# Enabling Cloud Services and Enhanced Data Discovery with earthdata-varinfo

Eni Awowale[1,2], Carlee Loeser[1,2], Owen Littlejohns[3,4], Jennifer Adams[1,2], David Auty[3,5], Lena Iredell[1,2], Brianna Rita Pagán[1,2], Mahabal Hegde[1], Nicholas Lenssen IV[1,2]

[1] NASA Goddard Earth Sciences Data and Information Services Center (NASA GES DISC), [2] ADNET Systems, Inc., [3] NASA Earth Observing System Data and Information System (EOSDIS) Evolution and Development (NASA EED-3), [4] INNOVIM, [5] Raytheon Company

## Introduction

NASA's Earth Observing System Data and Information System (EOSDIS) contains thousands of Earth science datasets from satellites, models, and field campaigns. Each of these collections can contain hundreds of variables that describe each measurement within the dataset, therefore an automated method for generating UMM-Var records is necessary. The Unified Metadata Model for Variables (UMM-Var) provides a framework for variable metadata records in NASA's Common Metadata Repository (CMR). The Python tool, earthdata-varinfo, was developed to solve this problem of automating the curation of UMM-Var records. Given either a collection DMR file or a netCDF-4 file, earthdata-varinfo can scrape variable metadata and return a CMR compliant UMM-Var record. Earthdata-varinfo can generate thousands of UMM-Var records in a matter of seconds, thus enabling subsetting capabilities and enhancing data discovery.

## Motivation

1) Many scientists search for data based on variables or measurements, rather than collections or projects. Curating metadata for these variables is important for discovery.
2) Subsetting, regridding, and reformatting are important data transformation tools that allow scientists to reduce the size of data files and transform them for comparison to other data. To enable subsetting by variable, curating variable metadata is necessary.
3) Since some data collections have hundreds of variables in a single file, it is challenging to curate this metadata manually. Earthdata-varinfo extracts the CF attributes for all of the variables in a data file and outputs them into CMR-ready metadata records.
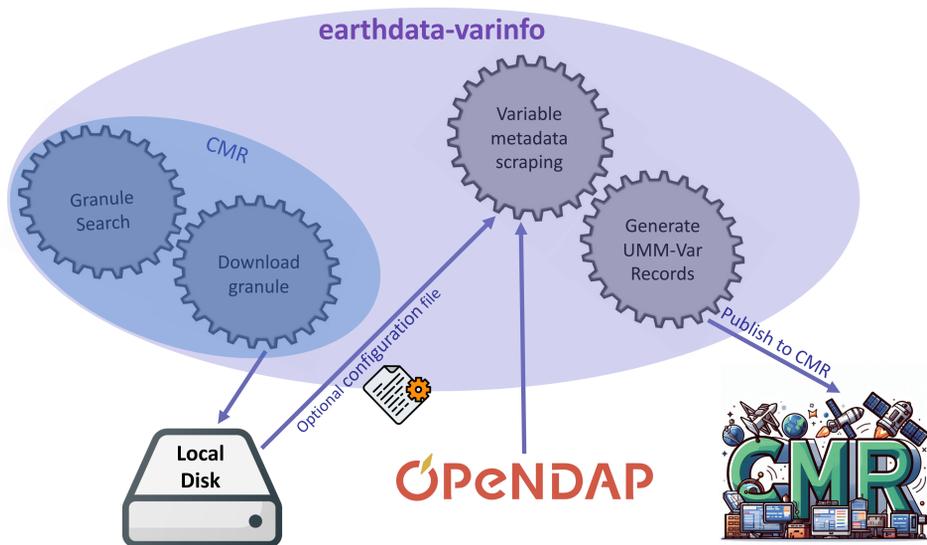
## Workflow and Capabilities



**Figure 1.** Workflow diagram for earthdata-varinfo

**How does it work?**

1) Earthdata-varinfo requires a single granule that is representative of the data collection and includes all its variables. There are three options:
   a) Conduct a CMR granule search to download the most recent granule from Earthdata Search locally.
   b) Choose the representative granule through OPeNDAP.
   c) Use a granule that is already stored locally.
2) Earthdata-varinfo can be run locally or in a Docker container. Collection level configuration files can be included to the workflow for known quirky datasets. Earthdata-varinfo returns a UMM-Var record for each variable in the granule that is CMR compliant.
3) UMM-Var records are published to CMR.

Alternatively, UMM-Var records can be generated via the CMR GraphQL API, which utilizes earthdata-varinfo.

**What CF attributes are collected with earthdata-varinfo?**

- long_name
- standard_name
- title (or definition, description)
- dtype (e.g. float32, int16)
- dimensions
- units
- _FillValue
- scale_factor
- add_offset
- valid_range



**Figure 2.** A global map of the Earth showing the locations of selected ground-based observation stations for OMI.

**What are the compatible granule file formats?**

- '.nc', '.nc4', '.grb', '.h5', '.he5', '.HDF5'
- OPeNDAP file types (e.g. '.dmr')

Find us and try it out on GitHub!
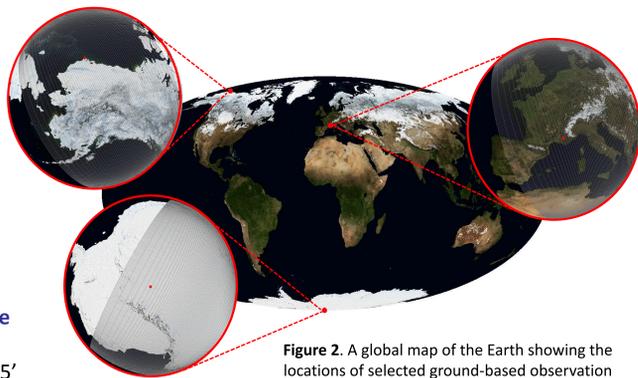
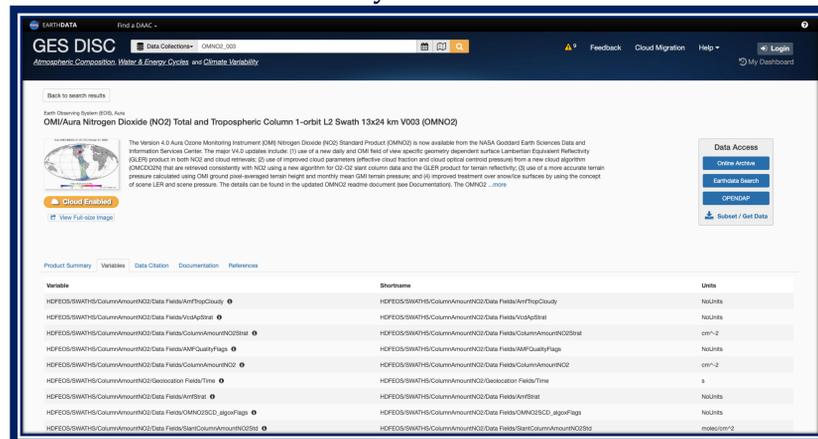## Data Discovery and Transformation



**Figure 3.** Dataset landing page for the *OMNO2_003* collection. Note the variables tab.

Once the UMM-Var records have been published, they are displayed on the GES DISC dataset landing pages under a "Variables" tab. Users can assess what variables are in the data before downloading it. Additionally, the short names are displayed for a more direct interpretation of the variables after downloading. In the future, this page will be expanded to provide more information for each variable.
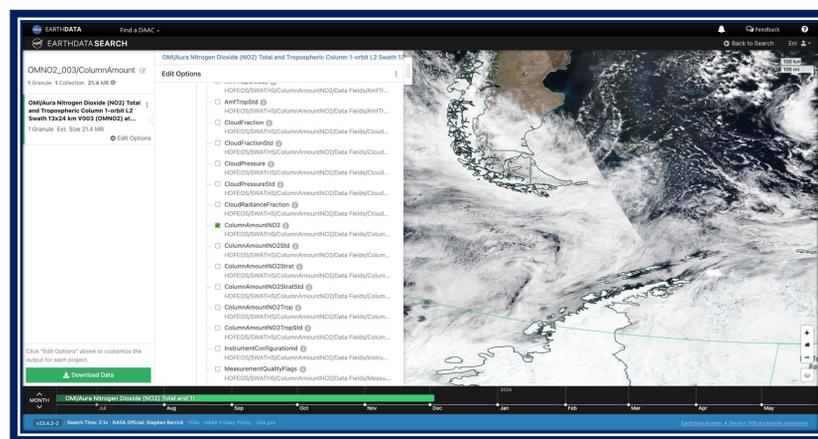


**Figure 4.** A variable subsetting request for ColumnAmountNO2 in Earthdata Search with the Harmony Level 2 subsetter.

Once the UMM-Var records have been published, it is possible to perform a variable subset using Harmony services within Earthdata Search. By selecting only the variables that the user needs, the subsequent data downloaded is smaller, more manageable, and often yields a faster download. The Earthdata Harmony service can be accessed through the NASA Earthdata Search website. Look for the customization icons to determine which datasets have these services available.
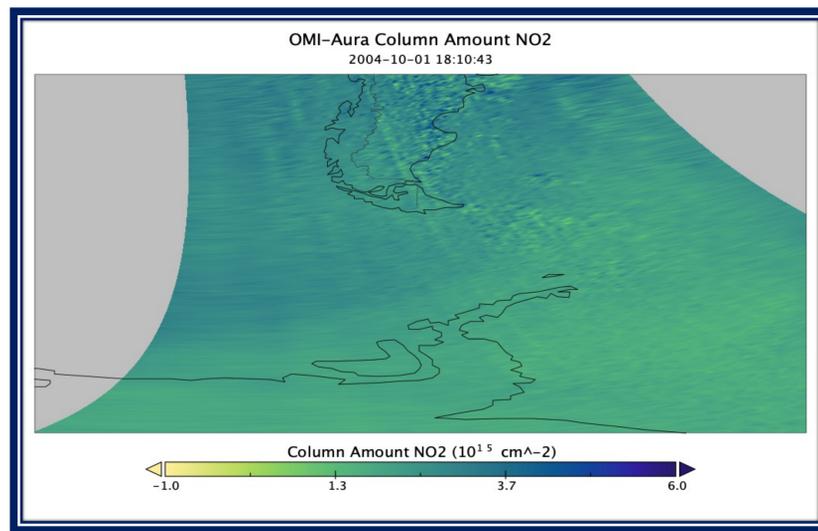


**Figure 5.** A spatial and variable subset for the variable ColumnAmountNO2 in the collection *OMNO2_003*.

Once the user selects the variable(s) and downloads the data, they can proceed with their data analysis as usual. This is one example of a Harmony subsetting output file, for the ColumnAmountNO2 variable from the collection *OMNO2_003*. This variable was also spatially subsetted to focus on southern Patagonia and the Antarctic peninsula region. The original file was ~ 20MB, and the output file after the variable and spatial subsetting was ~ 0.15MB.

## What's next?

To date, we have published 7000+ variables from ~58 GES DISC collections. We will continue to use earthdata-varinfo to publish more variables while enhancing and expanding its functionality and usability.

Future plans:

- Improvements for identifying the classification of variables into: science variables, dimension variables, quality variables, or ancillary variables, which is important for both human- and machine-readability.
- Mapping the variable to a GCMD science keyword to enable data discovery through the GES DISC website and Earthdata Search.
- Adding Launchpad token authentication for publishing to CMR as ingest via Earthdata Login (EDL) Bearer Tokens are being disabled.