datalad-registry

http://registry.datalad.org

https://github.com/datalad/datalad-registry

Goal

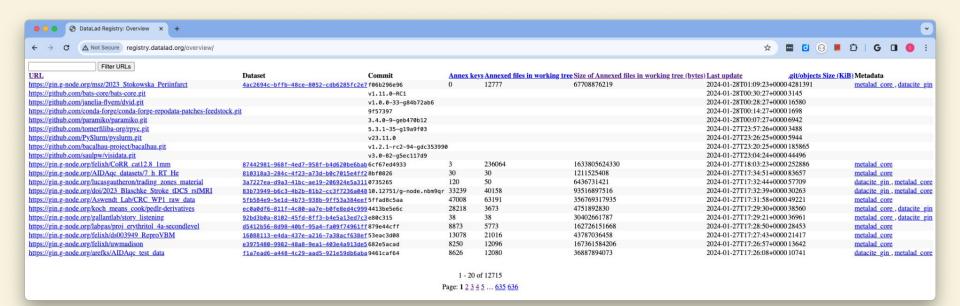
- 1. To allow discovery of data from a large collection of datasets, Datalad datasets and git repos touched by `datalad run`
- 2. To allow discovery of reuses of datasets as a submodule of a super dataset

Automation as the Defining Principle

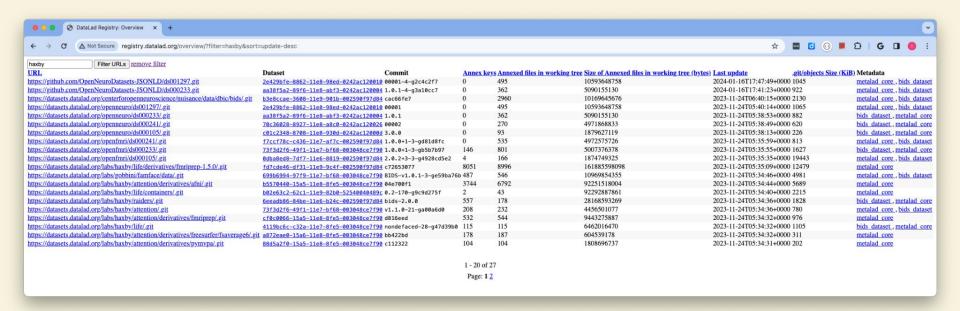
- In discovery of datasets
- In extraction of metadata from the datasets
- Updating of the dataset state and metadata

Providing up-to-date info on 12,000+ datasets

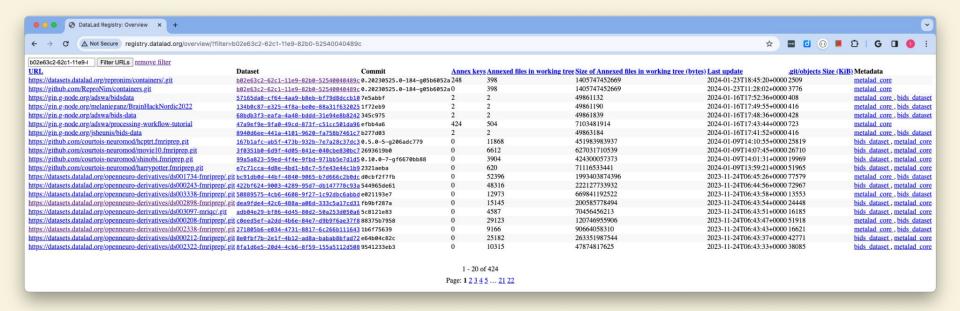
- All datasets on https://datasets.datalad.org/
- 2. Datasets in the wild as they are discovered by https://github.com/datalad/datalad-usage-dashboard
 - a. Datasets that reside on GitHub
 - b. Datasets that reside on GIN
- 3. Information of each dataset includes metadata of the dataset
 - a. "metalad core"
 - b. "metalad_studyminimeta"
 - c. "datacite_gin"
 - d. "bids_dataset"
 - e. "dandi" (not from datalad-metalad just loaded JSON dumps)



A simple interface that provides basic search



The search interface can also be used to search for forks and reuses of a dataset.

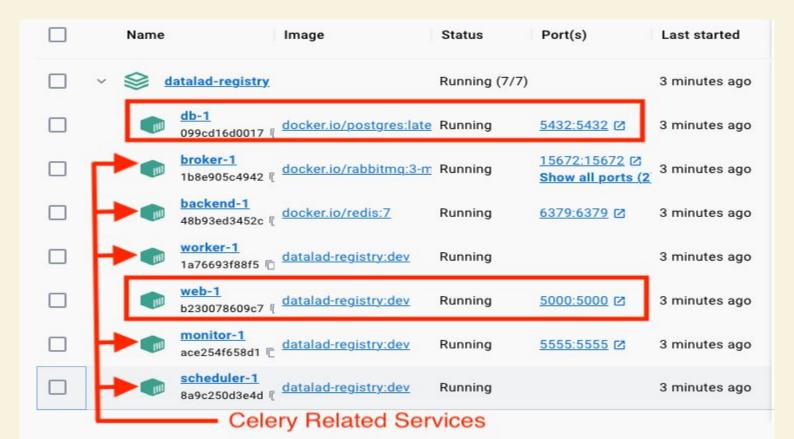


Current Developer View

An API with OpenAPI V3 documentation at http://registry.datalad.org/openapi/

- The OpenAPI V3 documentation is suitable for human consumption through Swagger, ReDoc, or RepiDoc
- The OpenAPI V3 documentation can be used by other parties generate code to access the API
- Supported operations
 - Registering new datasets (only on typhon, not on public <u>registry.datalad.org</u>)
 - Querying for datasets based on restrictions
 - Fetching datasets and metadata by internal IDs

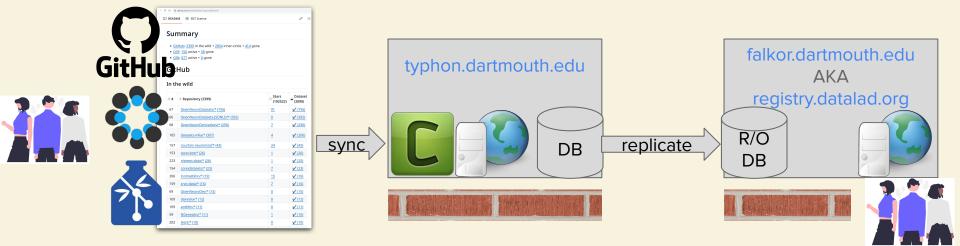
Current Developer View



Current Developer View

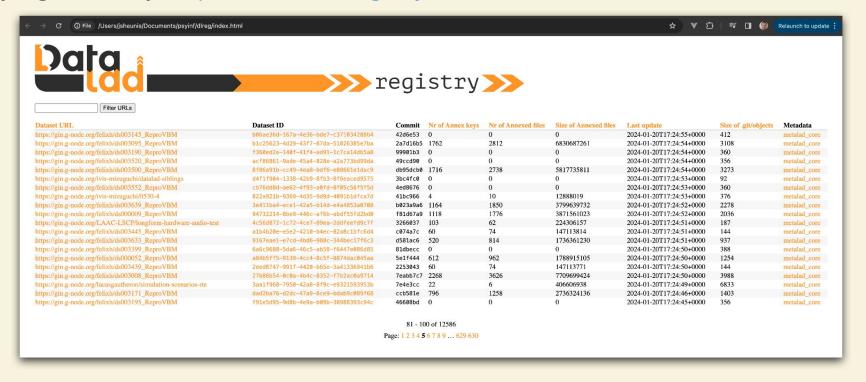
Ability to launch a **read-only instance** that consists of only the web and DB services

- Protect the main instance depicted in the previous slide from abuse
- Multiple read-only instances can split the workload
- A read-only instance is set up to be in sync with the main instance



Ongoing Development

Styling facelift by Stephan: datalad-registry/issues/290



Ongoing Development

Integration with datalad-catalog: datalad-registry/pull/278

- State: one metadata field at a time
- Goal: Provide datalad-catalog landing pages for datasets

Upcoming Development

- Elaborate search scheme (main focus)
- Dataset level metadata from other extractors
- File level metadata extractions
- Anticipated challenge: scalability
 - Picking the right database technology
 - Search through a large amount of semi-structured data and producing result in a reasonable amount of time.