

Data Search and Exploration using EarthCube Data Discovery Studio

Ilya Zaslavsky¹, David Valentine¹, Stephen Richard², and Ouida Meier³

¹UC San Diego

²US Geoscience Information Network


³University of Hawaii

November 23, 2022

Abstract

The EarthCube Data Discovery Studio (DDStudio) integrates several technical components into an end-to-end data discovery and exploration system. Beyond supporting dataset search across multiple data sources, it lets geoscientists explore the data using Jupyter notebooks; organize the discovered datasets into thematic collections which can be shared with other users; edit metadata records and contribute metadata describing additional datasets; and examine provenance and validate automated metadata enhancements. DDStudio provides access to 1.67 million metadata records from 40+ geoscience repositories, which are automatically enhanced and exposed via standard interfaces in both ISO-19115 and in schema.org markup; the latter can be used by commercial search engines (Google, Bing) to index DDStudio content. For geoscience end users, DDStudio provides a custom Geoportal-based user interface which enables spatio-temporal, faceted, and full-text search, and provides access to additional functions listed above. Key project accomplishments over the last year include: - User interface improvements, based on design advice from a Science Gateways Community Institute (SGCI) usability team, who conducted user interviews, performed usability testing, and analyzed a dozen of other search portals to identify the most useful features. This work resulted in a streamlined user interface, particularly in presentation of search results and in management of thematic collections. - The earlier effort to publish DDStudio content using schema.org markup resulted in significant usage increase. With over 900K records indexed by Google, nearly half of the roughly 1000 unique users per month are now accessing DDStudio via referrals from Google. - The added ability to harvest and process JSON-LD metadata makes it possible to integrate EarthCube GeoCodes content into DDStudio, and work with this content using DDStudio's user interface. - New application domains include joint work with the library community, and interoperation with DataMed, a similar system that indexes 2.3 million biomedical datasets.


Data Search and Exploration using EarthCube Data Discovery Studio



Data Search and Exploration using EarthCube Data Discovery Studio

Ilya Zaslavsky (1), David Valentine (1), Stephen Richard (2), Ouida Meier (3)

1: UC San Diego; 2: US Geoscience Information Network; 3: University of Hawaii



DDStudio: Beyond Search

EarthCube Data Discovery Studio is a platform for finding and exploring geoscience data.

DDStudio supports various modes of data discovery.

OPEN

Metadata Enhancement

DDStudio implements a metadata augmentation pipeline that uses text analytics and geoscience ontologies to generate keywords, spatial and temporal extents, and organization identifiers.

OPEN

Schema.org Markup

All metadata records in DDStudio are presented using schema.org markup and referenced in sitemaps submitted to Google. By now, Google indexed 900K records from DDStudio.

OPEN

Jupyter Integration

DDStudio prototyped workflows for a seamless transition from geoscience data discovery to research. Using the "Studio" link from search results in the discovery interface, users can launch Jupyter notebooks residing on

OPEN

Collection Management

DDStudio implemented a custom version of ESRi Geoportal Server, with the added capability that lets users save any discovered metadata records into collections. The found records can be added to new or existing collection. Users can also export collections and share them with collaborators, or import collections developed by others.

OPEN

Watch DataDiscoveryStudio.org in Action!

OPEN

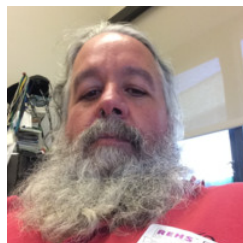
SGCI Collaboration

The Science Gateways Community Institute (SGCI) has served as an important partner to DDStudio over the last year.

OPEN

Ilya Zaslavsky (1), David Valentine (1), Stephen Richard (2), Ouida Meier (3)

1: UC San Diego; 2: US Geoscience Information Network; 3: University of Hawaii



PRESENTED AT:

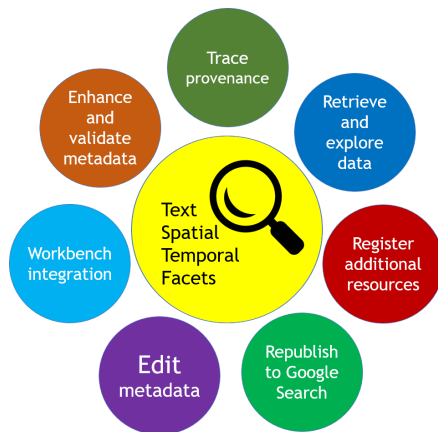


2020 EarthCube Annual Meeting

Virtual – June 18, 2020

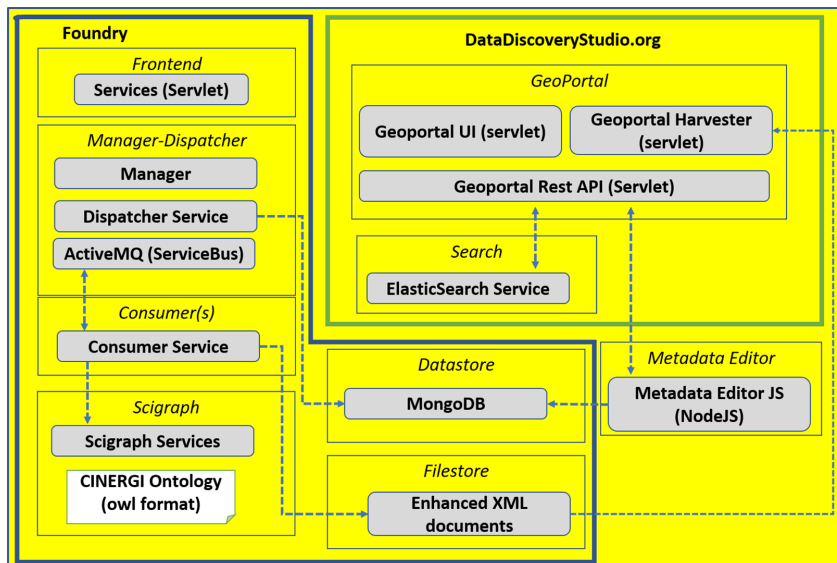
DDSTUDIO: BEYOND SEARCH

EarthCube Data Discovery Studio is a platform for finding and exploring geoscience data.



DDStudio supports various modes of data discovery, including spatio-temporal, faceted, and full-text search. Beyond search, it includes capabilities for metadata enhancement, editing, and data contribution. It lets users organize the discovered data into collections, and launch Jupyter notebooks for the datasets or dataset collections. All metadata records are exported as schema.org documents for indexing by commercial search engines.

Key DDStudio components

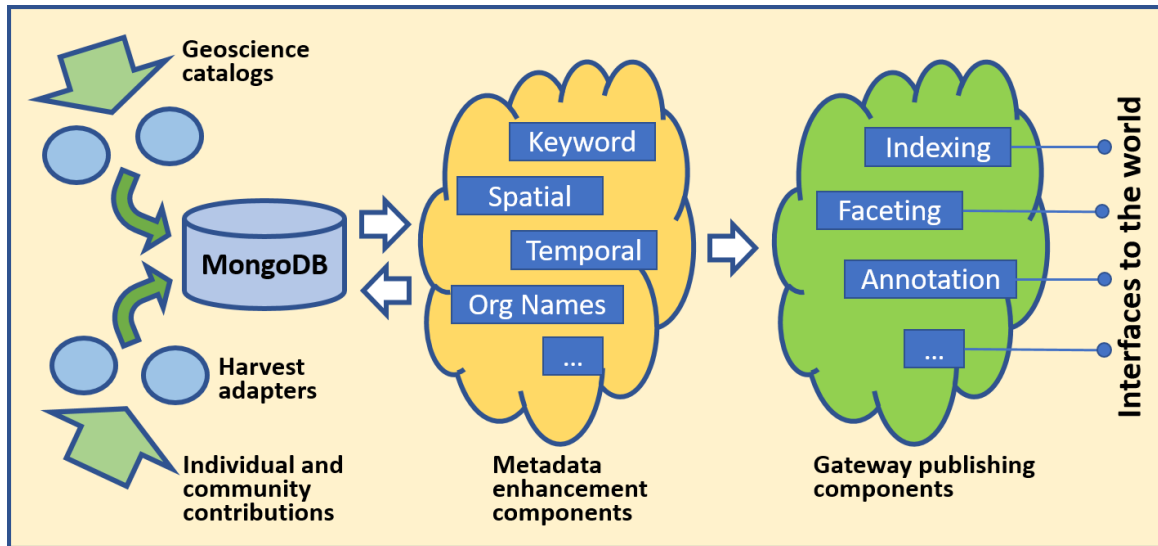


DDStudio is deployed on a set of virtual machines in OpenStack cloud platform. The key component is the DDStudio Foundry, which implements the metadata enhancement pipeline using service messaging bus, a producer/consumer architecture, and a MongoDB data store. Additional virtual machines are used for the user portal, Elasticsearch, and the Metadata editor.

Acknowledgment. NSF support (awards 1639764, 1639775) is gratefully acknowledged.

METADATA ENHANCEMENT

DDStudio implements a metadata augmentation pipeline that uses text analytics and geoscience ontologies to generate keywords, spatial and temporal extents, and organization identifiers.

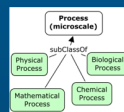
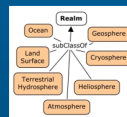


The keyword enhancer relies on GeoSciGraph services operating over 20+ geoscience ontologies, reconciled and organized into a faceted search hierarchy.

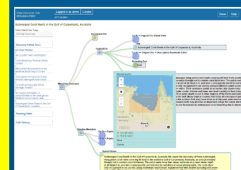
GeoSciGraph ontology management system provides semantic infrastructure. It relies on a cross-domain ontology of geoscience terms, integrating several independently developed ontologies or taxonomies

Some included ontologies:

- SWEET
- ENVO
- CHEBI
- YAGO (geo features)
- NASA GCMD (equipment, providers)
- GeoSciML
- Geochronology
- EDAM Bioinformatics (software terms and operations)
- Also: VIAF



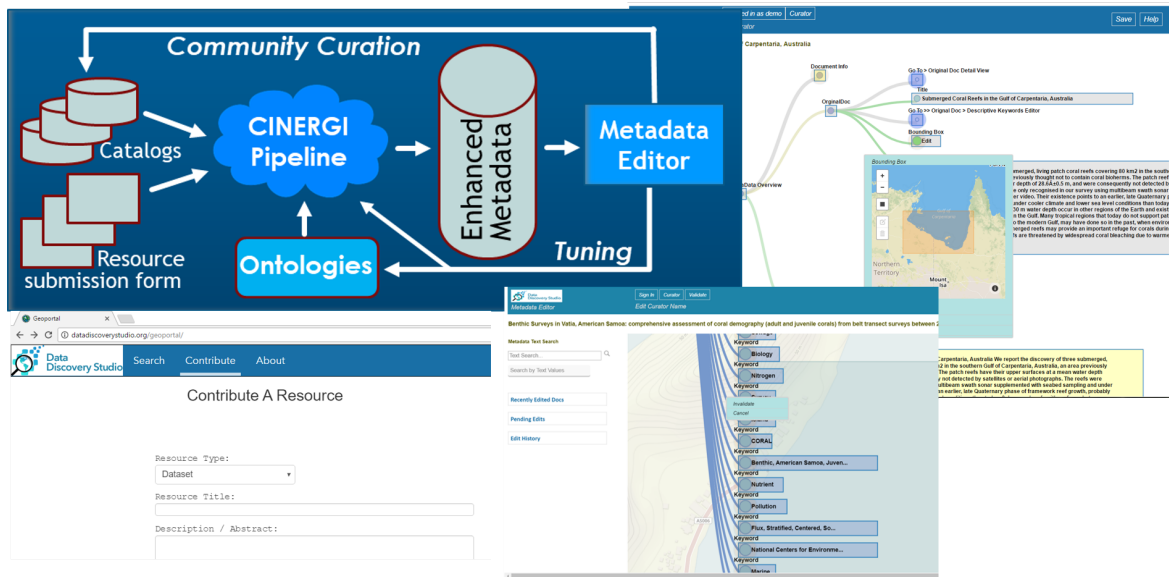
Added annotation properties for combining ontology fragments
(*cinergiFacet*, *cinergiParent*)



Metadata editor
Approve or discard semantic annotations

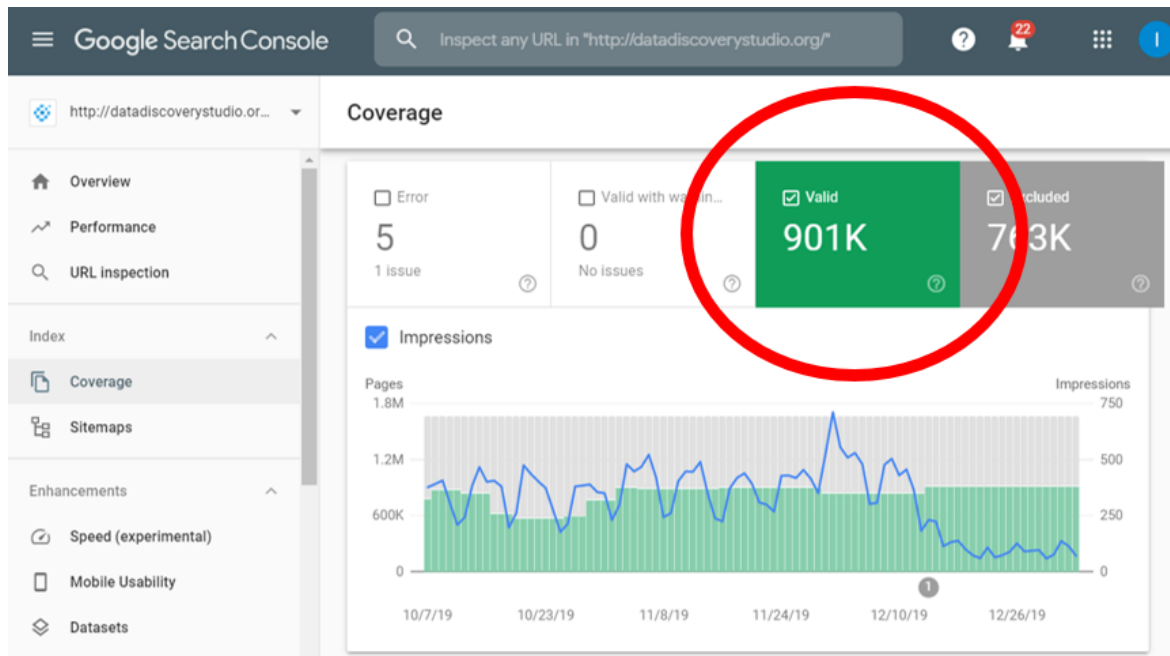


DDStudio experience suggests a scalable metadata curation model, which does not sacrifice domain semantics when improving metadata for discovery. In this model, the first curation step is performed by the automated metadata augmentation pipeline. Then a repository curator can examine the results, invalidating incorrect assignments. This review triggers ontology updates and re-processing, at the same time creating a labeled training set for supervised learning.

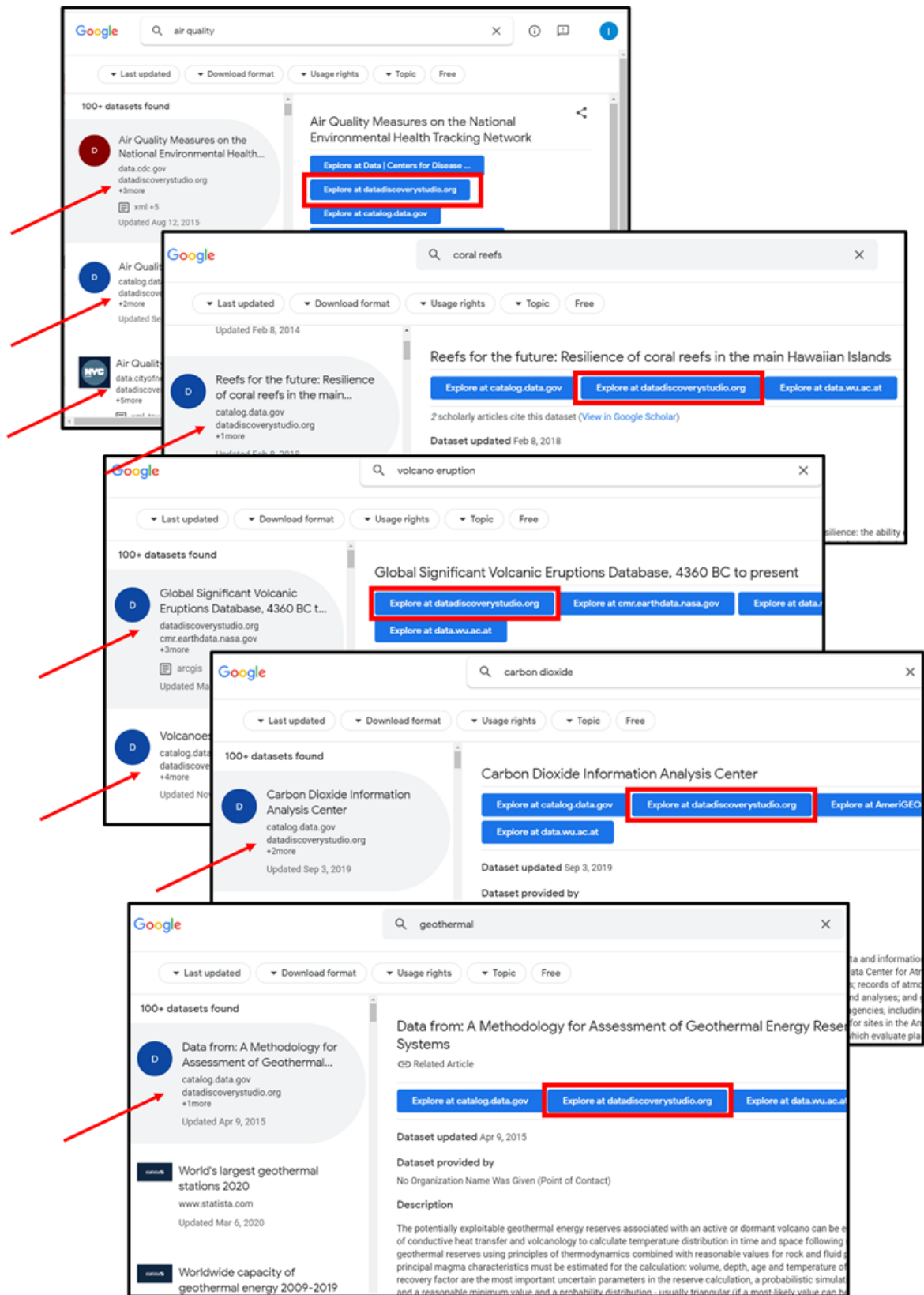


SCHEMA.ORG MARKUP

All metadata records in DDStudio are presented using schema.org markup and referenced in sitemaps submitted to Google. By now, Google indexed 900K records from DDStudio.

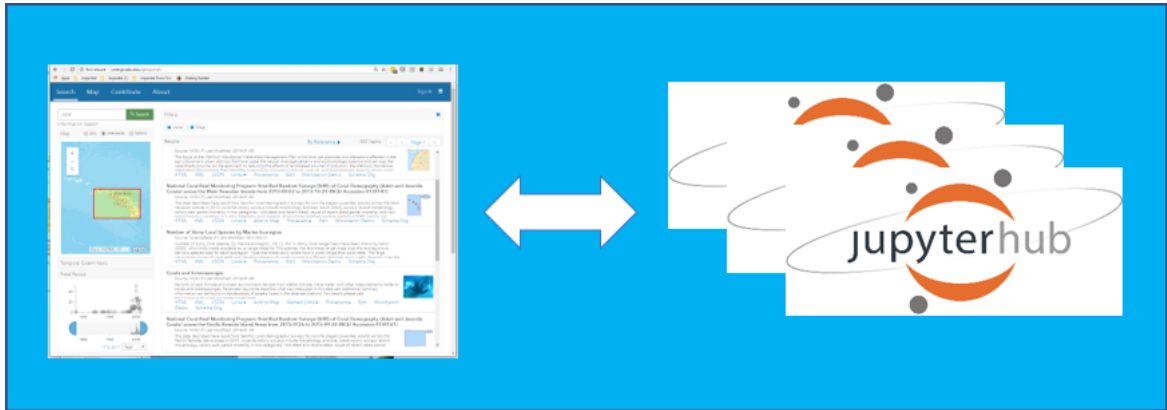


As a result, it is likely that you will be referred to DDStudio when you enter geoscience terms in Google Dataset Search. Please try it yourself!

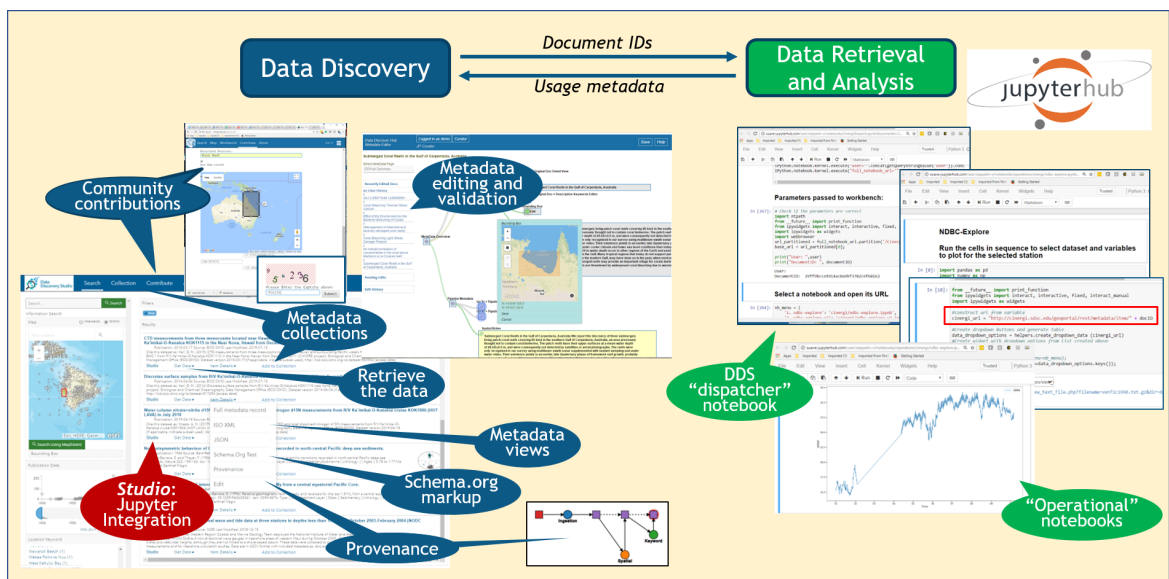


About 40% of DDStudio users have visited the platform via Google referrals.

JUPYTER INTEGRATION



DDStudio prototyped workflows for a seamless transition from geoscience data discovery to research. Using the "Studio" link from search results in the discovery interface, users can launch Jupyter notebooks residing on several JupyterHub servers. The initial link points to a "dispatcher" notebook, which calls other notebooks that implement visualizations, analytical techniques, or models.



The DDStudio system and the collection of Jupyter notebooks are loosely coupled and communicate via standard protocols, which can be used to integrate arbitrary discovery systems and research workbenches using a lightweight interface.

COLLECTION MANAGEMENT

DDStudio implemented a custom version of ESRI Geoportal Server, with the added capability that lets users save any discovered metadata records into collections. The found records can be added to new or existing collection. Users can also export collections and share them with collaborators, or import collections developed by others.

The screenshot displays the DDStudio web application's 'Collection' management page. The top navigation bar includes links for Search, Collection, Contribute, About, Tour, and Help. The main content area is divided into two panels. The left panel, titled 'Saved Collections', shows a list of collections with 'Unassigned Results' and a 'Show All' button. Below this are buttons for 'New Collection' and 'Export Collection'. The 'Collection Import' section includes a 'Collection Info (Merge or Overwrite)' dropdown set to 'Merge', a 'Select a CSV File' button, and an 'Import File' button. The right panel, titled 'Collection Items', shows a list of three items. The first item is 'A review of the Grenville orogen in its North American type area', which is a collection of 'Unassigned Records'. The description for this item states: 'In a reconstructed supercontinent assembly at ~0.9 Ga, the Grenville orogen extends from Scandinavia through North America and Antarctica to Australia. Part of it, the 2000 km long Grenville Province, exposed in the southeastern Canadian Shield, is large enough to allow a comprehensive view of its tectonic character. It has an orogen-parallel zonation: older, reworked crust, representing Archean and Palaeoproterozoic orogens exposed in adjacent parts of the shield, is restricted to its northwest side; supracrustal and plutonic rocks of Grenvillian age (~1.3-0.95 Ga) are limited to the southeastern half. The latter lie on or within late Palaeoproterozoic and earlier Mesoproterozoic crust, which is the deformed, temporal equivalent of terranes that form a substantial part of the buried North American craton south of the shield. A pre-Grenvillian period of quiescence at ~1.5 Ga may have followed an earlier continental assembly. Grenvillian calc-alkaline igneous rocks, limited in volume and distribution, represent arc accretion that terminated with ocean closure by ~1.2 Ga. New crust was added after continent-continent collision and attendant crustal thickening by emplacement'. Below the description are buttons for 'Add to Collection', 'Remove from Collection', and 'Remove Saved record'. The second item is 'Generalized Geology of Europe including Turkey (geo4_2l)', also a collection of 'Unassigned Records'. The description states: 'This coverage includes arcs, polygons, and polygon labels that describe the generalized geologic age of surface outcrops of bedrock of Europe including Turkey (Albania, Andorra, Austria, Belgium, Bosnia and Herzegovina, Bulgaria, Croatia, Cyprus, Czech Republic, Denmark, Finland, France, Germany, Greece, Hungary, Iceland, Ireland, Italy, Liechtenstein, Luxembourg, The Former Yugoslav Republic of Macedonia, Malta, Monaco, Netherlands, Norway, Poland, Portugal, Romania, San Marino, Serbia and Montenegro, Slovakia, Slovenia, Spain, Sweden, Switzerland, Turkey, United Kingdom and Vatican City.) It also includes shorelines and inland water bodies'. Below the description are buttons for 'Add to Collection', 'Remove from Collection', and 'Remove Saved record'. The third item is 'Grenville-age belts and associated older terranes in Australia and Antarctica', also a collection of 'Unassigned Records'. The description states: 'During 1.3-1.0 Ga, multiple magmatic/tectonic events, causally linked to extensive mantle melting, markedly affected the Musgrave Block, the Albany-Fraser belt and parts of the east Antarctic shield. It is possible that southern parts of the central Australian Arunta Block were also affected by this event. Equivalent terranes can be juxtaposed in conventional Gondwana reconstructions of southwestern Australia and Antarctica. Effects of a ~1060 Ma Gondwana-wide magmatism include areally extensive dyke swarms and volcanics of the central Australian Bentley Supergroup. It is unclear to what extent the effects of a ~550 Ma convergent event that resulted in significant crustal overthrusting and high-P recrystallisation in Musgrave Block rocks can be extrapolated to neighbouring terranes. Nonetheless, this event may be indicative of plate margin processes of latest Neoproterozoic age that controlled the assembly of central Australia'. Below the description are buttons for 'Add to Collection', 'Remove from Collection', and 'Remove Saved record'. At the bottom right of the page, there are links for 'Send Page to a Studio' and 'Send Collection to a Studio'.

Users can also launch Jupyter notebooks for collections. For example, a notebook implementing a landscape model and requiring multiple inputs that exist as separate datasets in DDStudio, can be launched from a collection that organizes such model inputs.

WATCH DATADISCOVERYSTUDIO.ORG IN ACTION!

[VIDEO] <https://www.youtube.com/embed/3opK1o8LgkI?feature=oembed&fs=1&modestbranding=1&rel=0&showinfo=0>

SGCI COLLABORATION

The Science Gateways Community Institute (SGCI) has served as an important partner to DDStudio over the last year.



Find the data you need

Data Discovery Studio, an NSF EarthCube Project, offers a large inventory of high-quality earth science resources with documentation that enables data discovery and re-use.

With over 1.4 million resources in our inventory, you'll have access to data from a wide variety of NSF-supported, government and international scientific repositories and catalogs at no cost to you.

 **Data Discovery Studio**
datadiscoverystudio.org

Key results of our collaboration with SGCI:

- Development of marketing and branding strategies, materials for distribution at professional meetings, and improved social media strategies
- User interface improvements, based on design advice from the SGCI usability team, who conducted user interviews, performed usability testing, and analyzed a dozen of other search portals to identify the most useful features.
- Cybersecurity assessment and single sign-on improvements.