing a `Repo` instance prevents detecting annex. This might be useful in case of a non-initialized annex, a `GitRepo` is pointing to.
- **annex** (*bool* or *None*) – explicitly set to `True` or `False` to indicate, that an annex is (not) expected; set to `None` to autodetect, whether there is an annex. Default: `None`.
- **untracked\_mode** (`{'no', 'normal', 'all'}`) – If and how untracked content is reported. The specification of untracked files that are OK to be found must match this mode. See `Repo.status()`
- **\*\*kwargs** – Files/directories that are OK to not be in 'clean' state. Each argument must be one of 'added', 'untracked', 'deleted', 'modified' and each value must be a list of filenames (relative to the root of the repository, in POSIX convention).

`datalad.tests.utils.assert_result_count` (*results*, *n*, *\*\*kwargs*)

Verify specific number of results (matching criteria, if any)

`datalad.tests.utils.assert_result_values_cond` (*results*, *prop*, *cond*)

Verify that the values of all results for a given key in the status dicts fulfill condition *cond*.

#### Parameters

- **results** –
- **prop** (*str*) –

- `cond` (*callable*) –

`datalad.tests.utils.assert_result_values_equal` (*results, prop, values*)

Verify that the values of all results for a given key in the status dicts match the given sequence

`datalad.tests.utils.assert_status` (*label, results*)

Verify that each status dict in the results has a given status label

*label* can be a sequence, in which case status must be one of the items in this sequence.

`datalad.tests.utils.assert_str_equal` (*s1, s2*)

Helper to compare two lines

`datalad.tests.utils.check_not_generatorfunction` (*func*)

Internal helper to verify that we are not decorating generator tests

`datalad.tests.utils.clone_url` (*url*)

`datalad.tests.utils.get_convoluted_situation` (*path, repocls=<class 'datalad.support.annexrepo.AnnexRepo'>*)

`datalad.tests.utils.get_datasets_topdir` ()

Delayed parsing so it could be monkey patched etc

`datalad.tests.utils.get_deeply_nested_structure` (*path*)

Here is what this does (assuming UNIX, locked): | . | |— directory\_untracked | | |— link2dir -> ../subdir | |— OBSCURE\_FILENAME\_file\_modified | |— link2dir -> subdir | |— link2subdsdir -> subds\_modified/subdir | |— link2subdsroot -> subds\_modified | |— subdir | | |— annexed\_file.txt -> ../git/annex/objects/... | | |— file\_modified | | |— git\_file.txt | | |— link2annex\_files.txt -> annexed\_file.txt | |— subds\_modified | |— link2superdsdir -> ../subdir | |— subdir | | |— annexed\_file.txt -> ../git/annex/objects/... | |— subds\_lvl1\_modified | |— OBSCURE\_FILENAME\_directory\_untracked | |— untracked\_file

When a system has no symlink support, the link2... components are not included.

`datalad.tests.utils.get_most_obscure_supported_name` (*tdir*)

Return the most obscure filename that the filesystem would support under TEMPDIR

TODO: we might want to use it as a function where we would provide *tdir*

`datalad.tests.utils.get_mtimes_and_digests` (*target\_path*)

Return digests (md5) and mtimes for all the files under *target\_path*

`datalad.tests.utils.has_symlink_capability` ()

`datalad.tests.utils.ignore_nose_capturing_stdout` (*func*)

DEPRECATED and will be removed soon. Does nothing!

Originally was intended as a decorator workaround for nose’s behaviour with redirecting `sys.stdout`, but now we monkey patch nose now so no test should no longer be skipped.

See issue reported here: <https://code.google.com/p/python-nose/issues/detail?id=243&can=1&sort=-id&colspec=ID%20Type%20Status%20Priority%20Stars%20Milestone%20Owner%20Summary>

`datalad.tests.utils.integration` (*f*)

Mark test as an “integration” test which generally is not needed to be run

Generally tend to be slower

`datalad.tests.utils.known_failure` (*func*)

Test decorator marking a test as known to fail

This combines *probe\_known\_failure* and *skip\_known\_failure* giving the skipping precedence over the probing.

`datalad.tests.utils.known_failure_appveyor` (*func*)

Test decorator marking a test as known to fail on AppVeyor.

`datalad.tests.utils.known_failure_direct_mode` (*func*)

DEPRECATED. Stop using. Does nothing

Test decorator marking a test as known to fail in a direct mode test run

If `datalad.repo.direct` is set to `True` behaves like *known\_failure*. Otherwise the original (undecorated) function is returned.

`datalad.tests.utils.known_failure_githubci_win` (*func*)

Test decorator for a known test failure on Github's Windows CI

`datalad.tests.utils.known_failure_v6` (*func*)

Test decorator marking a test as known to fail in a v6+ test run

If the default repository version is 6 or later behaves like *known\_failure*. Otherwise the original (undecorated) function is returned. The default repository version is controlled by the configured value of `DATA-LAD_REPO_VERSION` and whether v5 repositories are supported by the installed git-annex.

`datalad.tests.utils.known_failure_v6_or_later` (*func*)

Test decorator marking a test as known to fail in a v6+ test run

If the default repository version is 6 or later behaves like *known\_failure*. Otherwise the original (undecorated) function is returned. The default repository version is controlled by the configured value of `DATA-LAD_REPO_VERSION` and whether v5 repositories are supported by the installed git-annex.

`datalad.tests.utils.known_failure_windows` (*func*)

Test decorator marking a test as known to fail on windows

On Windows behaves like *known\_failure*. Otherwise the original (undecorated) function is returned.

`datalad.tests.utils.nok_startswith` (*s, prefix*)

`datalad.tests.utils.ok_annex_get` (*ar, files, network=True*)

Helper to run `.get` decorated checking for correct operation

`get` passes through `stderr` from the `ar` to the user, which pollutes screen while running tests

Note: Currently not true anymore, since usage of `-json` disables progressbars

`datalad.tests.utils.ok_archives_caches` (*repopath, n=1, persistent=None*)

Given a path to repository verify number of archives

#### Parameters

- **repopath** (*str*) – Path to the repository
- **n** (*int, optional*) – Number of archives directories to expect
- **persistent** (*bool or None, optional*) – If `None` – both persistent and not count.

`datalad.tests.utils.ok_broken_symlink` (*path*)

`datalad.tests.utils.ok_clean_git` (*path, annex=None, head\_modified=[], index\_modified=[], untracked=[], ignore\_submodules=False*)

Verify that under given path there is a clean git repository

it exists, `.git` exists, nothing is uncommitted/dirty/staged

---

**Note:** Parameters `head_modified` and `index_modified` currently work in pure git or indirect mode annex only. If they are given, no test of modification of known repo content is performed.

---

#### Parameters

- **path** (*str or Repo*) – in case of a str: path to the repository’s base dir; Note, that passing a Repo instance prevents detecting annex. This might be useful in case of a non-initialized annex, a GitRepo is pointing to.
- **annex** (*bool or None*) – explicitly set to True or False to indicate, that an annex is (not) expected; set to None to autodetect, whether there is an annex. Default: None.
- **ignore\_submodules** (*bool*) – if True, submodules are not inspected

`datalad.tests.utils.ok_endswith` (*s, suffix*)

`datalad.tests.utils.ok_exists` (*path*)

`datalad.tests.utils.ok_file_has_content` (*path, content, strip=False, re=False, decompress=False, \*\*kwargs*)

Verify that file exists and has expected content

`datalad.tests.utils.ok_file_under_git` (*path, filename=None, annexed=False*)

Test if file is present and under git/annex control

If relative path provided, then test from current directory

`datalad.tests.utils.ok_generator` (*gen*)

`datalad.tests.utils.ok_git_config_not_empty` (*ar*)

Helper to verify that nothing rewritten the config file

`datalad.tests.utils.ok_good_symlink` (*path*)

`datalad.tests.utils.ok_startswith` (*s, prefix*)

`datalad.tests.utils.ok_symlink` (*path*)

Checks whether path is either a working or broken symlink

`datalad.tests.utils.patch_config` (*vars*)

Patch our config with custom settings. Returns mock.patch cm

`datalad.tests.utils.probe_known_failure` (*func*)

Test decorator allowing the test to pass when it fails and vice versa

Setting config `datalad.tests.knownfailures.probe` to True tests, whether or not the test is still failing. If it’s not, an `AssertionError` is raised in order to indicate that the reason for failure seems to be gone.

`datalad.tests.utils.put_file_under_git` (*path, filename=None, content=None, annexed=False*)

Place file under git/annex and return used Repo

`datalad.tests.utils.set_annex_version` (*version*)

Override the git-annex version.

This temporarily masks the git-annex version present in `external_versions` and make `AnnexRepo` forget its cached version information.

`datalad.tests.utils.set_date` (*timestamp*)

Temporarily override environment variables for git/git-annex dates.

**Parameters** `timestamp` (*int*) – Unix timestamp.

`datalad.tests.utils.skip_httpretty_on_problematic_pythons` (*func*)

As discovered some `httpretty` bug causes a side-effect on other tests on some Pythons. So we skip the test if such problematic combination detected

References <https://travis-ci.org/datalad/datalad/jobs/94464988> <http://stackoverflow.com/a/29603206/1265472>

`datalad.tests.utils.skip_if_no_module` (*module*)

`datalad.tests.utils.skip_if_no_network` (*func=None*)

Skip test completely in NONETWORK settings

If not used as a decorator, and just a function, could be used at the module level

`datalad.tests.utils.skip_if_on_windows` (*func=None*)

Skip test completely under Windows

`datalad.tests.utils.skip_if_scrapy_without_selector` ()

A little helper to skip some tests which require recent scrapy

`datalad.tests.utils.skip_if_url_is_not_available` (*url, regex=None*)

`datalad.tests.utils.skip_ssh` (*func*)

Skips SSH tests if on windows or if environment variable DATALAD\_TESTS\_SSH was not set

`datalad.tests.utils.skip_wo_symlink_capability` (*func*)

Skip test when environment does not support symlinks

Perform a behavioral test instead of top-down logic, as on windows this could be on or off on a case-by-case basis.

`datalad.tests.utils.slow` (*f*)

Mark test as a slow, although not necessarily integration or usecase test

`datalad.tests.utils.usecase` (*f*)

Mark test as a usecase user ran into and which (typically) caused bug report to be filed/troubleshooted

## `datalad.tests.utils_testrepos`

**class** `datalad.tests.utils_testrepos.BasicAnnexTestRepo` (*path=None, puke\_if\_exists=True*)

Bases: `datalad.tests.utils_testrepos.TestRepo`

Creates a basic test git-annex repository

**REPO\_CLASS**

alias of `datalad.support.annexrepo.AnnexRepo`

**create\_info\_file** ()

**populate** ()

**class** `datalad.tests.utils_testrepos.BasicGitTestRepo` (*path=None, puke\_if\_exists=True*)

Bases: `datalad.tests.utils_testrepos.TestRepo`

Creates a basic test git repository.

**REPO\_CLASS**

alias of `datalad.support.gitrepo.GitRepo`

**create\_info\_file** ()

**populate** ()

**class** `datalad.tests.utils_testrepos.InnerSubmodule`

Bases: object

**create** ()

**path**

**url**

```
class datalad.tests.utils_testrepos.NestedDataset (path=None, puke_if_exists=True)
    Bases: datalad.tests.utils_testrepos.BasicAnnexTestRepo

    populate()

class datalad.tests.utils_testrepos.SubmoduleDataset (path=None,
                                                    puke_if_exists=True)
    Bases: datalad.tests.utils_testrepos.BasicAnnexTestRepo

    populate()

class datalad.tests.utils_testrepos.TestRepo (path=None, puke_if_exists=True)
    Bases: object

    REPO_CLASS = None

    create()

    create_file (name, content, add=True, annex=False)

    path

    populate()

    url
```

## datalad.tests.heavyoutput

Helper to provide heavy load on stdout and stderr

## Command line interface infrastructure

---

*cmdline.main*

---

*cmdline.helpers*

---

*cmdline.common\_args*

---

## datalad.cmdline.main

```
class datalad.cmdline.main.ArgumentParserDisableAbbrev (prog=None, usage=None,
                                                    description=None, epi-
                                                    log=None, parents=[],
                                                    formatter_class=<class
                                                    'argparse.HelpFormatter'>,
                                                    prefix_chars='-', from-
                                                    file_prefix_chars=None,
                                                    argument_default=None,
                                                    conflict_handler='error',
                                                    add_help=True, al-
                                                    low_abbrev=True)

    Bases: argparse.ArgumentParser

datalad.cmdline.main.add_entrpoints_to_interface_groups (interface_groups)

datalad.cmdline.main.fail_with_short_help (parser=None, msg=None, known=None,
                                                    provided=None, hint=None, exit_code=1,
                                                    what='command', out=None)

    Generic helper to fail with short help possibly hinting on what was intended if known were provided
```



`datalad.cmdline.main.get_commands_from_groups` (*groups*)

Get a dictionary of command: interface\_spec

`datalad.cmdline.main.get_description_with_cmd_summary` (*grp\_short\_descriptions*,  
*interface\_groups*,  
*parser\_description*)

`datalad.cmdline.main.main` (*args=None*)

`datalad.cmdline.main.setup_parser` (*cmdlineargs*, *formatter\_class=<class 'arg-  
parse.RawDescriptionHelpFormatter'>*, *re-  
turn\_subparsers=False*, *help\_ignore\_extensions=False*)

## datalad.cmdline.helpers

**class** `datalad.cmdline.helpers.HelpAction` (*option\_strings*, *dest*, *nargs=None*, *const=None*,  
*default=None*, *type=None*, *choices=None*, *re-  
quired=False*, *help=None*, *metavar=None*)

Bases: `argparse.Action`

**class** `datalad.cmdline.helpers.LogLevelAction` (*option\_strings*, *dest*, *nargs=None*,  
*const=None*, *default=None*, *type=None*,  
*choices=None*, *required=False*,  
*help=None*, *metavar=None*)

Bases: `argparse.Action`

`datalad.cmdline.helpers.get_repo_instance` (*path='.'*, *class\_=None*)

Returns an instance of appropriate datalad repository for path. Check whether a certain path is inside a known type of repository and returns an instance representing it. May also check for a certain type instead of detecting the type of repository.

### Parameters

- **path** (*str*) – path to check; default: current working directory
- **class** (*class*) – if given, check whether path is inside a repository, that can be represented as an instance of the passed class.

**Raises** `RuntimeError`, in case cwd is not inside a known repository.

`datalad.cmdline.helpers.parser_add_common_opt` (*parser*, *opt*, *names=None*, *\*\*kwargs*)

`datalad.cmdline.helpers.run_via_pbs` (*args*, *pbs*)

`datalad.cmdline.helpers.strip_arg_from_argv` (*args*, *value*, *opt\_names*)

Strip an originally listed option (with its value) from the list cmdline args

## datalad.cmdline.common\_args

### 1.5.3 Configuration

DataLad uses the same configuration mechanism and syntax as Git itself. Consequently, datalad can be configured using the `git config` command. Both a *global* user configuration (typically at `~/.gitconfig`), and a *local* repository-specific configuration (`.git/config`) are inspected.

In addition, datalad supports a persistent dataset-specific configuration. This configuration is stored at `.datalad/config` in any dataset. As it is part of a dataset, settings stored there will also be in effect for any consumer of such a dataset. Both *global* and *local* settings on a particular machine always override configuration shipped with a dataset.

All datalad-specific configuration variables are prefixed with `datalad.`.

It is possible to override or amend the configuration using environment variables. Any variable with a name that starts with `DATALAD_` will be available as the corresponding `datalad.` configuration variable, replacing any `__` (two underscores) with a hyphen, then any `_` (single underscore) with a dot, and finally converting all letters to lower case. Values from environment variables take precedence over configuration file settings.

The following sections provide a (non-exhaustive) list of settings honored by datalad. They are categorized according to the scope they are typically associated with.

### Global user configuration

**datalad.externals.nda.dbserver** NDA database server: Hostname of the database server

**datalad.locations.cache** Cache directory: Where should datalad cache files? Default: `~/cache/datalad`

**datalad.locations.default-dataset** Default dataset path: Where should datalad should look for (or install) a default dataset? Default: `~/datalad`

**datalad.locations.system-plugins** System plugin directory: Where should datalad search for system plugins? Default: `/etc/xdg/datalad/plugins`

**datalad.locations.system-procedures** System procedure directory: Where should datalad search for system procedures? Default: `/etc/xdg/datalad/procedures`

**datalad.locations.user-plugins** User plugin directory: Where should datalad search for user plugins? Default: `~/config/datalad/plugins`

**datalad.locations.user-procedures** User procedure directory: Where should datalad search for user procedures? Default: `~/config/datalad/procedures`

**datalad.ssh.identityfile** If set, pass this file as ssh's `-i` option.: Default: None

### Local repository configuration

**datalad.crawl.cache** Crawler download caching: Should the crawler cache downloaded files?

[bool]

**datalad.fake-dates** Fake (anonymize) dates: Should the dates in the logs be faked? Default: False

[value must be convertible to type bool]

### Sticky dataset configuration

**datalad.locations.dataset-procedures** Dataset procedure directory: Where should datalad search for dataset procedures (relative to a dataset root)? Default: `.datalad/procedures`

### Miscellaneous configuration

**datalad.cmd.protocol** Specifies the protocol number used by the Runner to note shell command or python function call times and allows for dry runs. “externals-time” for ExecutionTimeExternalsProtocol, “time” for ExecutionTimeProtocol and “null” for NullProtocol. Any new `DATALAD_CMD_PROTOCOL` has to implement `datalad.support.protocol.ProtocolInterface`:

**datalad.cmd.protocol.prefix** Sets a prefix to add before the command call times are noted by `DATALAD_CMD_PROTOCOL`..:

**datalad.exc.str.tblimit** This flag is used by the `datalad extract_tb` function which extracts and formats stack-traces. It caps the number of lines to `DATALAD_EXC_STR_TBLIMIT` of pre-processed entries from traceback.:

**datalad.fake-dates-start** Initial fake date: When faking dates and there are no commits in any local branches, generate the date by adding one second to this value (Unix epoch time). The value must be positive. Default: 1112911993

[value must be convertible to type 'int']

**datalad.github.token-note** Github token note: Description for a Personal access token to generate. Default: DataLad

**datalad.install.inherit-local-origin** Inherit local origin of dataset source: If enabled, a local 'origin' remote of a local dataset clone source is configured as an 'origin-2' remote to make its annex automatically available. The process is repeated recursively for any further qualifying 'origin' dataset thereof. Default: True

[value must be convertible to type bool]

**datalad.log.level** Used for control the verbosity of logs printed to stdout while running datalad commands/debugging:

**datalad.log.name** Include name of the log target in the log line:

**datalad.log.names** Which names (,-separated) to print log lines for:

**datalad.log.namesre** Regular expression for which names to print log lines for:

**datalad.log.outputs** Used to control whether both stdout and stderr of external commands execution are logged in detail (at DEBUG level):

**datalad.log.result-level** Log level for command result messages: Overrides the default behavior of logging 'impossible' results as a warning, 'error' results as errors, and everything else as 'debug' with a single alternative log level Default: None

[value must be one of ('debug', 'info', 'warning', 'error')]

**datalad.log.timestamp** Used to add timestamp to datalad logs: Default: False

[value must be convertible to type bool]

**datalad.log.traceback** Runs TraceBack function with `collide` set to True, if this flag is set to "collide". This replaces any common prefix between current traceback log and previous invocation with "...":

**datalad.metadata.create-aggregate-annex-limit** Limit configuration annexing aggregated metadata in new dataset: Git-annex large files expression (see <https://git-annex.branchable.com/tips/largefiles>; given expression will be wrapped in parentheses) Default: anything

**datalad.metadata.maxfieldsize** Maximum metadata field size: Metadata fields exceeding this size (in bytes/chars) are excluded from metadata extractio Default: 100000

[value must be convertible to type 'int']

**datalad.metadata.nativetype** Native dataset metadata scheme: Set this label to engage a particular metadata extraction parser

**datalad.metadata.store-aggregate-content** Aggregated content metadata storage: If this flag is enabled, content metadata is aggregated into superdataset to allow for discovery of individual files. If disable unique content metadata values are still aggregated to enable dataset discovery Default: True

[value must be convertible to type bool]

**datalad.repo.backend** git-annex backend: Backend to use when creating git-annex repositories Default: MD5E

**datalad.repo.direct** Direct Mode for git-annex repositories: Set this flag to create annex repositories in direct mode by default Default: False

[value must be convertible to type bool]

**datalad.repo.version** git-annex repository version: Specifies the repository version for git-annex to be used by default  
Default: 5

[value must be convertible to type 'int']

**datalad.runtime.raiseonerror** Error behavior: Set this flag to cause DataLad to raise an exception on errors that would have otherwise just get logged  
Default: False

[value must be convertible to type bool]

**datalad.runtime.report-status** Command line result reporting behavior: If set (to other than 'all'), constrains command result report to records matching the given status. 'success' is a synonym for 'ok' OR 'notneeded', 'failure' stands for 'impossible' OR 'error'  
Default: None

[value must be one of ('all', 'success', 'failure', 'ok', 'notneeded', 'impossible', 'error')]

**datalad.search.default-mode** Default search mode: Label of the mode to be used by default  
Default: egrep

[value must be one of ('egrep', 'textblob', 'autofield')]

**datalad.search.index-default-documenttype** Type of search index documents: Labels of document types to include in a default search index  
Default: datasets

[value must be one of ('all', 'datasets', 'files')]

**datalad.search.indexercachesize** Maximum cache size for search index (per process): Actual memory consumption can be twice as high as this value in MB (one process per CPU is used)  
Default: 256

[value must be convertible to type 'int']

**datalad.tests.dataladremote** Binary flag to specify whether each annex repository should get datalad special remote in every test repository:

[value must be convertible to type bool]

**datalad.tests.knownfailures.probe** Probes tests that are known to fail on whether or not they are actually still failing:  
Default: False

[value must be convertible to type bool]

**datalad.tests.knownfailures.skip** Skips tests that are known to currently fail: **Default: True**

[value must be convertible to type bool]

**datalad.tests.nonetwork** Skips network tests completely if this flag is set Examples include test for s3, git\_repositories, openfmri etc:

[value must be convertible to type bool]

**datalad.tests.nonlo** Specifies network interfaces to bring down/up for testing. Currently used by travis.:

**datalad.tests.noteardown** Does not execute teardown\_package which cleans up temp files and directories created by tests if this flag is set:

[value must be convertible to type bool]

**datalad.tests.protocolremote** Binary flag to specify whether to test protocol interactions of custom remote with annex:

[value must be convertible to type bool]

**datalad.tests.runcmdline** Binary flag to specify if shell testing using shunit2 to be carried out:

[value must be convertible to type bool]

**datalad.tests.ssh** Skips SSH tests if this flag is **not** set:

[value must be convertible to type bool]

**datalad.tests.temp.dir** Create a temporary directory at location specified by this flag. It is used by tests to create a temporary git directory while testing git annex archives etc:

**datalad.tests.temp.fs** Specify the temporary file system to use as loop device for testing DATA-LAD\_TESTS\_TEMP\_DIR creation:

**datalad.tests.temp.fssize** Specify the size of temporary file system to use as loop device for testing DATA-LAD\_TESTS\_TEMP\_DIR creation:

**datalad.tests.temp.keep** Function rmtmp will not remove temporary file/directory created for testing if this flag is set:

[value must be convertible to type bool]

**datalad.tests.ui.backend** Tests UI backend: Which UI backend to use Default: tests-noninteractive

**datalad.tests.usecassette** Specifies the location of the file to record network transactions by the VCR module. Currently used by when testing custom special remotes:

**datalad.ui.color** Colored terminal output: Enable or disable ANSI color codes in outputs; “on” overrides NO\_COLOR environment variable Default: auto

[value must be one of ('on', 'off', 'auto')]

**datalad.ui.progressbar** UI progress bars: Default backend for progress reporting Default: None

[value must be one of ('tqdm', 'tqdm-ipython', 'log', 'none')]

## 1.6 Extension packages

DataLad can be customized and additional functionality can be integrated via extensions. Each extension provides its own documentation:

- [Crawling web resources and automated data distributions](#)
- [Neuroimaging data and workflows](#)
- [Containerized computational environments](#)
- [Advanced metadata tooling with JSON-LD reporting and additional metadata extractors](#)

## 1.7 Indices and tables

- [genindex](#)
- [modindex](#)
- [search](#)



## d

- [datalad.auto](#), 210
- [datalad.cmd](#), 211
- [datalad.cmdline.common\\_args](#), 277
- [datalad.cmdline.helpers](#), 277
- [datalad.cmdline.main](#), 276
- [datalad.config](#), 266
- [datalad.consts](#), 216
- [datalad.customremotes.archives](#), 264
- [datalad.customremotes.base](#), 261
- [datalad.customremotes.main](#), 260
- [datalad.log](#), 216
- [datalad.plugin.add\\_readme](#), 54
- [datalad.plugin.addurls](#), 55
- [datalad.plugin.check\\_dates](#), 60
- [datalad.plugin.export\\_archive](#), 61
- [datalad.plugin.export\\_to\\_figshare](#), 62
- [datalad.plugin.no\\_annex](#), 64
- [datalad.plugin.wtf](#), 66
- [datalad.support.annexrepo](#), 245
- [datalad.support.archives](#), 258
- [datalad.support.configparserinc](#), 260
- [datalad.support.gitrepo](#), 228
- [datalad.tests.heavyoutput](#), 276
- [datalad.tests.utils](#), 270
- [datalad.tests.utils\\_testrepos](#), 275
- [datalad.utils](#), 216
- [datalad.version](#), 228





## Symbols

`__init__()` (*datalad.api.Dataset* method), 154

## A

`ac` (*datalad.plugin.check\_dates.CheckDates* attribute), 60

`activate()` (*datalad.auto.AutomagicIO* method), 210

`active` (*datalad.auto.AutomagicIO* attribute), 210

`add()` (*datalad.config.ConfigManager* method), 266

`add()` (*datalad.support.annexrepo.AnnexRepo* method), 245

`add()` (*datalad.support.gitrepo.GitRepo* method), 229

`add_()` (*datalad.support.annexrepo.AnnexRepo* method), 245

`add_()` (*datalad.support.gitrepo.GitRepo* method), 230

`add_archive_content()` (in module *datalad.api*), 208

`add_entrypoints_to_interface_groups()` (in module *datalad.cmdline.main*), 276

`add_extra_filename_values()` (in module *datalad.plugin.addurls*), 58

`add_fake_dates()` (*datalad.support.gitrepo.GitRepo* method), 230

`add_remote()` (*datalad.support.annexrepo.AnnexRepo* method), 245

`add_remote()` (*datalad.support.gitrepo.GitRepo* method), 230

`add_section()` (*datalad.customremotes.base.AnnexExchangeProtocol* method), 263

`add_submodule()` (*datalad.support.gitrepo.GitRepo* method), 230

`add_url_to_file()` (*datalad.support.annexrepo.AnnexRepo* method), 245

`add_urls()` (*datalad.support.annexrepo.AnnexRepo* method), 246

`AddReadme` (class in *datalad.plugin.add\_readme*), 54

`AddReadme.EnsureChoice` (class in *datalad.plugin.add\_readme*), 54

`AddReadme.EnsureDataset` (class in *datalad.plugin.add\_readme*), 54

`AddReadme.EnsureNone` (class in *datalad.plugin.add\_readme*), 54

`AddReadme.EnsureStr` (class in *datalad.plugin.add\_readme*), 54

`AddReadme.Parameter` (class in *datalad.plugin.add\_readme*), 54

`Addurls` (class in *datalad.plugin.addurls*), 55

`Addurls.EnsureChoice` (class in *datalad.plugin.addurls*), 56

`Addurls.EnsureDataset` (class in *datalad.plugin.addurls*), 56

`Addurls.EnsureNone` (class in *datalad.plugin.addurls*), 56

`Addurls.EnsureStr` (class in *datalad.plugin.addurls*), 57

`Addurls.Parameter` (class in *datalad.plugin.addurls*), 57

`adjust()` (*datalad.support.annexrepo.AnnexRepo* method), 246

`aggregate_metadata()` (in module *datalad.api*), 183

`all_same()` (in module *datalad.utils*), 217

`annex`, 71

`AnnexCustomRemote` (class in *datalad.customremotes.base*), 261

`AnnexExchangeProtocol` (class in *datalad.customremotes.base*), 263

`AnnexRemoteQuit`, 264

`AnnexRepo` (class in *datalad.support.annexrepo*), 245

`annexstatus()` (*datalad.support.annexrepo.AnnexRepo* method), 246

`annotate_paths()` (in module *datalad.api*), 194

`any_re_search()` (in module *datalad.utils*), 217

`anything2bool()` (in module *datalad.config*), 270

`API_URL` (*datalad.plugin.export\_to\_figshare.FigshareRESTLaison*

*attribute*), 64  
 ArchiveAnnexCustomRemote (*class in data-lad.customremotes.archives*), 264  
 ArchivesCache (*class in datalad.support.archives*), 258  
 args (*datalad.utils.ArgSpecFake attribute*), 216  
 ArgSpecFake (*class in datalad.utils*), 216  
 ArgumentParserDisableAbbrev (*class in data-lad.cmdline.main*), 276  
 as\_unicode() (*in module datalad.utils*), 217  
 assert\_dict\_equal() (*in module data-lad.tests.utils*), 271  
 assert\_in\_results() (*in module data-lad.tests.utils*), 271  
 assert\_is\_generator() (*in module data-lad.tests.utils*), 271  
 assert\_message() (*in module datalad.tests.utils*), 271  
 assert\_no\_errors\_logged() (*in module data-lad.tests.utils*), 271  
 assert\_no\_open\_files() (*in module data-lad.utils*), 217  
 assert\_not\_in\_results() (*in module data-lad.tests.utils*), 271  
 assert\_re\_in() (*in module datalad.tests.utils*), 271  
 assert\_repo\_status() (*in module data-lad.tests.utils*), 271  
 assert\_result\_count() (*in module data-lad.tests.utils*), 271  
 assert\_result\_values\_cond() (*in module data-lad.tests.utils*), 271  
 assert\_result\_values\_equal() (*in module datalad.tests.utils*), 272  
 assert\_status() (*in module datalad.tests.utils*), 272  
 assert\_str\_equal() (*in module data-lad.tests.utils*), 272  
 assure\_bool() (*in module datalad.utils*), 217  
 assure\_bytes() (*in module datalad.utils*), 217  
 assure\_dict\_from\_str() (*in module data-lad.utils*), 217  
 assure\_dir() (*in module datalad.utils*), 217  
 assure\_extracted() (*datalad.support.archives.ExtractedArchive method*), 259  
 assure\_iter() (*in module datalad.utils*), 218  
 assure\_list() (*in module datalad.utils*), 218  
 assure\_list\_from\_str() (*in module data-lad.utils*), 218  
 assure\_tuple\_or\_list() (*in module data-lad.utils*), 218  
 assure\_unicode() (*in module datalad.utils*), 218  
 auto\_repr() (*in module datalad.utils*), 218  
 autoget (*datalad.auto.AutomagicIO attribute*), 210  
 AutomagicIO (*class in datalad.auto*), 210

AVAILABILITY (*datalad.customremotes.archives.ArchiveAnnexCustomRemote attribute*), 264

AVAILABILITY (*datalad.customremotes.base.AnnexCustomRemote attribute*), 261

## B

BasicAnnexTestRepo (*class in data-lad.tests.utils\_testrepos*), 275

BasicGitTestRepo (*class in data-lad.tests.utils\_testrepos*), 275

BatchedAnnex (*class in datalad.support.annexrepo*), 258

BatchedAnnexes (*class in data-lad.support.annexrepo*), 258

BatchedCommand (*class in datalad.cmd*), 211

BEGIN (*datalad.support.gitrepo.GitProgress attribute*), 229

better\_wraps() (*in module datalad.utils*), 218

bytes2human() (*in module datalad.utils*), 218

## C

cache (*datalad.customremotes.archives.ArchiveAnnexCustomRemote attribute*), 264

call() (*datalad.cmd.Runner method*), 212

call\_git() (*datalad.support.gitrepo.GitRepo method*), 230

call\_git\_items() (*datalad.support.gitrepo.GitRepo method*), 231

call\_git\_online() (*datalad.support.gitrepo.GitRepo method*), 231

call\_git\_success() (*datalad.support.gitrepo.GitRepo method*), 231

check\_direct\_mode\_support() (*datalad.support.annexrepo.AnnexRepo class method*), 246

check\_not\_generatorfunction() (*in module datalad.tests.utils*), 272

check\_repository\_versions() (*datalad.support.annexrepo.AnnexRepo class method*), 246

check\_symlink\_capability() (*in module datalad.utils*), 219

CheckDates (*class in datalad.plugin.check\_dates*), 60

CheckDates.EnsureChoice (*class in datalad.plugin.check\_dates*), 60

CheckDates.EnsureNone (*class in datalad.plugin.check\_dates*), 60

CheckDates.EnsureStr (*class in datalad.plugin.check\_dates*), 60

CheckDates.Parameter (*class in datalad.plugin.check\_dates*), 60

CHECKING\_OUT (*datalad.support.gitrepo.GitProgress attribute*), 229

checkout() (*datalad.support.gitrepo.GitRepo method*), 232

cherry\_pick() (*datalad.support.gitrepo.GitRepo method*), 232

chpwd (*class in datalad.utils*), 219

clean() (*datalad.support.archives.ArchivesCache method*), 259

clean() (*datalad.support.archives.ExtractedArchive method*), 259

clean() (*in module datalad.api*), 196

clean\_meta\_args() (*in module datalad.plugin.addurls*), 58

clear() (*datalad.support.annexrepo.BatchedAnnexes method*), 258

clone() (*datalad.support.gitrepo.GitRepo class method*), 232

clone() (*in module datalad.api*), 197

clone\_url() (*in module datalad.tests.utils*), 272

close() (*datalad.cmd.BatchedCommand method*), 211

close() (*datalad.support.annexrepo.BatchedAnnexes method*), 258

ColorFormatter (*class in datalad.log*), 216

commands (*datalad.cmd.Runner attribute*), 213

commit() (*datalad.support.gitrepo.GitRepo method*), 232

commit\_exists() (*datalad.support.gitrepo.GitRepo method*), 232

COMPRESSING (*datalad.support.gitrepo.GitProgress attribute*), 229

config (*datalad.support.gitrepo.GitRepo attribute*), 233

ConfigManager (*class in datalad.config*), 266

configure\_fake\_dates() (*datalad.support.gitrepo.GitRepo method*), 233

connection\_made() (*datalad.cmd.WitlessProtocol method*), 214

connection\_made() (*datalad.support.gitrepo.GitProgress method*), 229

convert\_field() (*datalad.plugin.addurls.Formatter method*), 58

copy\_to() (*datalad.support.annexrepo.AnnexRepo method*), 246

COST (*datalad.customremotes.archives.ArchiveAnnexCustomRemote attribute*), 264

COST (*datalad.customremotes.base.AnnexCustomRemote attribute*), 261

count\_objects (*datalad.support.gitrepo.GitRepo attribute*), 233

COUNTING (*datalad.support.gitrepo.GitProgress attribute*), 229

create() (*datalad.tests.utils\_testrepos.InnerSubmodule method*), 275

create() (*datalad.tests.utils\_testrepos.TestRepo method*), 276

create() (*in module datalad.api*), 155

create\_article() (*datalad.plugin.export\_to\_figshare.FigshareRESTLiaison method*), 64

create\_file() (*datalad.tests.utils\_testrepos.TestRepo method*), 276

create\_info\_file() (*datalad.tests.utils\_testrepos.BasicAnnexTestRepo method*), 275

create\_info\_file() (*datalad.tests.utils\_testrepos.BasicGitTestRepo method*), 275

create\_sibling() (*in module datalad.api*), 158

create\_sibling\_github() (*in module datalad.api*), 160

create\_sibling\_gitlab() (*in module datalad.api*), 161

create\_test\_dataset() (*in module datalad.api*), 199

create\_tree() (*in module datalad.utils*), 219

create\_tree\_archive() (*in module datalad.utils*), 219

CUSTOM\_REMOTE\_NAME (*datalad.customremotes.archives.ArchiveAnnexCustomRemote attribute*), 264

CUSTOM\_REMOTE\_NAME (*datalad.customremotes.base.AnnexCustomRemote attribute*), 261

custom\_result\_renderer() (*datalad.plugin.check\_dates.CheckDates static method*), 60

custom\_result\_renderer() (*datalad.plugin.wtf.WTF static method*), 67

cwd (*datalad.cmd.Runner attribute*), 213

cwd (*datalad.cmd.WitlessRunner attribute*), 215

## D

DataLad extension, 71

datalad.auto (*module*), 210

datalad.cmd (*module*), 211

datalad.cmd.protocol, 278

datalad.cmd.protocol.prefix, 278

datalad.cmdline.common\_args (*module*), 277

datalad.cmdline.helpers (*module*), 277

datalad.cmdline.main (*module*), 276

datalad.config (*module*), 266

datalad.consts (*module*), 216

datalad.crawl.cache, 278

datalad.customremotes.archives (*module*), 264

datalad.customremotes.base (*module*), 261  
 datalad.customremotes.main (*module*), 260  
 datalad.exc.str.tblimit, 279  
 datalad.externals.nda.dbserver, 278  
 datalad.fake-dates, 278  
 datalad.fake-dates-start, 279  
 datalad.github.token-note, 279  
 datalad.install.inherit-local-origin, 279  
 datalad.locations.cache, 278  
 datalad.locations.dataset-procedures, 278  
 datalad.locations.default-dataset, 278  
 datalad.locations.system-plugins, 278  
 datalad.locations.system-procedures, 278  
 datalad.locations.user-plugins, 278  
 datalad.locations.user-procedures, 278  
 datalad.log (*module*), 216  
 datalad.log.level, 279  
 datalad.log.name, 279  
 datalad.log.names, 279  
 datalad.log.namesre, 279  
 datalad.log.outputs, 279  
 datalad.log.result-level, 279  
 datalad.log.timestamp, 279  
 datalad.log.traceback, 279  
 datalad.metadata.create-aggregate-annex-dataset, 279  
 datalad.metadata.maxfieldsize, 279  
 datalad.metadata.nativetype, 279  
 datalad.metadata.store-aggregate-content, 279  
 datalad.plugin.add\_readme (*module*), 54  
 datalad.plugin.addurls (*module*), 55  
 datalad.plugin.check\_dates (*module*), 60  
 datalad.plugin.export\_archive (*module*), 61  
 datalad.plugin.export\_to\_figshare (*module*), 62  
 datalad.plugin.no\_annex (*module*), 64  
 datalad.plugin.wtf (*module*), 66  
 datalad.repo.backend, 279  
 datalad.repo.direct, 279  
 datalad.repo.version, 280  
 datalad.runtime.raiseonerror, 280  
 datalad.runtime.report-status, 280  
 datalad.search.default-mode, 280  
 datalad.search.index-default-documenttype, 280  
 datalad.search.indexercachesize, 280  
 datalad.ssh.identityfile, 278  
 datalad.support.annexrepo (*module*), 245  
 datalad.support.archives (*module*), 258  
 datalad.support.configparserinc (*module*), 260  
 datalad.support.gitrepo (*module*), 228  
 datalad.tests.dataladremote, 280  
 datalad.tests.heavyoutput (*module*), 276  
 datalad.tests.knownfailures.probe, 280  
 datalad.tests.knownfailures.skip, 280  
 datalad.tests.nonetwork, 280  
 datalad.tests.nonlo, 280  
 datalad.tests.noteardown, 280  
 datalad.tests.protocolremote, 280  
 datalad.tests.runcmdline, 280  
 datalad.tests.ssh, 280  
 datalad.tests.temp.dir, 281  
 datalad.tests.temp.fs, 281  
 datalad.tests.temp.fssize, 281  
 datalad.tests.temp.keep, 281  
 datalad.tests.ui.backend, 281  
 datalad.tests.usecassette, 281  
 datalad.tests.utils (*module*), 270  
 datalad.tests.utils\_testrepos (*module*), 275  
 datalad.ui.color, 281  
 datalad.ui.progressbar, 281  
 datalad.utils (*module*), 216  
 datalad.version (*module*), 228  
 dataset, 71  
 Dataset (*class in datalad.api*), 154  
 datasetmethod() (*data-lad.plugin.add\_readme.AddReadme method*), 54  
 datasetmethod() (*data-lad.plugin.addurls.Addurls method*), 57  
 datasetmethod() (*data-lad.plugin.export\_archive.ExportArchive method*), 62  
 datasetmethod() (*data-lad.plugin.export\_to\_figshare.ExportToFigshare method*), 63  
 datasetmethod() (*data-lad.plugin.no\_annex.NoAnnex method*), 65  
 datasetmethod() (*data-lad.plugin.wtf.WTF method*), 67  
 deactivate() (*data-lad.auto.AutomagicIO method*), 210  
 debug() (*data-lad.customremotes.base.AnnexCustomRemote method*), 261  
 decode\_input() (*in module datalad.utils*), 219  
 decompress\_file() (*in module data-lad.support.archives*), 259  
 default\_backends (*data-lad.support.annexrepo.AnnexRepo attribute*), 247  
 defaults (*data-lad.utils.ArgSpecFake attribute*), 216

deinit\_submodule() (*datalad.support.gitrepo.GitRepo* method), 233  
 DELETED (*datalad.support.gitrepo.PushInfo* attribute), 243  
 describe() (*datalad.support.gitrepo.GitRepo* method), 233  
 diff() (*datalad.support.gitrepo.GitRepo* method), 233  
 diff() (*in module datalad.api*), 199  
 diffstatus() (*datalad.support.gitrepo.GitRepo* method), 234  
 dirty (*datalad.support.gitrepo.GitRepo* attribute), 234  
 disable\_logger() (*in module datalad.utils*), 220  
 dlabspath() (*in module datalad.utils*), 220  
 do\_execute\_callables (*datalad.customremotes.base.AnnexExchangeProtocol* attribute), 263  
 do\_execute\_ext\_commands (*datalad.customremotes.base.AnnexExchangeProtocol* attribute), 263  
 DONE\_TOKEN (*datalad.support.gitrepo.GitProgress* attribute), 229  
 download\_url() (*in module datalad.api*), 201  
 drop() (*datalad.support.annexrepo.AnnexRepo* method), 247  
 drop() (*in module datalad.api*), 163  
 drop\_key() (*datalad.support.annexrepo.AnnexRepo* method), 247  
 dry (*datalad.cmd.Runner* attribute), 213

## E

enable\_remote() (*datalad.support.annexrepo.AnnexRepo* method), 247  
 encode\_filename() (*in module datalad.utils*), 220  
 END (*datalad.support.gitrepo.GitProgress* attribute), 229  
 end\_section() (*datalad.customremotes.base.AnnexExchangeProtocol* method), 263  
 ensure\_bool() (*in module datalad.utils*), 220  
 ensure\_bytes() (*in module datalad.utils*), 220  
 ensure\_dict\_from\_str() (*in module datalad.utils*), 220  
 ensure\_dir() (*in module datalad.utils*), 220  
 ensure\_iter() (*in module datalad.utils*), 220  
 ensure\_list() (*in module datalad.utils*), 220  
 ensure\_list\_from\_str() (*in module datalad.utils*), 221  
 ensure\_tuple\_or\_list() (*in module datalad.utils*), 221  
 ensure\_unicode() (*in module datalad.utils*), 221  
 ENUMERATING (*datalad.support.gitrepo.GitProgress* attribute), 229  
 env (*datalad.cmd.Runner* attribute), 213  
 env (*datalad.cmd.WitlessRunner* attribute), 215  
 ERROR (*datalad.support.gitrepo.FetchInfo* attribute), 228  
 ERROR (*datalad.support.gitrepo.PushInfo* attribute), 243  
 error() (*datalad.customremotes.base.AnnexCustomRemote* method), 261  
 escape\_filename() (*in module datalad.utils*), 221  
 eval\_results() (*datalad.plugin.add\_readme.AddReadme* method), 55  
 eval\_results() (*datalad.plugin.addurls.Addurls* method), 57  
 eval\_results() (*datalad.plugin.check\_dates.CheckDates* method), 60  
 eval\_results() (*datalad.plugin.export\_archive.ExportArchive* method), 62  
 eval\_results() (*datalad.plugin.export\_to\_figshare.ExportToFigshare* method), 64  
 eval\_results() (*datalad.plugin.no\_annex.NoAnnex* method), 65  
 eval\_results() (*datalad.plugin.wtf.WTF* method), 67  
 expandpath() (*in module datalad.utils*), 221  
 ExportArchive (class *in datalad.plugin.export\_archive*), 61  
 ExportArchive.EnsureChoice (class *in datalad.plugin.export\_archive*), 61  
 ExportArchive.EnsureDataset (class *in datalad.plugin.export\_archive*), 61  
 ExportArchive.EnsureNone (class *in datalad.plugin.export\_archive*), 61  
 ExportArchive.EnsureStr (class *in datalad.plugin.export\_archive*), 62  
 ExportArchive.Parameter (class *in datalad.plugin.export\_archive*), 62  
 ExportToFigshare (class *in datalad.plugin.export\_to\_figshare*), 62  
 ExportToFigshare.EnsureChoice (class *in datalad.plugin.export\_to\_figshare*), 63  
 ExportToFigshare.EnsureDataset (class *in datalad.plugin.export\_to\_figshare*), 63  
 ExportToFigshare.EnsureInt (class *in datalad.plugin.export\_to\_figshare*), 63  
 ExportToFigshare.EnsureNone (class *in datalad.plugin.export\_to\_figshare*), 63  
 ExportToFigshare.EnsureStr (class *in datalad.plugin.export\_to\_figshare*), 63  
 ExportToFigshare.Parameter (class *in datalad.plugin.export\_to\_figshare*), 63  
 extract() (*in module datalad.plugin.addurls*), 58  
 extract\_metadata() (*in module datalad.api*), 186  
 ExtractedArchive (class *in datalad.support.archives*), 259

## F

`fail_with_short_help()` (in module `datalad.cmdline.main`), 276

`fake_dates_enabled` (`datalad.support.gitrepo.GitRepo` attribute), 234

`FAST_FORWARD` (`datalad.support.gitrepo.FetchInfo` attribute), 228

`FAST_FORWARD` (`datalad.support.gitrepo.PushInfo` attribute), 244

`FD_NAMES` (`datalad.cmd.WitlessProtocol` attribute), 214

`fetch()` (`datalad.support.gitrepo.GitRepo` method), 234

`fetch_()` (`datalad.support.gitrepo.GitRepo` method), 234

`FetchInfo` (class in `datalad.support.gitrepo`), 228

`FigshareRESTLaison` (class in `datalad.plugin.export_to_figshare`), 64

`File` (class in `datalad.utils`), 216

`file_basename()` (in module `datalad.utils`), 221

`file_has_content()` (`datalad.support.annexrepo.AnnexRepo` method), 247

`filter_legal_metafield()` (in module `datalad.plugin.addurls`), 58

`find()` (`datalad.support.annexrepo.AnnexRepo` method), 247

`find_files()` (in module `datalad.utils`), 221

`FINDING_SOURCES` (`datalad.support.gitrepo.GitProgress` attribute), 229

`finish()` (`datalad.support.annexrepo.ProcessAnnexProgressIndicator` method), 258

`fmt_to_name()` (in module `datalad.plugin.addurls`), 58

`for_each_ref_()` (`datalad.support.gitrepo.GitRepo` method), 234

`FORCED_UPDATE` (`datalad.support.gitrepo.FetchInfo` attribute), 228

`FORCED_UPDATE` (`datalad.support.gitrepo.PushInfo` attribute), 244

`format()` (`datalad.log.ColorFormatter` method), 216

`format()` (`datalad.plugin.addurls.Formatter` method), 58

`format()` (`datalad.plugin.addurls.RepFormatter` method), 58

`format_commit()` (`datalad.support.gitrepo.GitRepo` method), 235

`format_element()` (`datalad.utils.SequenceFormatter` method), 217

`format_field()` (`datalad.utils.SequenceFormatter` method), 217

`Formatter` (class in `datalad.plugin.addurls`), 57

`fsck()` (`datalad.support.annexrepo.AnnexRepo` method), 248

## G

`gc()` (`datalad.support.gitrepo.GitRepo` method), 235

`generate_chunks()` (in module `datalad.utils`), 221

`generate_file_chunks()` (in module `datalad.utils`), 221

`generate_uuids()` (in module `datalad.customremotes.base`), 264

`get()` (`datalad.config.ConfigManager` method), 267

`get()` (`datalad.plugin.export_to_figshare.FigshareRESTLaison` method), 64

`get()` (`datalad.support.annexrepo.AnnexRepo` method), 248

`get()` (`datalad.support.annexrepo.BatchedAnnexes` method), 258

`get()` (in module `datalad.api`), 165

`get_active_branch()` (`datalad.support.gitrepo.GitRepo` method), 235

`get_annexed_files()` (`datalad.support.annexrepo.AnnexRepo` method), 248

`get_archive()` (`datalad.support.archives.ArchivesCache` method), 259

`get_article_ids()` (`datalad.plugin.export_to_figshare.FigshareRESTLaison` method), 64

`get_autodoc()` (`datalad.plugin.add_readme.AddReadme.Parameter` method), 54

`get_autodoc()` (`datalad.plugin.addurls.Addurls.Parameter` method), 57

`get_autodoc()` (`datalad.plugin.check_dates.CheckDates.Parameter` method), 60

`get_autodoc()` (`datalad.plugin.export_archive.ExportArchive.Parameter` method), 62

`get_autodoc()` (`datalad.plugin.export_to_figshare.ExportToFigshare.Parameter` method), 63

`get_autodoc()` (`datalad.plugin.no_annex.NoAnnex.Parameter` method), 65

`get_autodoc()` (`datalad.plugin.wtf.WTF.Parameter` method), 66

`get_branch_commits_()` (`datalad.support.gitrepo.GitRepo` method), 235

`get_branches()` (`datalad.support.gitrepo.GitRepo` method), 235

`get_commands_from_groups()` (in module `datalad.cmdline.main`), 276

`get_commit_date()` (`datalad.support.gitrepo.GitRepo` method), 235

`get_content_annexinfo()` (*data-lad.support.annexrepo.AnnexRepo method*), 248  
`get_content_info()` (*data-lad.support.gitrepo.GitRepo method*), 235  
`get_contentlocation()` (*data-lad.customremotes.base.AnnexCustomRemote method*), 261  
`get_contentlocation()` (*data-lad.support.annexrepo.AnnexRepo method*), 249  
`get_convoluted_situation()` (*in module data-lad.tests.utils*), 272  
`get_corresponding_branch()` (*data-lad.support.annexrepo.AnnexRepo method*), 249  
`get_dataset_root()` (*in module datalad.utils*), 222  
`get_datasets_topdir()` (*in module data-lad.tests.utils*), 272  
`get_deeply_nested_structure()` (*in module datalad.tests.utils*), 272  
`get_description()` (*data-lad.support.annexrepo.AnnexRepo method*), 250  
`get_description_with_cmd_summary()` (*in module datalad.cmdline.main*), 277  
`get_DIRHASH()` (*data-lad.customremotes.base.AnnexCustomRemote method*), 261  
`get_encoding_info()` (*in module datalad.utils*), 222  
`get_envvars_info()` (*in module datalad.utils*), 222  
`get_extracted_file()` (*data-lad.support.archives.ExtractedArchive method*), 259  
`get_extracted_filename()` (*data-lad.support.archives.ExtractedArchive method*), 259  
`get_extracted_files()` (*data-lad.support.archives.ExtractedArchive method*), 259  
`get_file_backend()` (*data-lad.support.annexrepo.AnnexRepo method*), 250  
`get_file_key()` (*data-lad.support.annexrepo.AnnexRepo method*), 250  
`get_file_parts()` (*in module data-lad.plugin.addurls*), 59  
`get_file_size()` (*data-lad.support.annexrepo.AnnexRepo method*), 250  
`get_file_url()` (*data-lad.customremotes.archives.ArchiveAnnexCustomRemote method*), 265  
`get_files()` (*data-lad.support.gitrepo.GitRepo method*), 236  
`get_fmt_names()` (*in module data-lad.plugin.addurls*), 59  
`get_func_kwargs_doc()` (*in module datalad.utils*), 222  
`get_function_nargs()` (*in module data-lad.customremotes.base*), 264  
`get_git_attributes()` (*data-lad.support.gitrepo.GitRepo method*), 236  
`get_git_dir()` (*data-lad.support.gitrepo.GitRepo static method*), 237  
`get_git_enviro_nadjusted()` (*data-lad.cmd.GitRunner static method*), 211  
`get_git_version()` (*in module datalad.config*), 270  
`get_gitattributes()` (*data-lad.support.gitrepo.GitRepo method*), 237  
`get_groupwanted()` (*data-lad.support.annexrepo.AnnexRepo method*), 250  
`get_hexsha()` (*data-lad.support.gitrepo.GitRepo method*), 237  
`get_indexed_files()` (*data-lad.support.gitrepo.GitRepo method*), 237  
`get_ipython_shell()` (*in module datalad.utils*), 222  
`get_key_backend()` (*data-lad.support.annexrepo.AnnexRepo class method*), 250  
`get_last_commit_hexsha()` (*data-lad.support.gitrepo.GitRepo method*), 237  
`get_leading_directory()` (*data-lad.support.archives.ExtractedArchive method*), 259  
`get_linux_distribution()` (*in module data-lad.utils*), 222  
`get_logfilename()` (*in module datalad.utils*), 222  
`get_max_path_length()` (*in module data-lad.plugin.wtf*), 67  
`get_merge_base()` (*data-lad.support.gitrepo.GitRepo method*), 238  
`get_metadata()` (*data-lad.support.annexrepo.AnnexRepo method*), 250  
`get_most_obscure_supported_name()` (*in module datalad.tests.utils*), 272  
`get_mtimes_and_digests()` (*in module data-lad.tests.utils*), 272  
`get_open_files()` (*in module datalad.utils*), 222  
`get_path_prefix()` (*in module datalad.utils*), 222  
`get_preferred_content()` (*data-lad.support.annexrepo.AnnexRepo method*), 251

- `get_remote_branches()` (*datalad.support.gitrepo.GitRepo method*), 238
  - `get_remote_url()` (*datalad.support.gitrepo.GitRepo method*), 238
  - `get_remotes()` (*datalad.support.annexrepo.AnnexRepo method*), 251
  - `get_remotes()` (*datalad.support.gitrepo.GitRepo method*), 238
  - `get_repo_instance()` (*in module datalad.cmdline.helpers*), 277
  - `get_revisions()` (*datalad.support.gitrepo.GitRepo method*), 238
  - `get_size_from_key()` (*datalad.support.annexrepo.AnnexRepo static method*), 251
  - `get_special_remotes()` (*datalad.support.annexrepo.AnnexRepo method*), 251
  - `get_staged_paths()` (*datalad.support.gitrepo.GitRepo method*), 238
  - `get_submodules()` (*datalad.support.gitrepo.GitRepo method*), 238
  - `get_submodules_()` (*datalad.support.gitrepo.GitRepo method*), 239
  - `get_subpaths()` (*in module datalad.plugin.addurls*), 59
  - `get_suggestions_msg()` (*in module datalad.utils*), 222
  - `get_tags()` (*datalad.support.gitrepo.GitRepo method*), 239
  - `get_tempfile_kwargs()` (*in module datalad.utils*), 222
  - `get_timestamp_suffix()` (*in module datalad.utils*), 222
  - `get_toppath()` (*datalad.support.annexrepo.AnnexRepo class method*), 252
  - `get_toppath()` (*datalad.support.gitrepo.GitRepo class method*), 239
  - `get_trace()` (*in module datalad.utils*), 223
  - `get_tracking_branch()` (*datalad.support.annexrepo.AnnexRepo method*), 252
  - `get_tracking_branch()` (*datalad.support.gitrepo.GitRepo method*), 239
  - `get_url_parts()` (*in module datalad.plugin.addurls*), 59
  - `get_URLS()` (*datalad.customremotes.base.AnnexCustomRemote method*), 261
  - `get_urls()` (*datalad.support.annexrepo.AnnexRepo method*), 252
  - `get_value()` (*datalad.config.ConfigManager method*), 267
  - `get_value()` (*datalad.plugin.addurls.Formatter method*), 58
  - `get_value()` (*datalad.plugin.addurls.RepFormatter method*), 58
  - `get_wrapped_class()` (*in module datalad.utils*), 223
  - `getargspec()` (*in module datalad.utils*), 223
  - `getbool()` (*datalad.config.ConfigManager method*), 267
  - `getfloat()` (*datalad.config.ConfigManager method*), 267
  - `getIncludes()` (*datalad.support.configparserinc.SafeConfigParserWithIncludes static method*), 260
  - `getint()` (*datalad.config.ConfigManager method*), 267
  - `getpwd()` (*in module datalad.utils*), 223
  - `GIT_ANNEX_MIN_VERSION` (*datalad.support.annexrepo.AnnexRepo attribute*), 245
  - `git_annex_version` (*datalad.support.annexrepo.AnnexRepo attribute*), 252
  - `GitProgress` (*class in datalad.support.gitrepo*), 229
  - `GitRepo` (*class in datalad.support.gitrepo*), 229
  - `GitRunner` (*class in datalad.cmd*), 211
  - `guard_BadName()` (*in module datalad.support.gitrepo*), 244
- ## H
- `has_option()` (*datalad.config.ConfigManager method*), 267
  - `has_section()` (*datalad.config.ConfigManager method*), 267
  - `has_symlink_capability()` (*in module datalad.tests.utils*), 272
  - `HEAD_UPTODATE` (*datalad.support.gitrepo.FetchInfo attribute*), 228
  - `HEADER` (*datalad.customremotes.base.AnnexExchangeProtocol attribute*), 263
  - `heavydebug()` (*datalad.customremotes.base.AnnexCustomRemote method*), 261
  - `HelpAction` (*class in datalad.cmdline.helpers*), 277
  - `HTTPPath` (*class in datalad.tests.utils*), 270
- ## I
- `ignore_nose_capturing_stdout()` (*in module datalad.tests.utils*), 272
  - `import_module_from_file()` (*in module datalad.utils*), 223
  - `import_modules()` (*in module datalad.utils*), 223
  - `info()` (*datalad.customremotes.base.AnnexCustomRemote method*), 261



- [info\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 252  
[init\\_datalad\\_remote\(\)](#) (*in module datalad.customremotes.base*), 264  
[init\\_remote\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[initiate\(\)](#) (*datalad.customremotes.base.AnnexExchangeProtocol method*), 263  
[InnerSubmodule](#) (*class in datalad.tests.utils\_testrepos*), 275  
[install\(\)](#) (*in module datalad.api*), 167  
[integration\(\)](#) (*in module datalad.tests.utils*), 272  
[is\\_ancestor\(\)](#) (*datalad.support.gitrepo.GitRepo method*), 240  
[is\\_available\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_crippled\\_fs\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_direct\\_mode\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_explicit\\_path\(\)](#) (*in module datalad.utils*), 224  
[is\\_extracted](#) (*datalad.support.archives.ExtractedArchive attribute*), 259  
[is\\_initialized\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_interactive\(\)](#) (*in module datalad.utils*), 224  
[is\\_legal\\_metafield\(\)](#) (*in module datalad.plugin.addurls*), 59  
[is\\_managed\\_branch\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_remote\\_annex\\_ignored\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_special\\_annex\\_remote\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 253  
[is\\_under\\_annex\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 254  
[is\\_valid\\_annex\(\)](#) (*datalad.support.annexrepo.AnnexRepo method*), 254  
[is\\_valid\\_git\(\)](#) (*datalad.support.gitrepo.GitRepo method*), 240  
[is\\_valid\\_repo\(\)](#) (*datalad.support.annexrepo.AnnexRepo class method*), 254  
[is\\_valid\\_repo\(\)](#) (*datalad.support.gitrepo.GitRepo class method*), 240  
[is\\_with\\_annex\(\)](#) (*datalad.support.gitrepo.GitRepo method*), 240  
[items\(\)](#) (*datalad.config.ConfigManager method*), 267
- ## K
- [keyring\(\)](#) (*datalad.config.ConfigManager method*), 267  
[keywords](#) (*datalad.utils.ArgSpecFake attribute*), 216  
[KillOutput](#) (*class in datalad.cmd*), 212  
[known\\_failure\(\)](#) (*in module datalad.tests.utils*), 272  
[known\\_failure\\_appveyor\(\)](#) (*in module datalad.tests.utils*), 272  
[known\\_failure\\_direct\\_mode\(\)](#) (*in module datalad.tests.utils*), 272  
[known\\_failure\\_githubci\\_win\(\)](#) (*in module datalad.tests.utils*), 273  
[known\\_failure\\_v6\(\)](#) (*in module datalad.tests.utils*), 273  
[known\\_failure\\_v6\\_or\\_later\(\)](#) (*in module datalad.tests.utils*), 273  
[known\\_failure\\_windows\(\)](#) (*in module datalad.tests.utils*), 273  
[knows\\_annex\(\)](#) (*in module datalad.utils*), 224
- ## L
- [line\\_profile\(\)](#) (*in module datalad.utils*), 224  
[link\\_file\\_load\(\)](#) (*in module datalad.customremotes.archives*), 265  
[lmtime\(\)](#) (*in module datalad.utils*), 224  
[log\(\)](#) (*datalad.cmd.Runner method*), 213  
[log\\_cwd](#) (*datalad.cmd.Runner attribute*), 213  
[log\\_env](#) (*datalad.cmd.Runner attribute*), 213  
[log\\_message\(\)](#) (*datalad.tests.utils.SilentHTTPHandler method*), 270  
[log\\_outputs](#) (*datalad.cmd.Runner attribute*), 213  
[log\\_stdin](#) (*datalad.cmd.Runner attribute*), 213  
[LogLevelAction](#) (*class in datalad.cmdline.helpers*), 277  
[long\\_description\(\)](#) (*datalad.plugin.add\_readme.AddReadme.EnsureChoice method*), 54  
[long\\_description\(\)](#) (*datalad.plugin.add\_readme.AddReadme.EnsureDataset method*), 54  
[long\\_description\(\)](#) (*datalad.plugin.add\_readme.AddReadme.EnsureNone method*), 54  
[long\\_description\(\)](#) (*datalad.plugin.add\_readme.AddReadme.EnsureStr method*), 54  
[long\\_description\(\)](#) (*datalad.plugin.addurls.Addurls.EnsureChoice*

method), 56

long\_description() (data-lad.plugin.addurls.Addurls.EnsureDataset method), 56

long\_description() (data-lad.plugin.addurls.Addurls.EnsureNone method), 56

long\_description() (data-lad.plugin.addurls.Addurls.EnsureStr method), 57

long\_description() (data-lad.plugin.check\_dates.CheckDates.EnsureChoice method), 60

long\_description() (data-lad.plugin.check\_dates.CheckDates.EnsureNone method), 60

long\_description() (data-lad.plugin.check\_dates.CheckDates.EnsureStr method), 60

long\_description() (data-lad.plugin.export\_archive.ExportArchive.EnsureChoice method), 61

long\_description() (data-lad.plugin.export\_archive.ExportArchive.EnsureDataset method), 61

long\_description() (data-lad.plugin.export\_archive.ExportArchive.EnsureNone method), 62

long\_description() (data-lad.plugin.export\_archive.ExportArchive.EnsureStr method), 62

long\_description() (data-lad.plugin.export\_to\_figshare.ExportToFigshare.EnsureChoice method), 63

long\_description() (data-lad.plugin.export\_to\_figshare.ExportToFigshare.EnsureDataset method), 63

long\_description() (data-lad.plugin.export\_to\_figshare.ExportToFigshare.EnsureNone method), 63

long\_description() (data-lad.plugin.export\_to\_figshare.ExportToFigshare.EnsureStr method), 63

long\_description() (data-lad.plugin.no\_annex.NoAnnex.EnsureDataset method), 65

long\_description() (data-lad.plugin.no\_annex.NoAnnex.EnsureNone method), 65

long\_description() (data-lad.plugin.wtf.WTF.EnsureChoice method), 66

long\_description() (data-lad.plugin.wtf.WTF.EnsureDataset method), 66

long\_description() (data-lad.plugin.wtf.WTF.EnsureNone method), 66

ls() (in module datalad.api), 202

## M

main() (datalad.customremotes.base.AnnexCustomRemote method), 261

main() (in module datalad.cmdline.main), 277

main() (in module datalad.customremotes.archives), 266

main() (in module datalad.customremotes.main), 260

make\_tempfile() (in module datalad.utils), 224

map\_items() (in module datalad.utils), 225

md5sum() (in module datalad.utils), 225

merge() (datalad.support.gitrepo.GitRepo method), 240

merge\_annex() (data-lad.support.annexrepo.AnnexRepo method), 254

metadata() (in module datalad.api), 182

migrate\_backend() (data-lad.support.annexrepo.AnnexRepo method), 254

## N

name (datalad.support.gitrepo.Submodule attribute), 244

NestedDataset (class in data-lad.tests.utils\_testrepos), 275

never\_fail() (in module datalad.utils), 225

NEW\_HEAD (datalad.support.gitrepo.FetchInfo attribute), 228

NEW\_HEAD (datalad.support.gitrepo.PushInfo attribute), 244

NEW\_TAG (datalad.support.gitrepo.FetchInfo attribute), 228

NEW\_TAG (datalad.support.gitrepo.PushInfo attribute), 244

NO\_MATCH (datalad.support.gitrepo.PushInfo attribute), 244

NoAnnex (class in datalad.plugin.no\_annex), 64

NoAnnex.EnsureDataset (class in data-lad.plugin.no\_annex), 65

NoAnnex.EnsureNone (class in data-lad.plugin.no\_annex), 65

NoAnnex.Parameter (class in data-lad.plugin.no\_annex), 65

NoCapture (class in datalad.cmd), 212

nok\_startswith() (in module datalad.tests.utils), 273

not\_supported\_on\_windows() (in module datalad.utils), 225

nothing\_cm() (in module datalad.utils), 225

## O

- obtain() (*datalad.config.ConfigManager* method), 267
  - ok\_annex\_get() (in module *datalad.tests.utils*), 273
  - ok\_archives\_caches() (in module *datalad.tests.utils*), 273
  - ok\_broken\_symlink() (in module *datalad.tests.utils*), 273
  - ok\_clean\_git() (in module *datalad.tests.utils*), 273
  - ok\_endswith() (in module *datalad.tests.utils*), 274
  - ok\_exists() (in module *datalad.tests.utils*), 274
  - ok\_file\_has\_content() (in module *datalad.tests.utils*), 274
  - ok\_file\_under\_git() (in module *datalad.tests.utils*), 274
  - ok\_generator() (in module *datalad.tests.utils*), 274
  - ok\_git\_config\_not\_empty() (in module *datalad.tests.utils*), 274
  - ok\_good\_symlink() (in module *datalad.tests.utils*), 274
  - ok\_startswith() (in module *datalad.tests.utils*), 274
  - ok\_symlink() (in module *datalad.tests.utils*), 274
  - OP\_MASK (*datalad.support.gitrepo.GitProgress* attribute), 229
  - open\_r\_encdetect() (in module *datalad.utils*), 225
  - optional\_args() (in module *datalad.utils*), 225
  - options() (*datalad.config.ConfigManager* method), 268
- ## P
- parser\_add\_common\_opt() (in module *datalad.cmdline.helpers*), 277
  - partition() (in module *datalad.utils*), 225
  - patch\_config() (in module *datalad.tests.utils*), 274
  - path (*datalad.support.archives.ArchivesCache* attribute), 259
  - path (*datalad.support.archives.ExtractedArchive* attribute), 259
  - path (*datalad.support.gitrepo.Submodule* attribute), 244
  - path (*datalad.tests.utils\_testrepos.InnerSubmodule* attribute), 275
  - path (*datalad.tests.utils\_testrepos.TestRepo* attribute), 276
  - path\_is\_subpath() (in module *datalad.utils*), 226
  - path\_startswith() (in module *datalad.utils*), 226
  - pipe\_data\_received() (*datalad.cmd.KillOutput* method), 212
  - pipe\_data\_received() (*datalad.cmd.WitlessProtocol* method), 214
  - pipe\_data\_received() (*datalad.support.gitrepo.GitProgress* method), 229
  - populate() (*datalad.tests.utils\_testrepos.BasicAnnexTestRepo* method), 275
  - populate() (*datalad.tests.utils\_testrepos.BasicGitTestRepo* method), 275
  - populate() (*datalad.tests.utils\_testrepos.NestedDataset* method), 276
  - populate() (*datalad.tests.utils\_testrepos.SubmoduleDataset* method), 276
  - populate() (*datalad.tests.utils\_testrepos.TestRepo* method), 276
  - posix\_relpath() (in module *datalad.utils*), 226
  - post() (*datalad.plugin.export\_to\_figshare.FigshareRESTLaison* method), 64
  - precommit() (*datalad.support.annexrepo.AnnexRepo* method), 254
  - precommit() (*datalad.support.gitrepo.GitRepo* method), 240
  - probe\_known\_failure() (in module *datalad.tests.utils*), 274
  - procl() (*datalad.cmd.BatchedCommand* method), 211
  - proc\_err (*datalad.cmd.KillOutput* attribute), 212
  - proc\_err (*datalad.cmd.StdErrCapture* attribute), 214
  - proc\_err (*datalad.cmd.StdOutErrCapture* attribute), 214
  - proc\_err (*datalad.cmd.WitlessProtocol* attribute), 214
  - proc\_err (*datalad.support.gitrepo.GitProgress* attribute), 229
  - proc\_out (*datalad.cmd.KillOutput* attribute), 212
  - proc\_out (*datalad.cmd.StdOutCapture* attribute), 214
  - proc\_out (*datalad.cmd.StdOutErrCapture* attribute), 214
  - proc\_out (*datalad.cmd.WitlessProtocol* attribute), 215
  - proc\_out (*datalad.support.gitrepo.StdOutCaptureWithGitProgress* attribute), 244
  - process\_exited() (*datalad.cmd.WitlessProtocol* method), 215
  - process\_exited() (*datalad.support.gitrepo.GitProgress* method), 229
  - ProcessAnnexProgressIndicators (class in *datalad.support.annexrepo*), 258
  - progress() (*datalad.customremotes.base.AnnexCustomRemote* method), 261
  - protocol (*datalad.cmd.Runner* attribute), 213
  - publish() (in module *datalad.api*), 169
  - pull() (*datalad.support.gitrepo.GitRepo* method), 240
  - push() (*datalad.support.gitrepo.GitRepo* method), 240
  - push\_() (*datalad.support.gitrepo.GitRepo* method), 241
  - PushInfo (class in *datalad.support.gitrepo*), 243
  - put() (*datalad.plugin.export\_to\_figshare.FigshareRESTLaison* method), 64
  - put\_file\_under\_git() (in module *datalad.tests.utils*), 274

**Q**

quote\_cmdlinearg() (in module *datalad.utils*), 226

**R**

re\_op\_absolute (data-lad.support.gitrepo.GitProgress attribute), 229

re\_op\_relative (data-lad.support.gitrepo.GitProgress attribute), 229

read() (*datalad.customremotes.base.AnnexCustomRemote* method), 261

read() (*datalad.support.configparserinc.SafeConfigParserWithIncludes* method), 260

read\_csv\_lines() (in module *datalad.utils*), 226

readline\_json() (in module *datalad.support.annexrepo*), 258

readline\_rstripped() (in module *datalad.cmd*), 215

readlines\_until\_ok\_or\_failed() (in module *datalad.support.annexrepo*), 258

RECEIVING (*datalad.support.gitrepo.GitProgress* attribute), 229

records\_callable (data-lad.customremotes.base.AnnexExchangeProtocol attribute), 263

records\_ext\_commands (data-lad.customremotes.base.AnnexExchangeProtocol attribute), 264

REJECTED (*datalad.support.gitrepo.FetchInfo* attribute), 229

REJECTED (*datalad.support.gitrepo.PushInfo* attribute), 244

reload() (*datalad.config.ConfigManager* method), 268

REMOTE\_FAILURE (*datalad.support.gitrepo.PushInfo* attribute), 244

REMOTE\_REJECTED (*datalad.support.gitrepo.PushInfo* attribute), 244

remove() (*datalad.support.annexrepo.AnnexRepo* method), 254

remove() (*datalad.support.gitrepo.GitRepo* method), 241

remove() (in module *datalad.api*), 170

remove\_branch() (*datalad.support.gitrepo.GitRepo* method), 241

remove\_remote() (*datalad.support.gitrepo.GitRepo* method), 241

remove\_section() (*datalad.config.ConfigManager* method), 268

rename\_section() (*datalad.config.ConfigManager* method), 268

RepFormatter (class in *datalad.plugin.addurls*), 58

repo (*datalad.support.gitrepo.GitRepo* attribute), 241

Repo() (in module *datalad.support.gitrepo*), 244

REPO\_CLASS (*datalad.tests.utils\_testrepos.BasicAnnexTestRepo* attribute), 275

REPO\_CLASS (*datalad.tests.utils\_testrepos.BasicGitTestRepo* attribute), 275

REPO\_CLASS (*datalad.tests.utils\_testrepos.TestRepo* attribute), 276

repo\_info() (*datalad.support.annexrepo.AnnexRepo* method), 254

repository\_versions (*datalad.support.annexrepo.AnnexRepo* attribute), 255

req\_CHECKPRESENT() (*datalad.customremotes.archives.ArchiveAnnexCustomRemote* method), 265

req\_CHECKPRESENT() (*datalad.customremotes.base.AnnexCustomRemote* method), 261

req\_CHECKURL() (*datalad.customremotes.archives.ArchiveAnnexCustomRemote* method), 265

req\_CHECKURL() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_CLAIMURL() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_EXPORTSUPPORTED() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_GETAVAILABILITY() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_GETCOST() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_INITREMOTE() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_PREPARE() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_REMOVE() (*datalad.customremotes.archives.ArchiveAnnexCustomRemote* method), 265

req\_REMOVE() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_TRANSFER() (*datalad.customremotes.base.AnnexCustomRemote* method), 262

req\_WHEREIS() (*datalad.customremotes.archives.ArchiveAnnexCustomRemote* method), 265

req\_WHEREIS () (*data-lad.customremotes.base.AnnexCustomRemote method*), 262  
 rerun () (*in module datalad.api*), 189  
 RESOLVING (*datalad.support.gitrepo.GitProgress attribute*), 229  
 result\_renderer (*data-lad.plugin.check\_dates.CheckDates attribute*), 61  
 result\_renderer (*datalad.plugin.wtf.WTF attribute*), 67  
 rewrite\_url () (*datalad.config.ConfigManager method*), 269  
 rewrite\_url () (*in module datalad.config*), 270  
 RFC  
     RFC 822, 49  
 rm\_url () (*datalad.support.annexrepo.AnnexRepo method*), 255  
 rmdir () (*in module datalad.utils*), 226  
 rmtree () (*in module datalad.utils*), 226  
 rotree () (*in module datalad.utils*), 227  
 run () (*datalad.cmd.GitRunner method*), 211  
 run () (*datalad.cmd.Runner method*), 213  
 run () (*datalad.cmd.WitlessRunner method*), 215  
 run () (*in module datalad.api*), 187  
 run\_async\_cmd () (*in module datalad.cmd*), 215  
 run\_gitcommand\_on\_file\_list\_chunks () (*in module datalad.cmd*), 216  
 run\_procedure () (*in module datalad.api*), 191  
 run\_via\_pbs () (*in module datalad.cmdline.helpers*), 277  
 Runner (*class in datalad.cmd*), 212

**S**

safe\_print () (*in module datalad.utils*), 227  
 SafeConfigParserWithIncludes (*class in datalad.support.configparserinc*), 260  
 SafeDelCloseMixin (*class in datalad.cmd*), 214  
 save () (*datalad.support.gitrepo.GitRepo method*), 241  
 save () (*in module datalad.api*), 171  
 save\_ () (*datalad.support.gitrepo.GitRepo method*), 241  
 saved\_generator () (*in module datalad.utils*), 227  
 search () (*in module datalad.api*), 178  
 SECTION\_NAME (*data-lad.support.configparserinc.SafeConfigParserWithIncludes attribute*), 260  
 sections () (*datalad.config.ConfigManager method*), 269  
 send () (*datalad.customremotes.base.AnnexCustomRemote method*), 262  
 send\_unsupported () (*data-lad.customremotes.base.AnnexCustomRemote method*), 263  
 SequenceFormatter (*class in datalad.utils*), 217  
 set () (*datalad.config.ConfigManager method*), 269  
 set\_annex\_version () (*in module datalad.tests.utils*), 274  
 set\_date () (*in module datalad.tests.utils*), 274  
 set\_default\_backend () (*data-lad.support.annexrepo.AnnexRepo method*), 255  
 set\_gitattributes () (*data-lad.support.gitrepo.GitRepo method*), 241  
 set\_groupwanted () (*data-lad.support.annexrepo.AnnexRepo method*), 255  
 set\_metadata () (*data-lad.support.annexrepo.AnnexRepo method*), 255  
 set\_metadata\_ () (*data-lad.support.annexrepo.AnnexRepo method*), 256  
 set\_preferred\_content () (*data-lad.support.annexrepo.AnnexRepo method*), 256  
 set\_remote\_dead () (*data-lad.support.annexrepo.AnnexRepo method*), 256  
 set\_remote\_url () (*data-lad.support.annexrepo.AnnexRepo method*), 256  
 set\_remote\_url () (*data-lad.support.gitrepo.GitRepo method*), 242  
 setup\_exceptionhook () (*in module datalad.utils*), 227  
 setup\_parser () (*in module datalad.cmdline.main*), 277  
 setup\_parser () (*in module data-lad.customremotes.main*), 260  
 short\_description () (*data-lad.plugin.add\_readme.AddReadme.EnsureChoice method*), 54  
 short\_description () (*data-lad.plugin.add\_readme.AddReadme.EnsureDataset method*), 54  
 short\_description () (*data-lad.plugin.add\_readme.AddReadme.EnsureNone method*), 54  
 short\_description () (*data-lad.plugin.add\_readme.AddReadme.EnsureStr method*), 54  
 short\_description () (*data-lad.plugin.addurls.Addurls.EnsureChoice method*), 56  
 short\_description () (*data-lad.plugin.addurls.Addurls.EnsureDataset method*), 56



- subdataset, [71](#)
  - subdatasets() (in module *datalad.api*), [206](#)
  - Submodule (class in *datalad.support.gitrepo*), [244](#)
  - SubmoduleDataset (class in *datalad.tests.utils\_testrepos*), [276](#)
  - superdataset, [71](#)
  - SUPPORTED\_SCHEMES (data-lad.customremotes.archives.ArchiveAnnexCustomRemote attribute), [264](#)
  - SUPPORTED\_SCHEMES (data-lad.customremotes.base.AnnexCustomRemote attribute), [261](#)
  - supports\_direct\_mode (data-lad.support.annexrepo.AnnexRepo attribute), [256](#)
  - supports\_unlocked\_pointers (data-lad.support.annexrepo.AnnexRepo attribute), [256](#)
  - swallow\_logs() (in module *datalad.utils*), [227](#)
  - swallow\_outputs() (in module *datalad.utils*), [227](#)
  - sync() (*datalad.support.annexrepo.AnnexRepo* method), [256](#)
- T**
- tag() (*datalad.support.gitrepo.GitRepo* method), [242](#)
  - TAG\_UPDATE (*datalad.support.gitrepo.FetchInfo* attribute), [229](#)
  - test() (in module *datalad.api*), [209](#)
  - TestRepo (class in *datalad.tests.utils\_testrepos*), [276](#)
  - to\_options() (in module *datalad.support.gitrepo*), [244](#)
  - token (*datalad.plugin.export\_to\_figshare.FigshareRESTLiaison* attribute), [64](#)
  - TOKEN\_SEPARATOR (*datalad.support.gitrepo.GitProgress* attribute), [229](#)
  - try\_multiple() (in module *datalad.utils*), [227](#)
  - try\_multiple\_dec() (in module *datalad.utils*), [227](#)
- U**
- unannex() (*datalad.support.annexrepo.AnnexRepo* method), [257](#)
  - uninstall() (in module *datalad.api*), [175](#)
  - unique() (in module *datalad.utils*), [228](#)
  - unlink() (in module *datalad.utils*), [228](#)
  - unlock() (*datalad.support.annexrepo.AnnexRepo* method), [257](#)
  - unlock() (in module *datalad.api*), [176](#)
  - unset() (*datalad.config.ConfigManager* method), [269](#)
  - untracked\_files (*datalad.support.gitrepo.GitRepo* attribute), [243](#)
  - UP\_TO\_DATE (*datalad.support.gitrepo.PushInfo* attribute), [244](#)
  - update() (in module *datalad.api*), [173](#)
  - update\_ref() (*datalad.support.gitrepo.GitRepo* method), [243](#)
  - update\_remote() (*datalad.support.gitrepo.GitRepo* method), [243](#)
  - update\_submodule() (*datalad.support.gitrepo.GitRepo* method), [243](#)
  - updated() (in module *datalad.utils*), [228](#)
  - updated\_file() (*datalad.plugin.export\_to\_figshare.FigshareRESTLiaison* method), [64](#)
  - url (*datalad.support.gitrepo.Submodule* attribute), [244](#)
  - url (*datalad.tests.utils\_testrepos.InnerSubmodule* attribute), [275](#)
  - url (*datalad.tests.utils\_testrepos.TestRepo* attribute), [276](#)
  - URL\_PREFIX (*datalad.customremotes.archives.ArchiveAnnexCustomRemote* attribute), [264](#)
  - URL\_SCHEME (*datalad.customremotes.archives.ArchiveAnnexCustomRemote* attribute), [264](#)
  - usecase() (in module *datalad.tests.utils*), [275](#)
  - uuid (*datalad.support.annexrepo.AnnexRepo* attribute), [257](#)
- V**
- varargs (*datalad.utils.ArgSpecFake* attribute), [216](#)
- W**
- WEB\_UUID (*datalad.support.annexrepo.AnnexRepo* attribute), [245](#)
  - whereis() (*datalad.support.annexrepo.AnnexRepo* method), [257](#)
  - whereis\_pathsep() (in module *datalad.utils*), [228](#)
  - WitlessProtocol (class in *datalad.cmd*), [214](#)
  - WitlessRunner (class in *datalad.cmd*), [215](#)
  - write\_entries() (*datalad.customremotes.base.AnnexExchangeProtocol* method), [264](#)
  - write\_section() (*datalad.customremotes.base.AnnexExchangeProtocol* method), [264](#)
  - WRITING (*datalad.support.gitrepo.GitProgress* attribute), [229](#)
  - WTF (class in *datalad.plugin.wtf*), [66](#)
  - WTF.EnsureChoice (class in *datalad.plugin.wtf*), [66](#)
  - WTF.EnsureDataset (class in *datalad.plugin.wtf*), [66](#)
  - WTF.EnsureNone (class in *datalad.plugin.wtf*), [66](#)
  - WTF.Parameter (class in *datalad.plugin.wtf*), [66](#)
- Y**
- yield\_() (*datalad.cmd.BatchedCommand* method), [211](#)