

Big data in Earth science: Using Pangeo to work with OOI data in the Cloud

Tim Crone

Rich Signel, Dax Soule, Filipe Pires
Alvarenga Fernandes

15 December 2021

Lamont-Doherty Earth Observatory
COLUMBIA UNIVERSITY | EARTH INSTITUTE



Google Cloud Platform



Big data in the geosciences



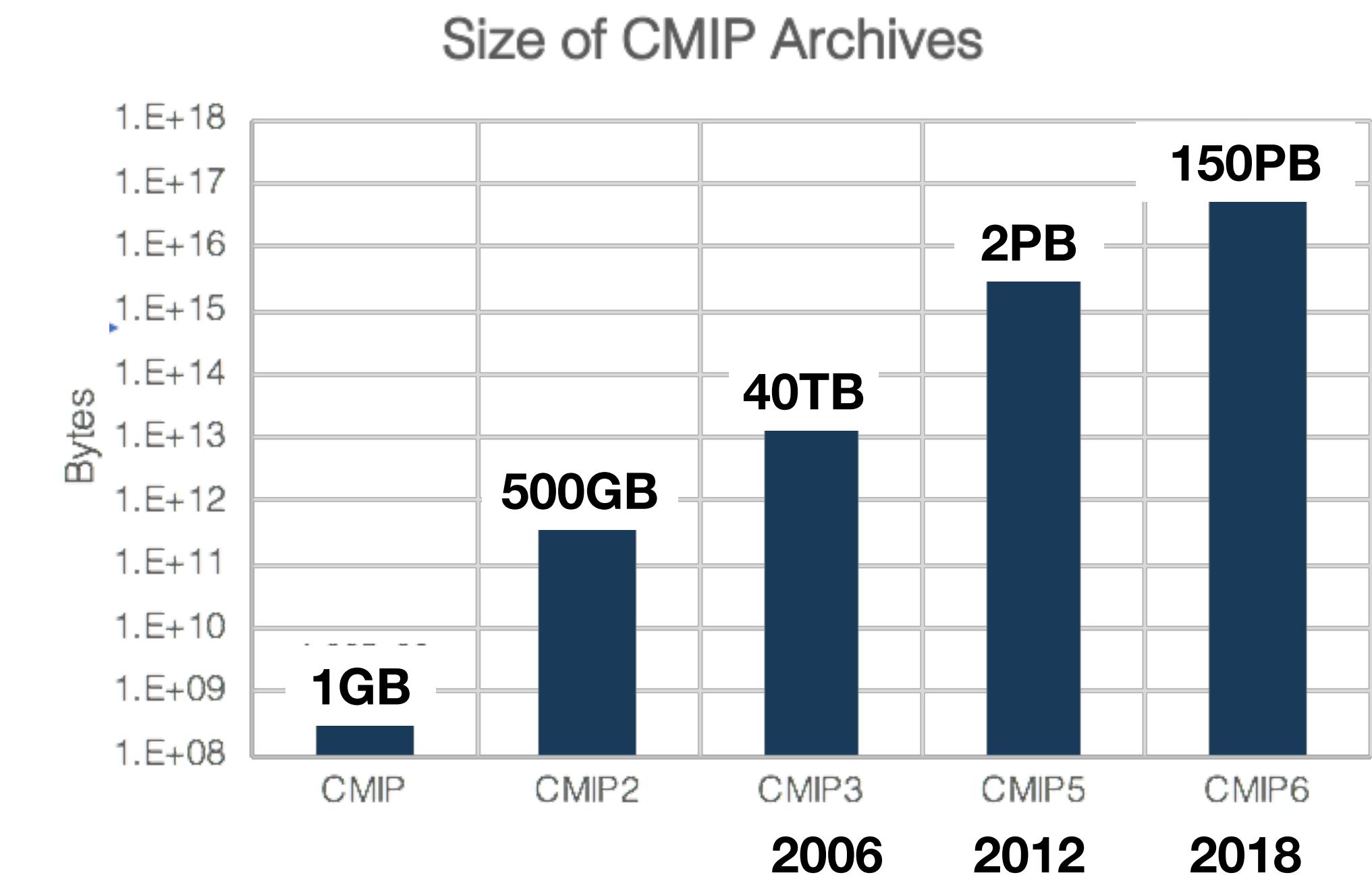
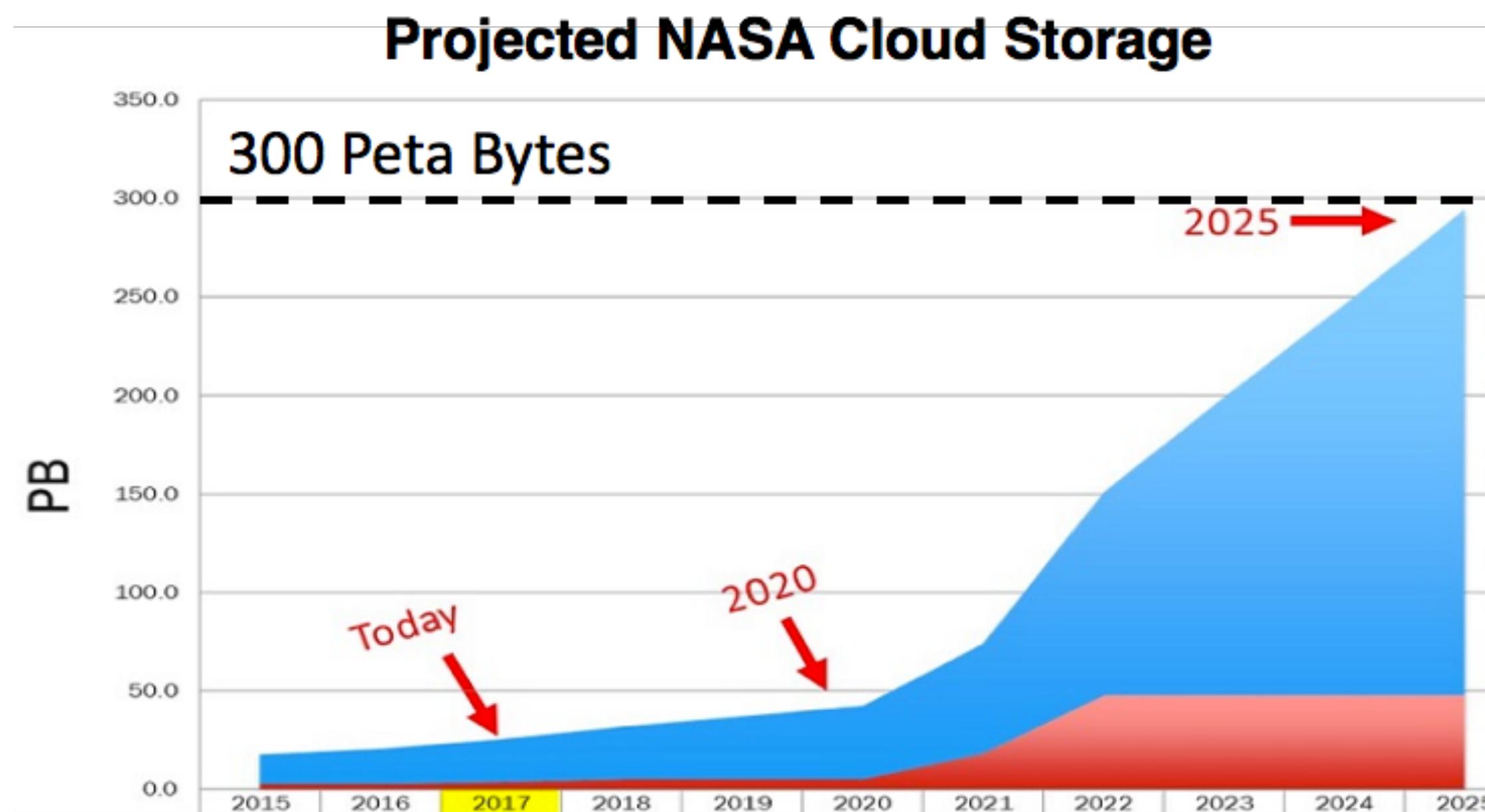
Increase in Earth Observations

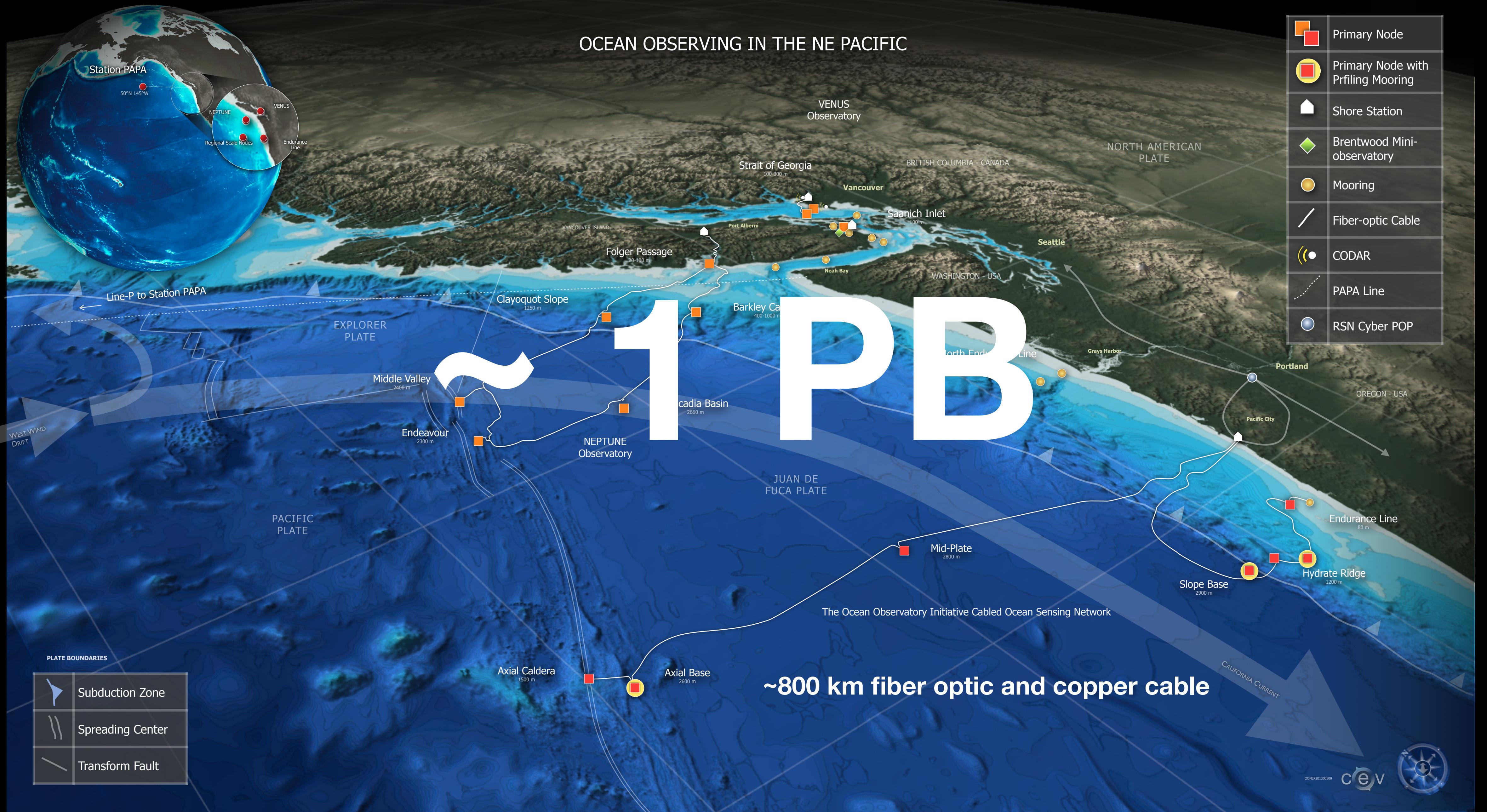
- New sensors / platforms
- Continuous observations
- Multiple versions of derived datasets

Increase in Model Data

- Higher resolution
- More process representation
- Larger ensembles

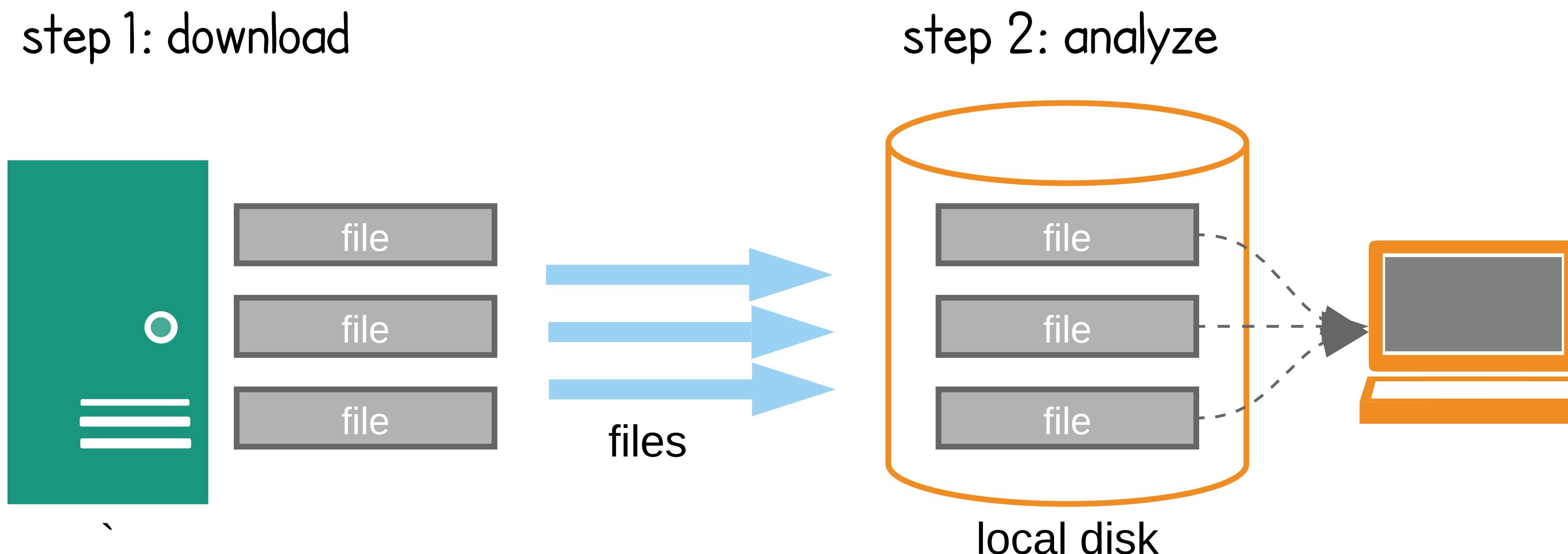
<https://earthdata.nasa.gov/about/eosdis-cloud-evolution>





*How should scientific data analysis
infrastructure be organized for the
petabyte era?*

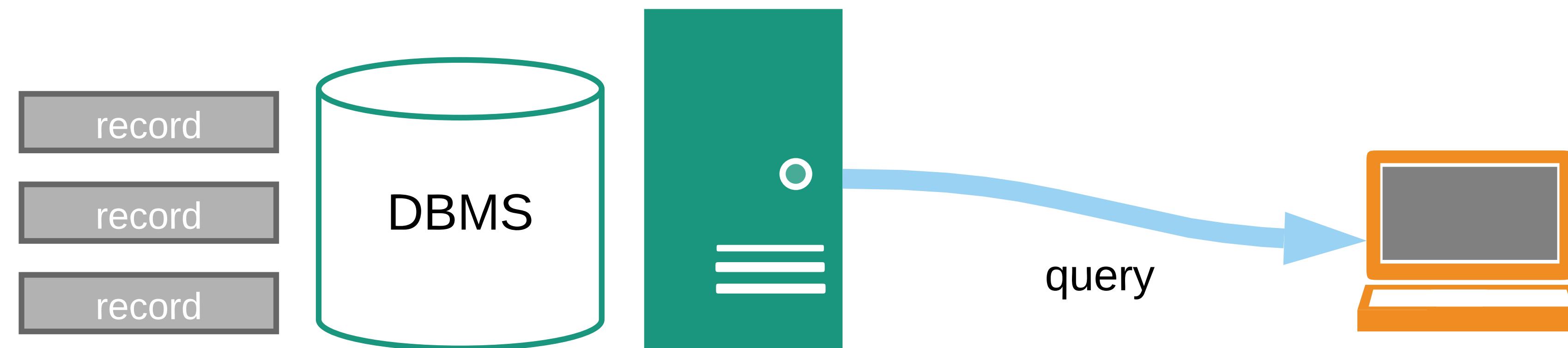
File-based Approach



Data provider's responsibilities

End user's responsibilities

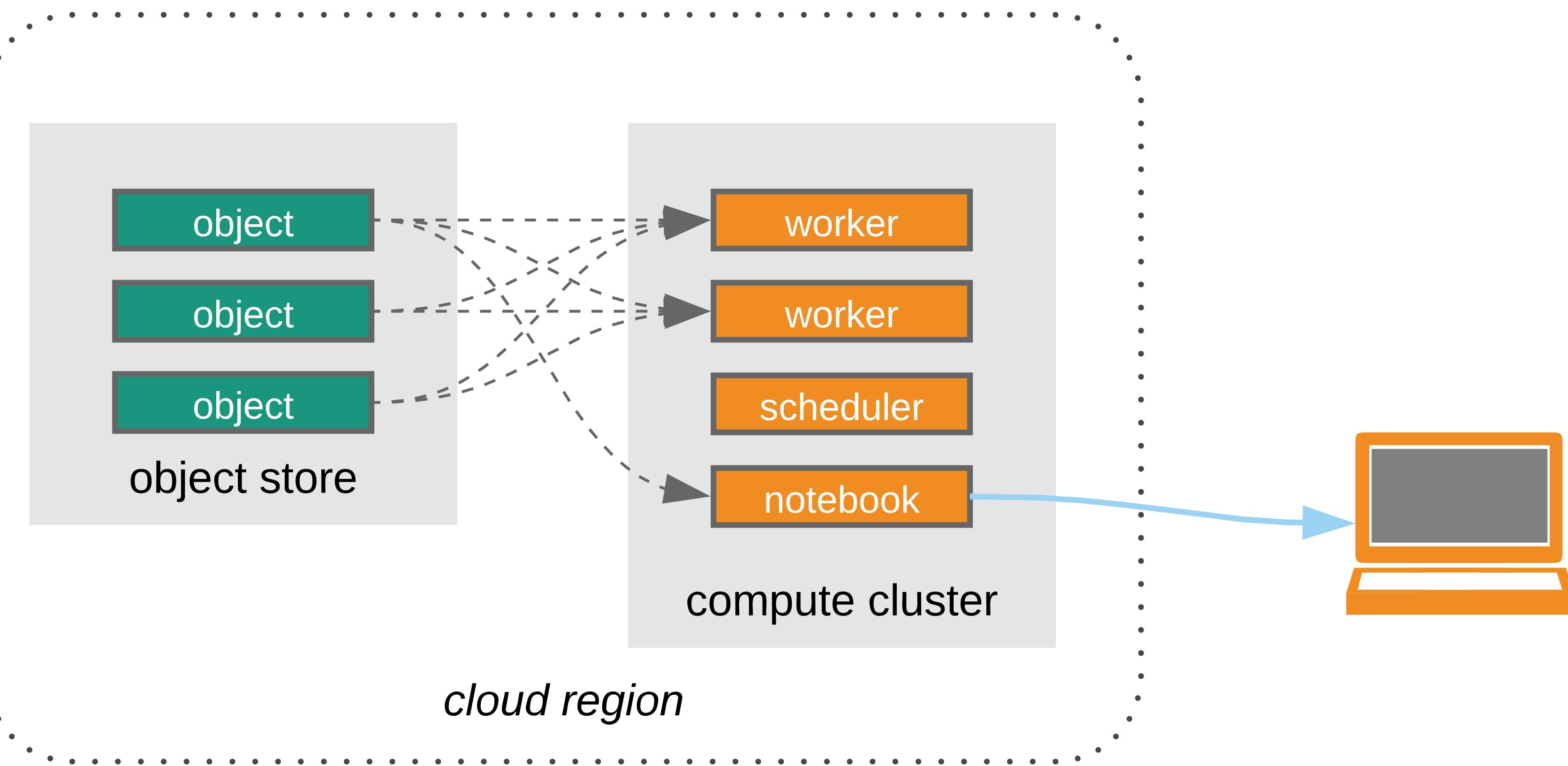
Server-Side Database



Data provider's responsibilities

End user's responsibilities

Cloud-Native Approach



Data provider's responsibilities

End user's responsibilities

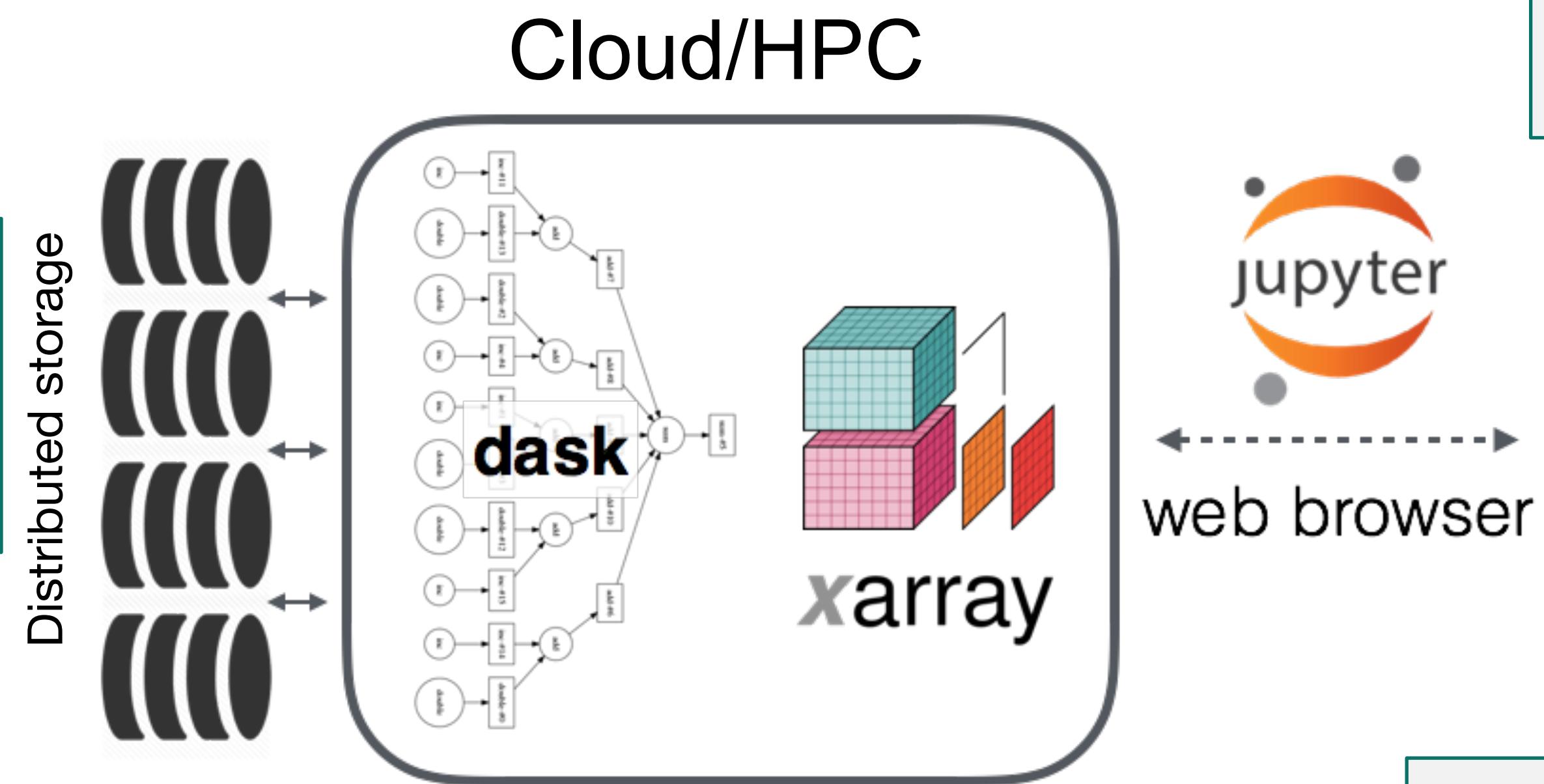


Pangeo

A *community-driven* effort for
Big Data geoscience

The Pangeo Architecture

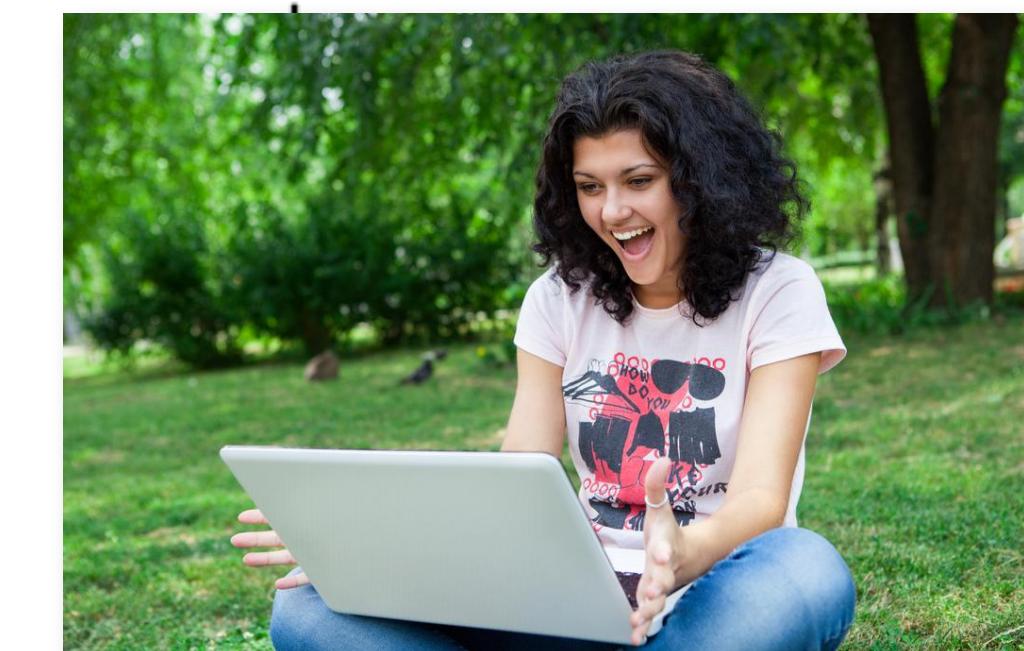
Analysis Ready Data
Stored and cataloged on
globally-available distributed
storage (**e.g. S3, GCS**)



Parallel computing system built on
top of **Kubernetes or HPC**.

Dask tells the nodes what to do.

Jupyter for interactive access
on remote systems

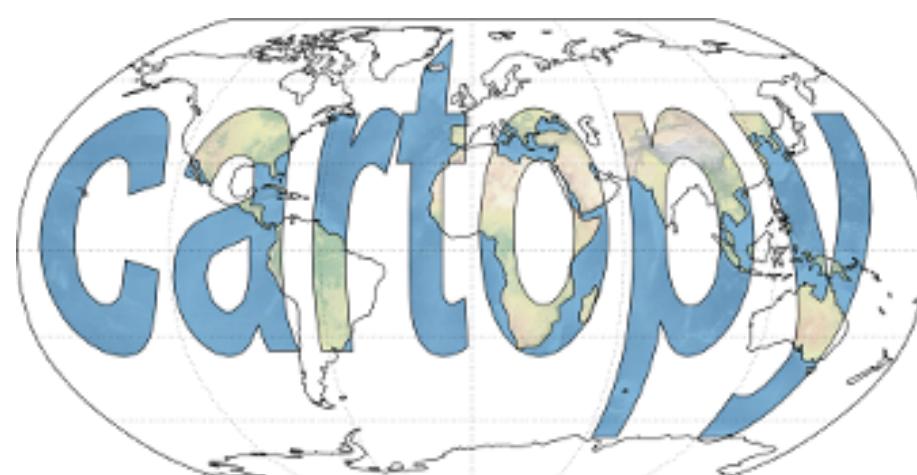


Xarray provides data structures and
intuitive interface for interacting with
datasets

Scientific Python for Geoscience



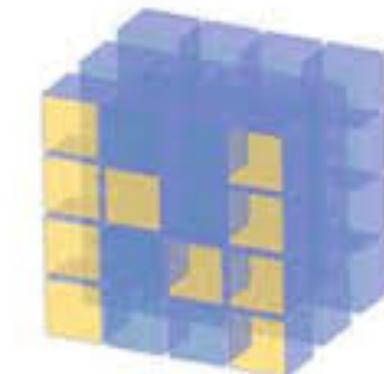
Iris



matplotlib



SciPy

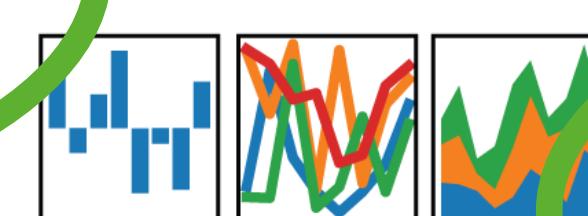


NumPy



python™

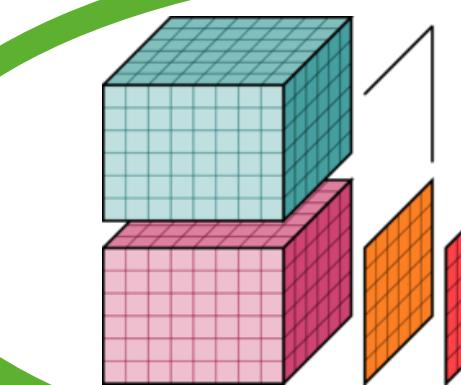
pandas
 $y_t = \beta' x_{it} + \mu_i + \epsilon_{it}$



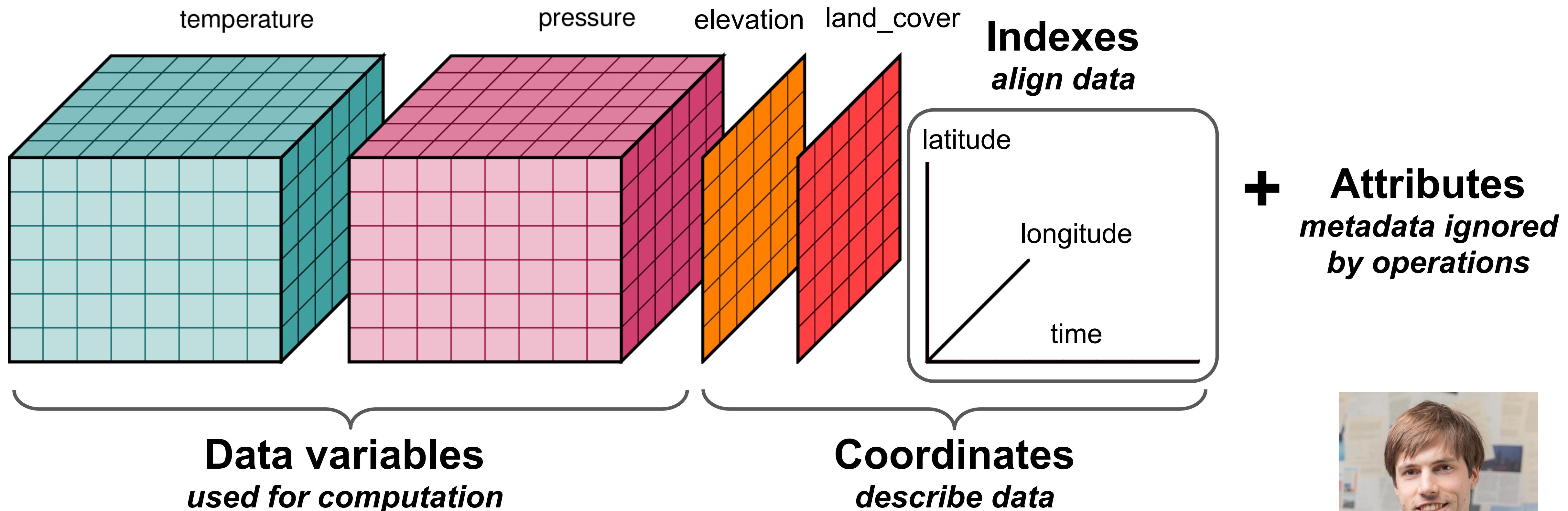
IP[y]:
IPython



xarray



Xarray: Multidimensional variables with coordinates and metadata



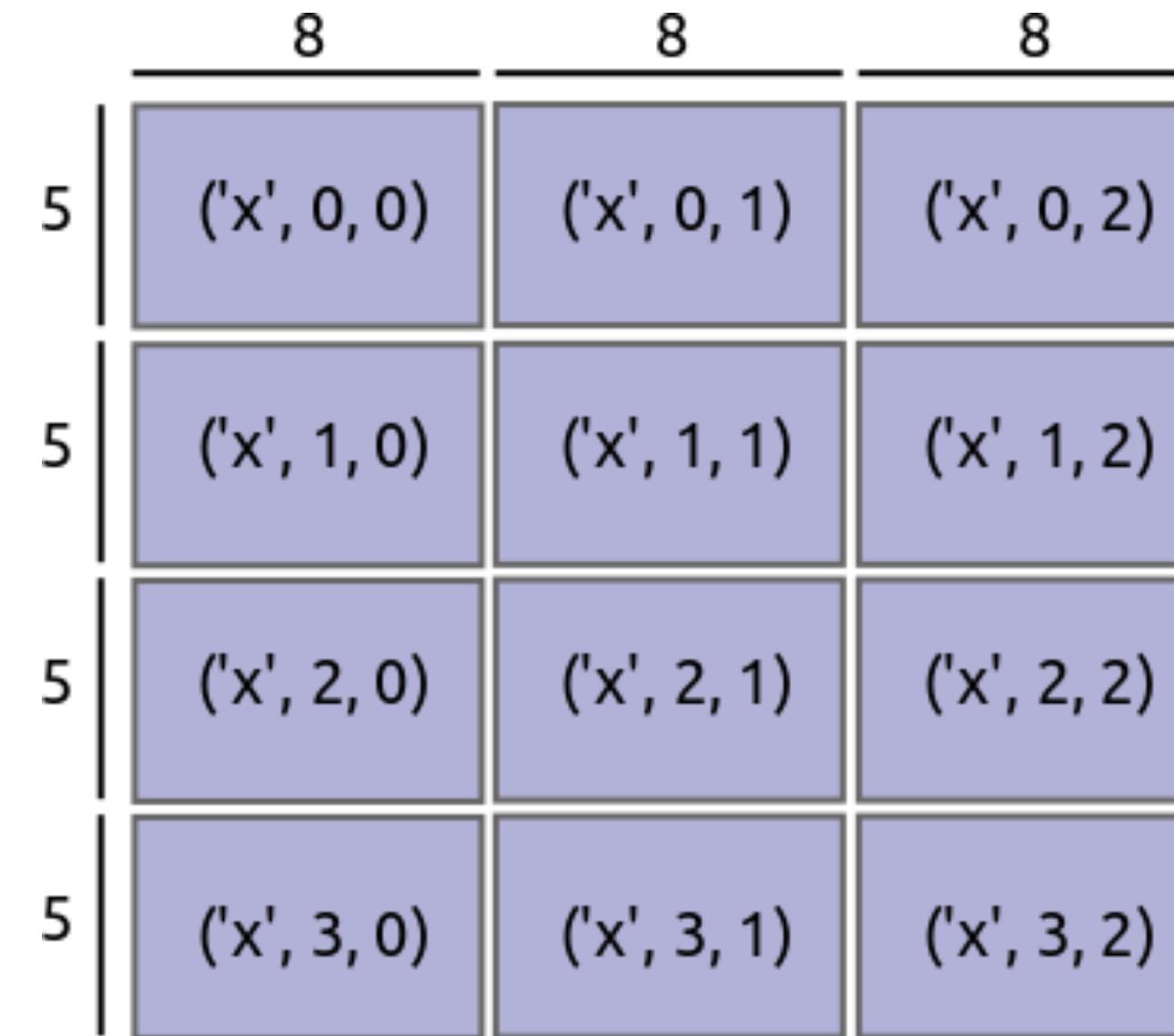
“netCDF meets pandas.DataFrame” plus! Dask integration



Credit: Stephan Hoyer

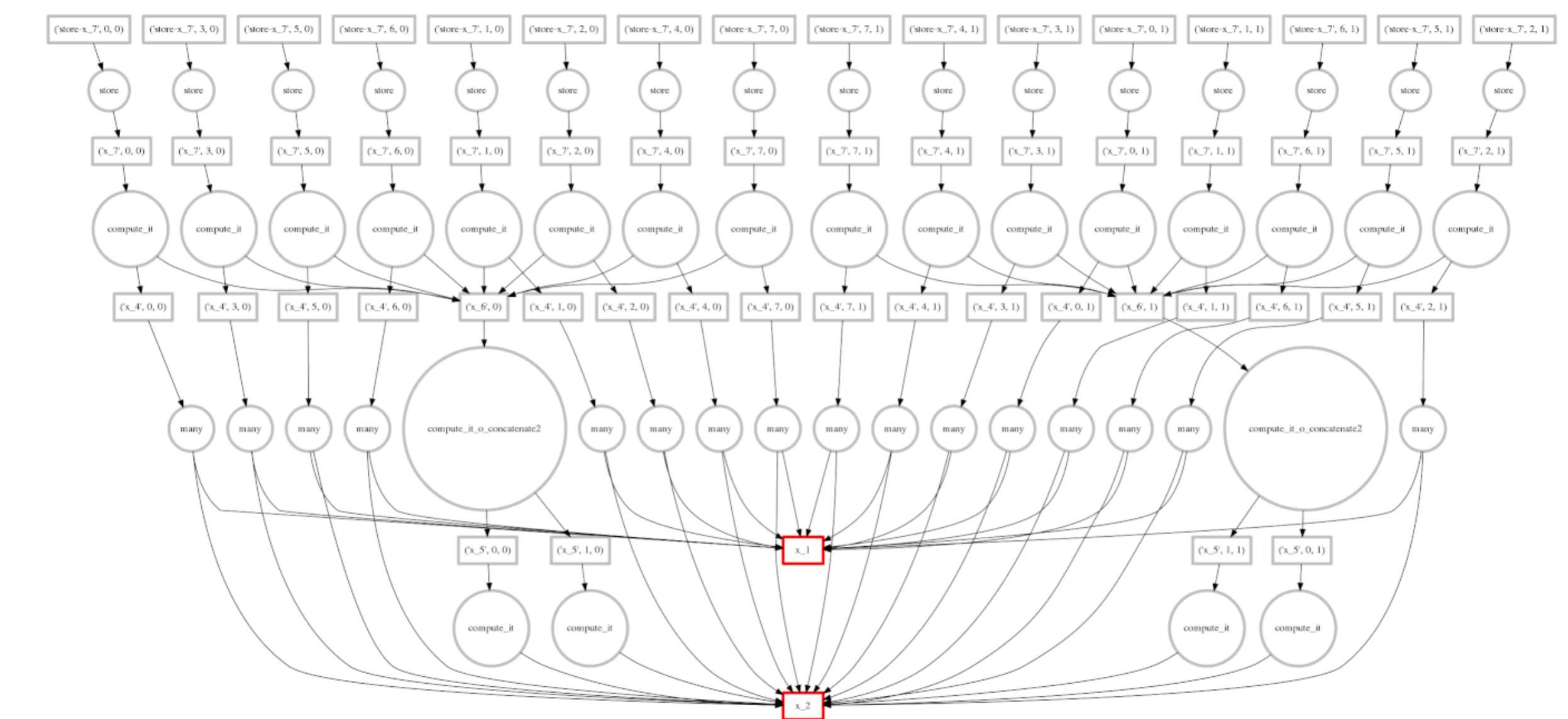
Dask: Simplified parallel and distributed processing in Python

<https://github.com/dask/dask/>



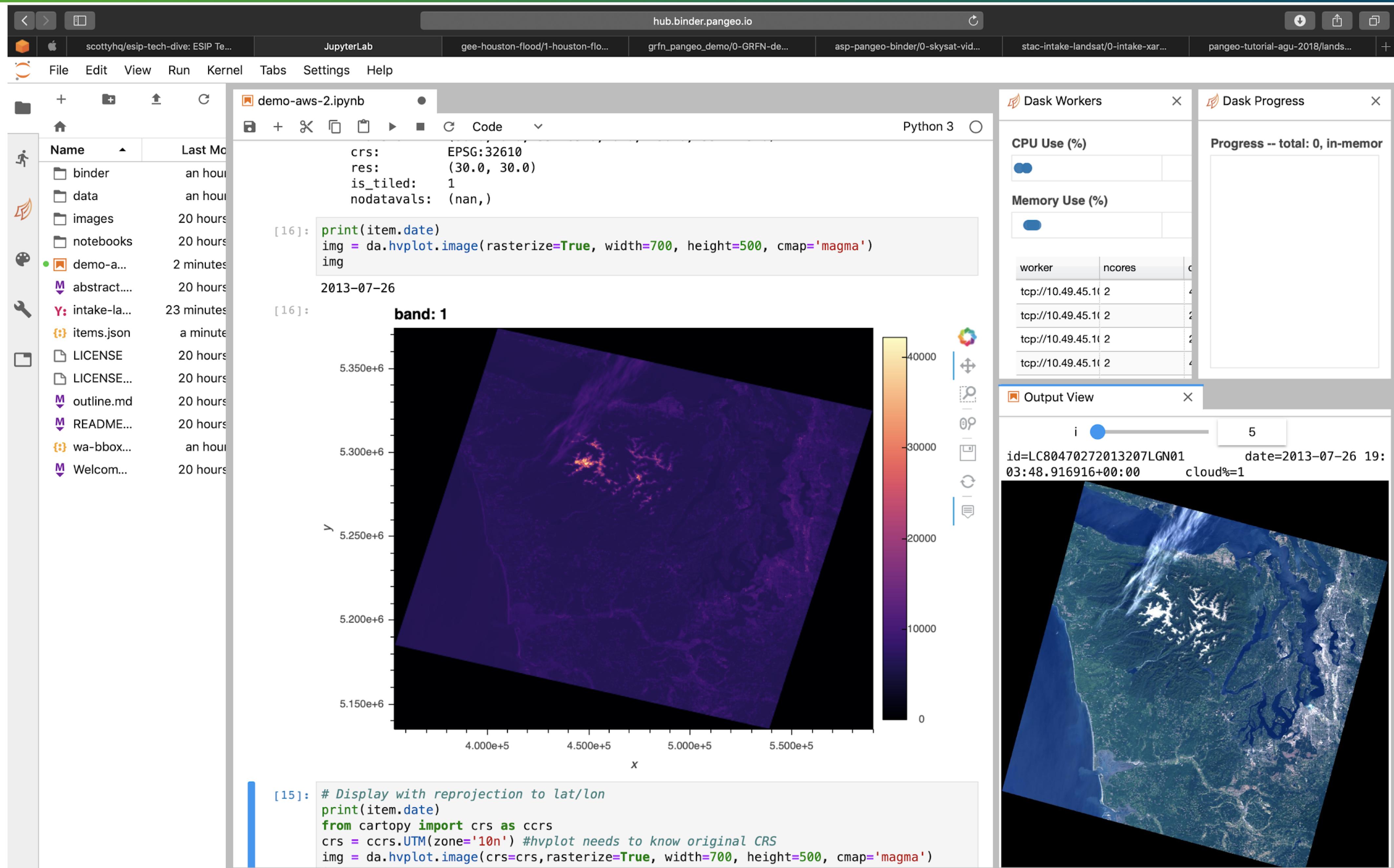
ND-Arrays are split into chunks that comfortably fit in memory

Dask.delayed allows the “lazy” loading of chunks which are not pulled into memory until required

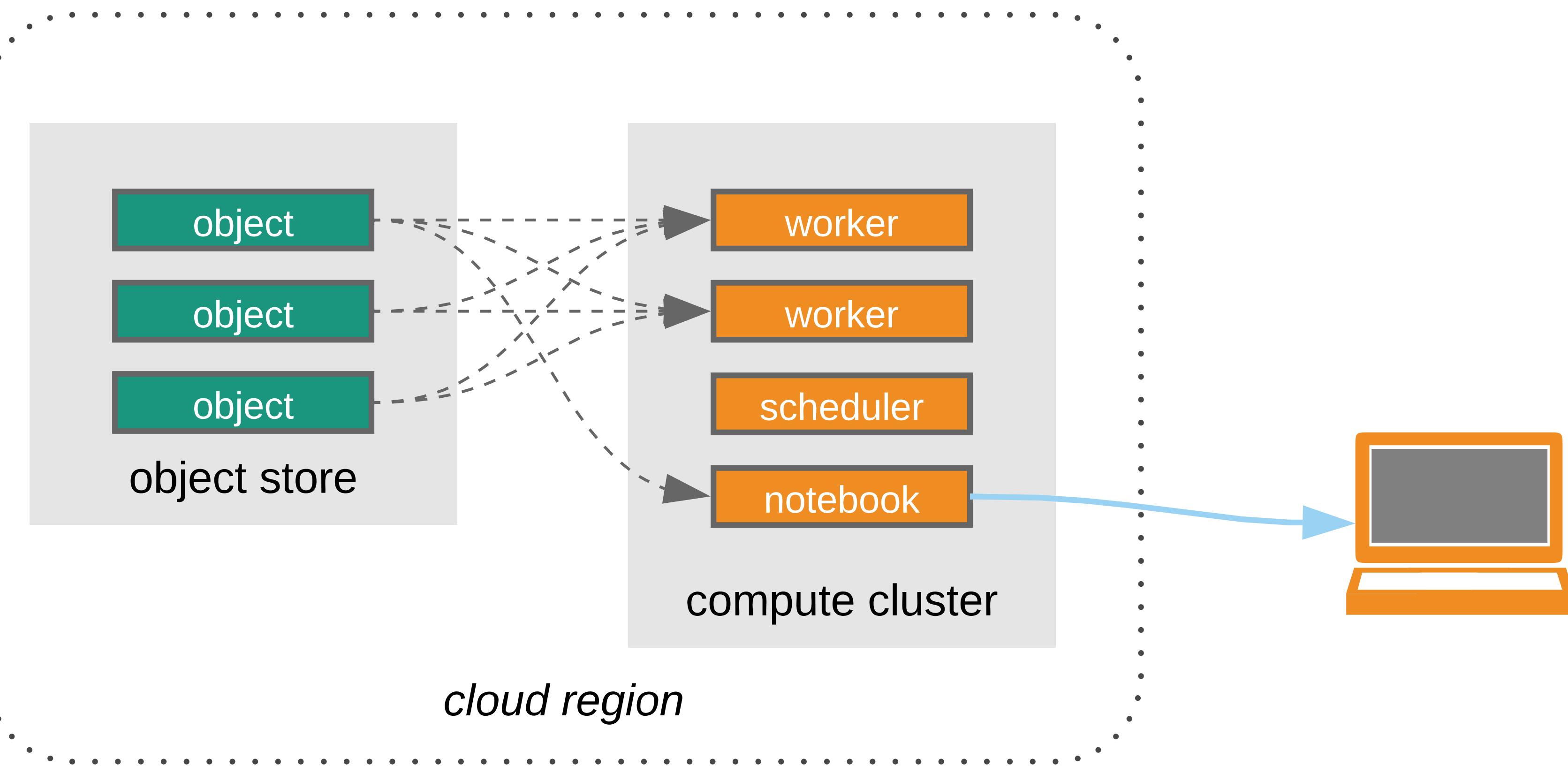


Complex computations represented as a graph of individual tasks.

Scheduler optimizes execution of graph.



Cloud-Native Approach



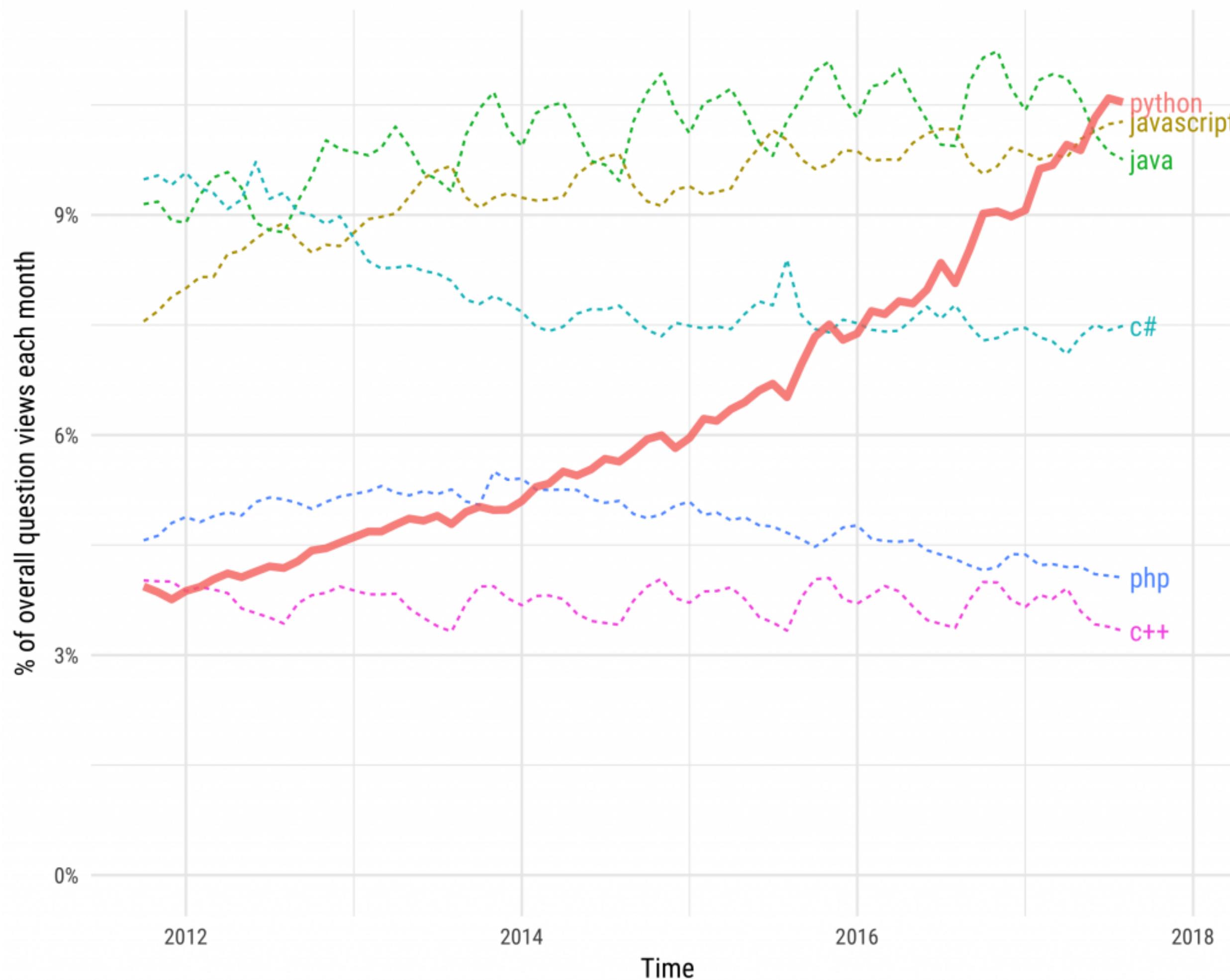
Data provider's responsibilities

End user's responsibilities

Scientific Python for Data Science

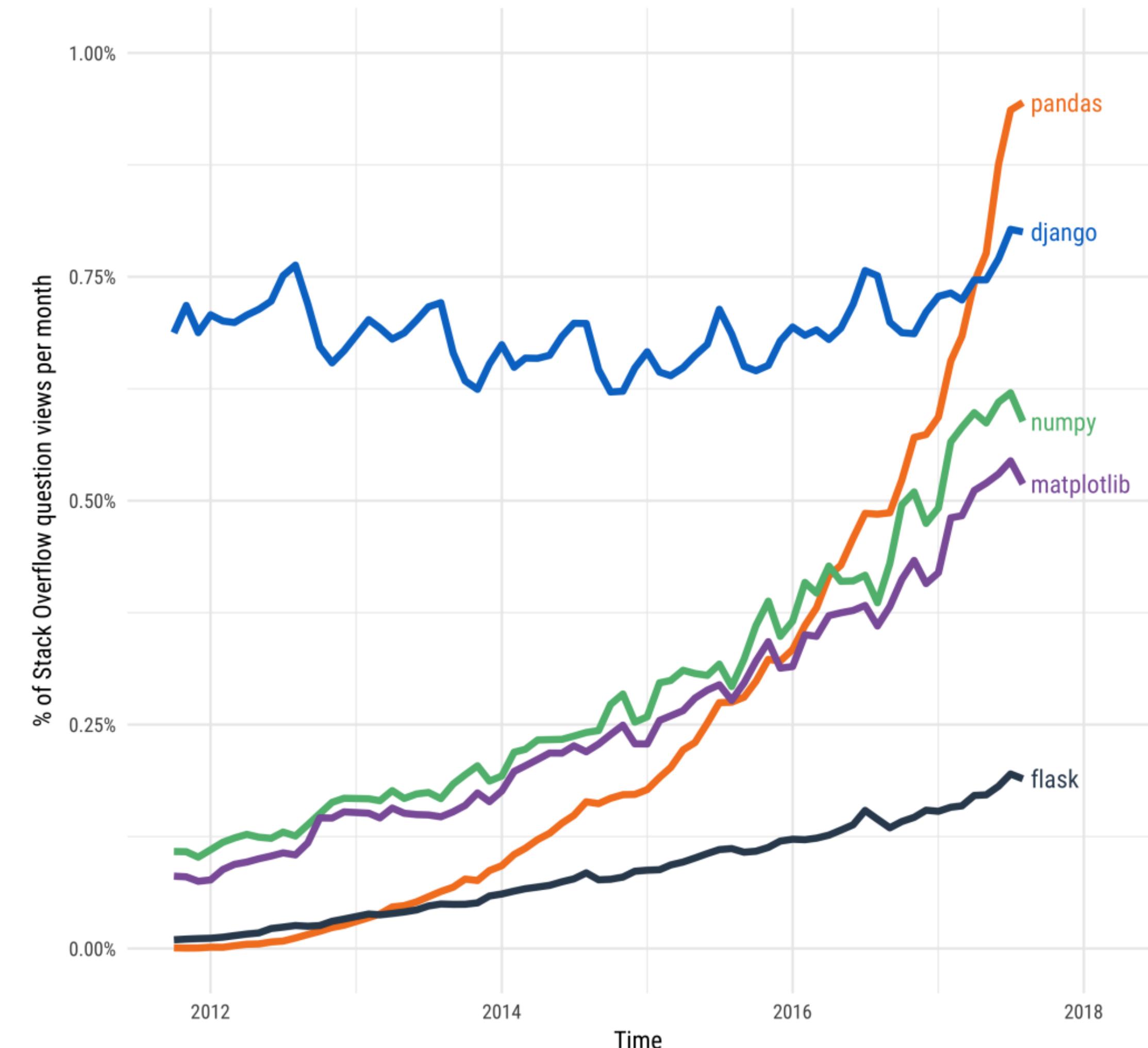
Growth of major programming languages

Based on Stack Overflow question views in World Bank high-income countries



Stack Overflow Traffic to Questions About Selected Python Packages

Based on visits to Stack Overflow questions from World Bank high-income countries



all categories ▶

Latest

Top

Categories

Topic	Replies	Views	Activity
<p>TOPIC</p> <p>Welcome to Pangeo Discourse</p> <p>Meta</p> <p>Pangeo is a community of scientists and software developers working together to improve the way we do scientific research. We work on software tools, such as the python packages Xarray and Dask, as well as customized env... read more</p>	5	215	Oct 18
<p>Call for speakers @ ESIP Pangeo session</p> <p>Science</p>	4	34	20h
<p>Access to Pangeo GCS Bucket to push model output from pre-CMIP6 experiments?</p> <p>Cloud</p>	6	19	21h
<p>Interpolation and regridding</p> <p>Science</p>	7	169	1d
<p>Pangeo Weekly Telecon</p> <p>News & Announcements</p>	17	173	2d

Pangeo Funding and Contributors



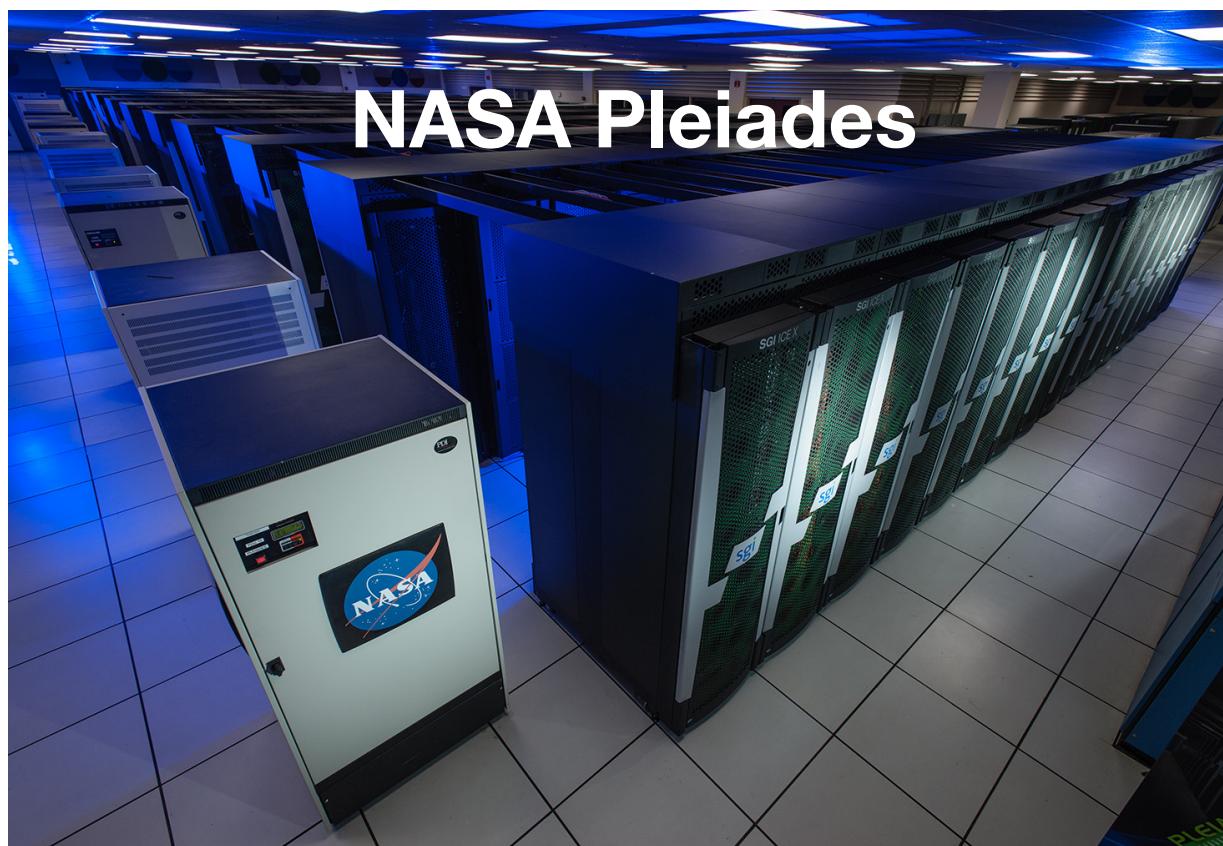
Google Cloud Platform



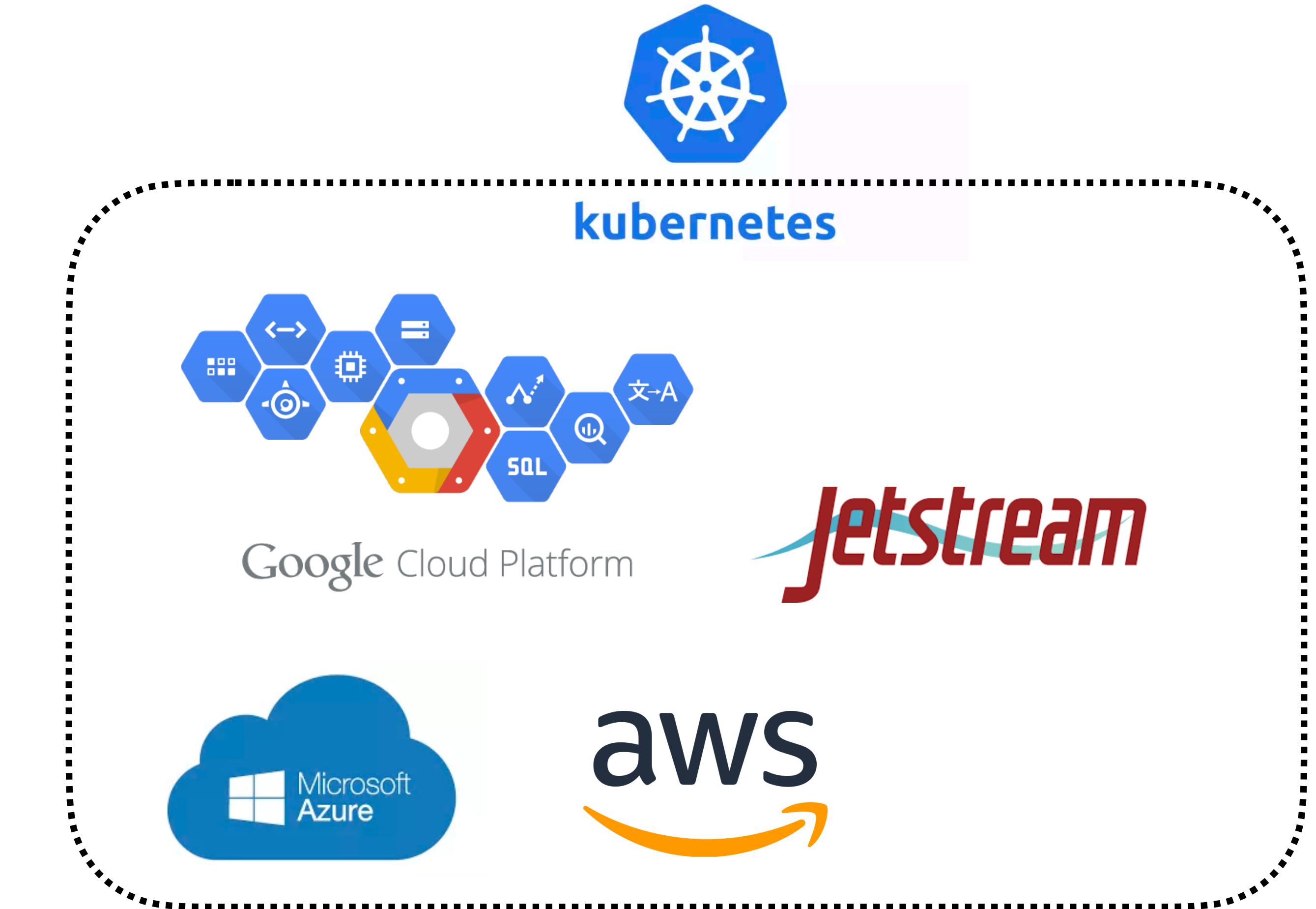
Pangeo Principles for Cloud-Native Science Infrastructure

- **Community-driven** - Our needs are no different from those of our peer institutions. By developing infrastructure collaboratively, we can accomplish much more than any one institution can alone.
- **Open source** - Because infrastructure is code, the code should be licensed in a way that enables the entire research community to reuse and build upon it.
- **Modular** - “all in one” solutions are impossible to maintain long term. Separation of concerns is a key principle of good software and systems engineering.
- **Vendor neutral** - Academic research infrastructure should use only vendor-neutral services APIs. If this principle is followed, it means we can redeploy our infrastructure anywhere.

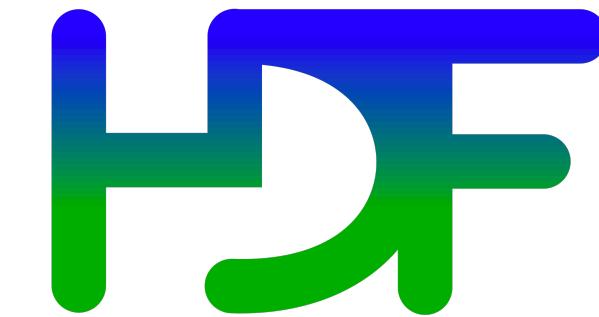
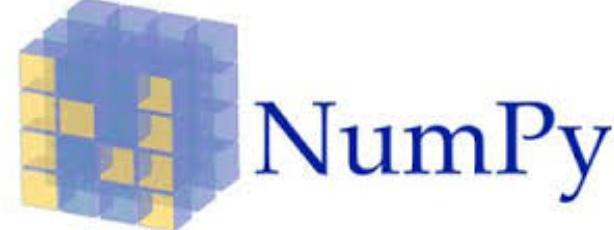
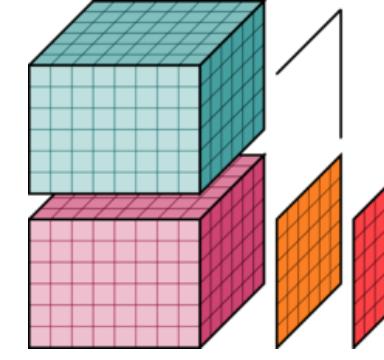
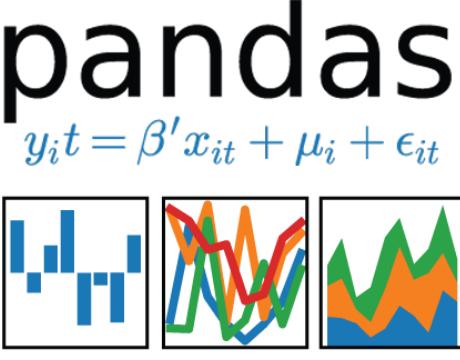
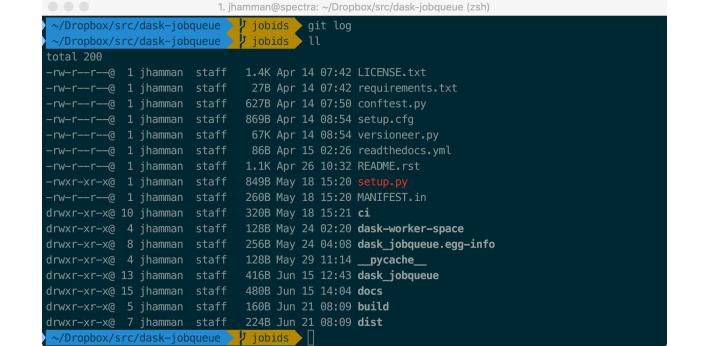
Pangeo Deployments



NCAR Cheyenne



Build your own pangeo

Storage Formats			Cloud Optimized Zarr/TileDB/Parquet/etc.
ND-Arrays			More coming...
Data Models			
Processing Mode	 Interactive	Batch 	 Serverless
Compute Platform	HPC 	Cloud  Google Cloud Platform	Local 

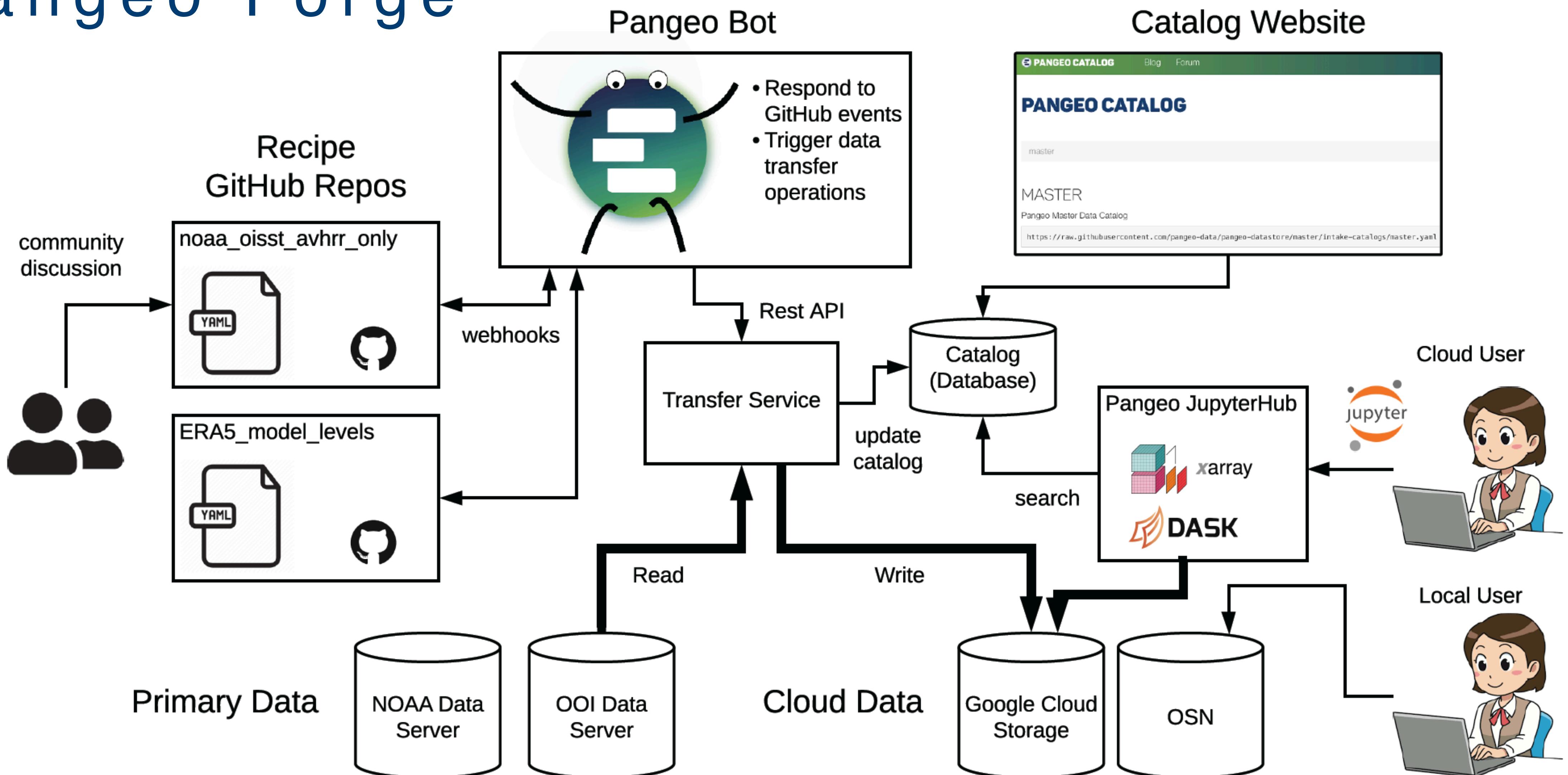
Toil

Adapted from the Google Site Reliability Engineering Guide

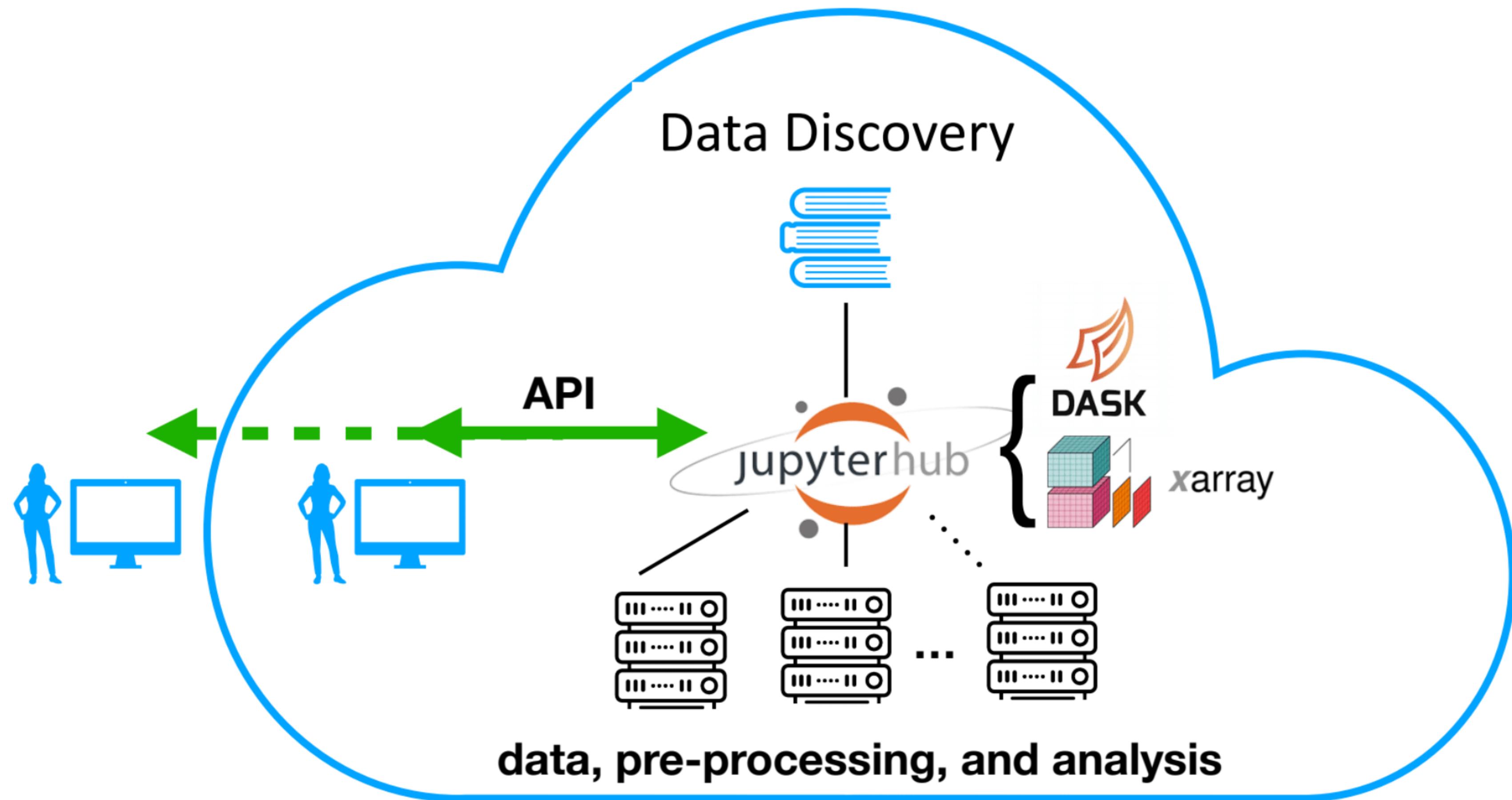
"Toil is the kind of work tied to scientific research that tends to be manual, repetitive, automatable, devoid of enduring value, and that scales linearly with the number of papers published."

B.Beyer, C. Jones, J. Petoff, and N. Murphy (2016)
Site Reliability Engineering: How Google Runs
Production Systems, O'Reilly Media, Incorporated.

Pangeo Forge



The Future is Now



Robinson et al., 2019

Some Things about the cloud that we have come to understand in the last couple of years

1. Compute is actually really inexpensive in the cloud
2. Storage is expensive, if you have to pay retail, but its possibly worth it
3. There is magic in the data centers. The providers typically don't say much about how their data centers are organized, but the data in the cloud can move really fast around the data centers.

Learn More and Get Involved!

- Presentation: [ESIP Tech Dive](#) (01/2018) by Ryan Abernathy & Matthew Rocklin
- Website: <http://pangeo.io/>
- Github: <https://github.com/pangeo-data>
- Blog: <https://medium.com/pangeo>
- Gitter: <https://gitter.im/pangeo-data>
- Discourse: <https://discourse.pangeo.io/>

The Ocean Observatories Initiative (OOI) and OOI Data

<https://oceanobservatories.org>

OOI Organizations



W
UNIVERSITY *of*
WASHINGTON

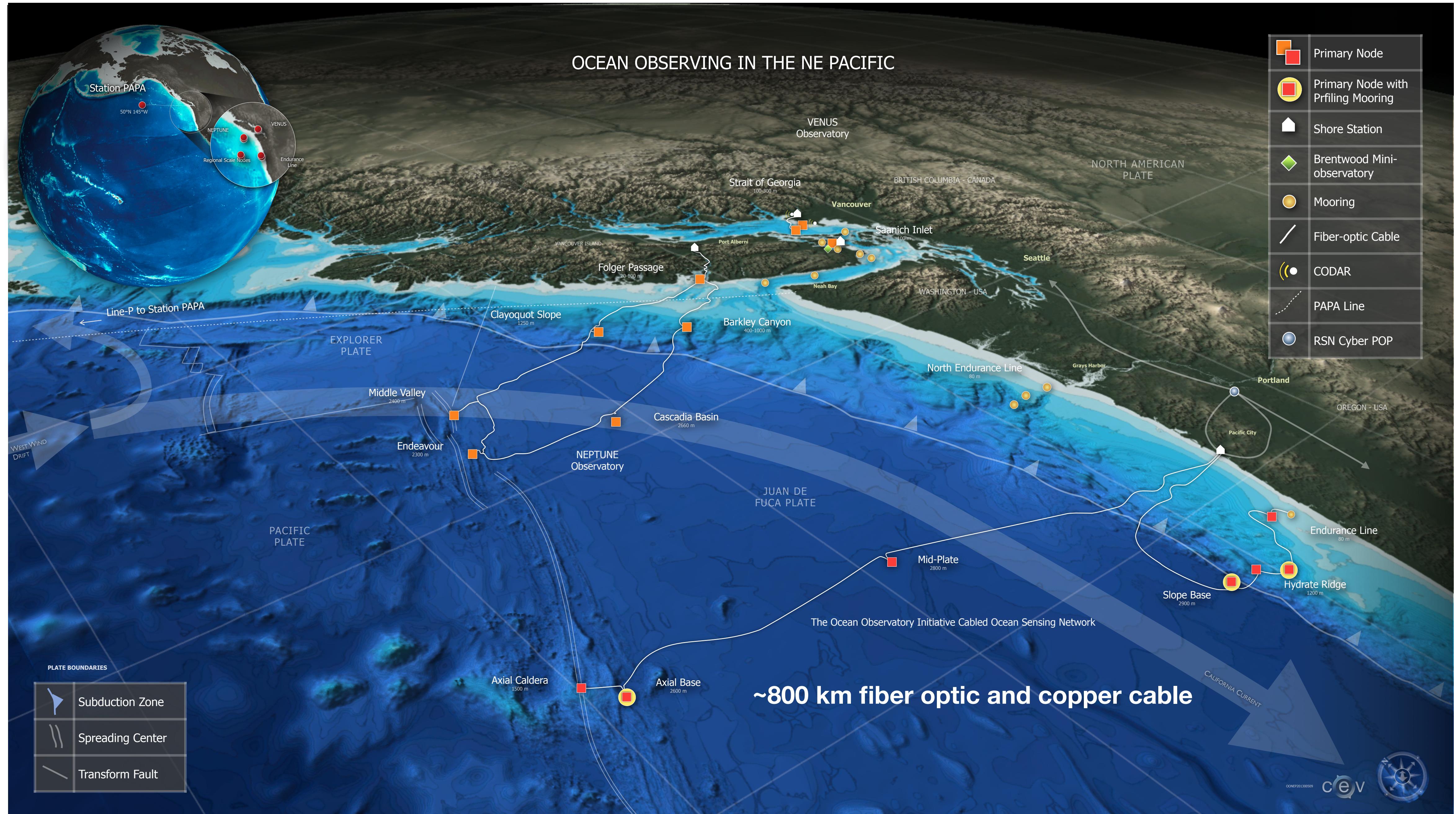


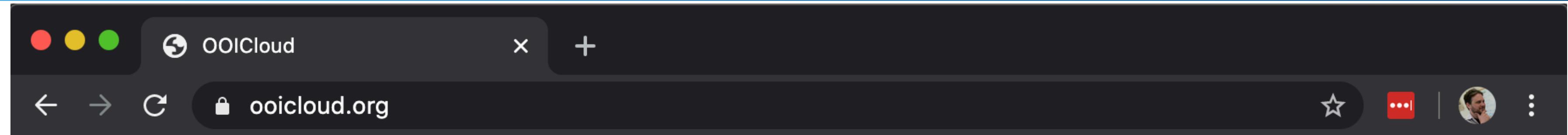
Lamont-Doherty Earth Observatory
COLUMBIA UNIVERSITY | EARTH INSTITUTE



Oregon State
University

RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY





About

The [OOICloud Project](#) is working to make data from the [Ocean Observatories Initiative \(OOI\)](#) publically available in the cloud and accessible through a [Pangeo](#) interface. A primary goal is to provide these data to the scientific community using a cloud-performant object storage model, and to provide large-scale remote compute capabilities for research investigations.

With a generous gift from the [Microsoft AI for Earth](#) program, OOICloud resides in the [Azure Cloud](#), and currently contains all of the data from the [OOI HD video camera](#) deployed at a hydrothermal vent in the caldera of Axial Volcano, a submarine volcano on the Juan de Fuca Ridge.

OOICloud will begin incorporating other OOI datasets in the coming year, starting with data from the geodetic instruments and the echosounders.

OOICloud  github.com/ooicloud    

Repositories 8 Packages People 19 Teams Projects Settings



OOICloud

Tools for working with OOI data in the cloud using Pangeo

Find a repository... Type: All Language: All Customize pins 

[ooicloud.github.io](#)
OOICloud website
HTML MIT 0 stars 0 forks Updated 30 minutes ago



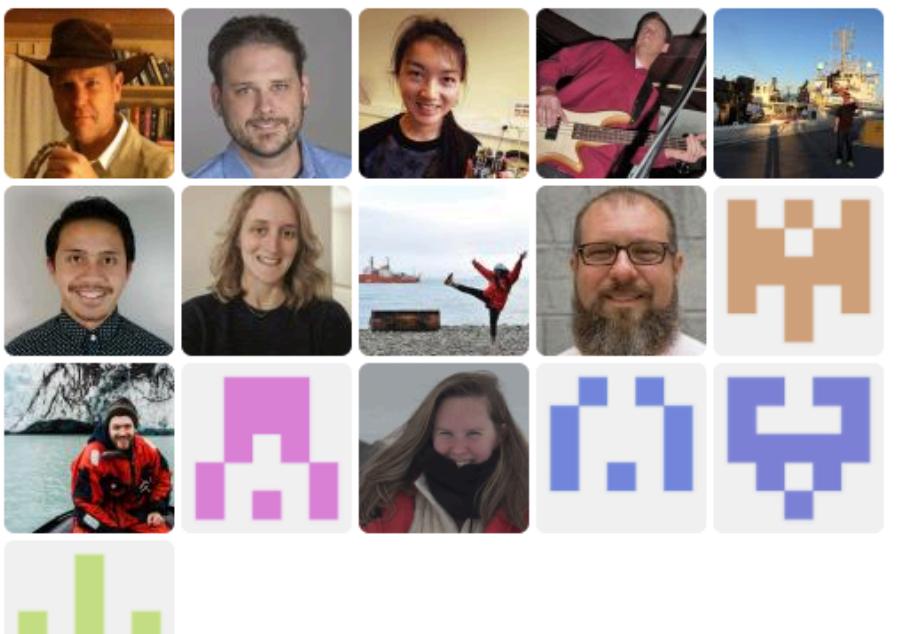
Top languages

- Jupyter Notebook
- Shell
- Python
- HTML
- Dockerfile

[axial-drilling](#)
Notebooks and data related to the Axial Drilling project
Jupyter Notebook MIT 2 stars 0 forks Updated 8 days ago



People 19 >



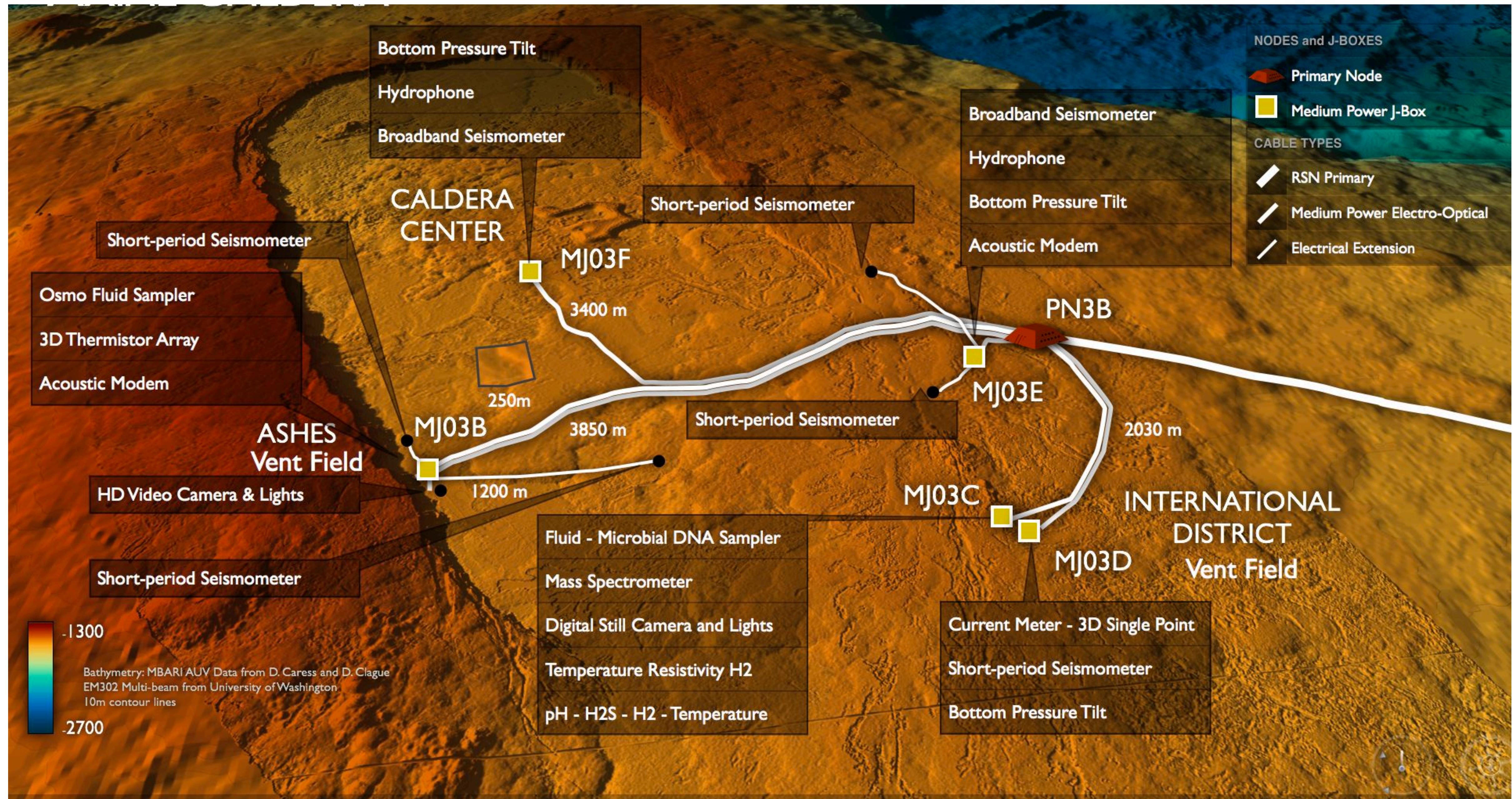
[pangeo-docker-images](#)
Forked from pangeo-data/pangeo-docker-images
An experiment to simplify pangeo docker images
Dockerfile MIT 7 stars 0 forks Updated on Mar 30

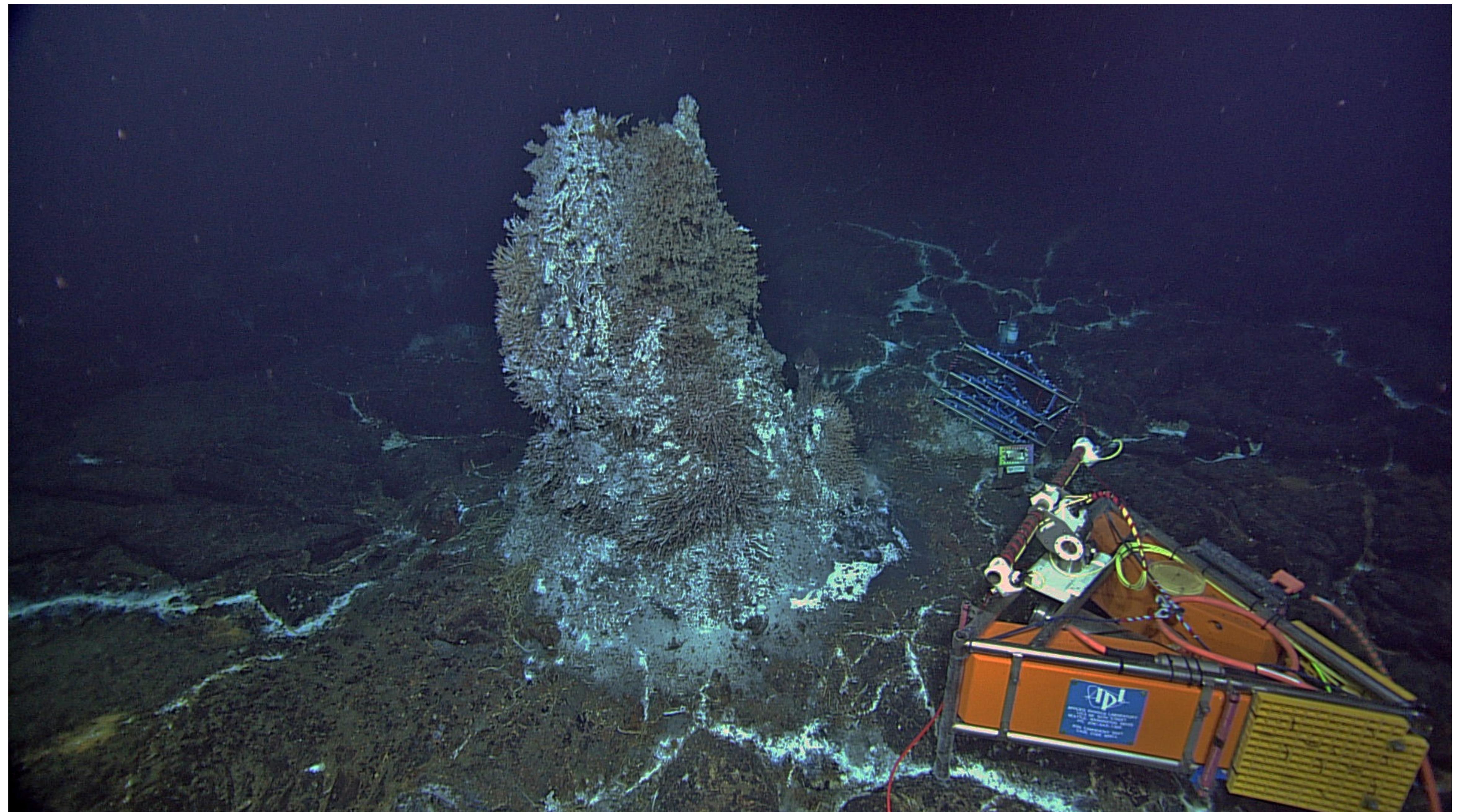


<https://oceanobservatories.org>

<https://ooicloud.org>

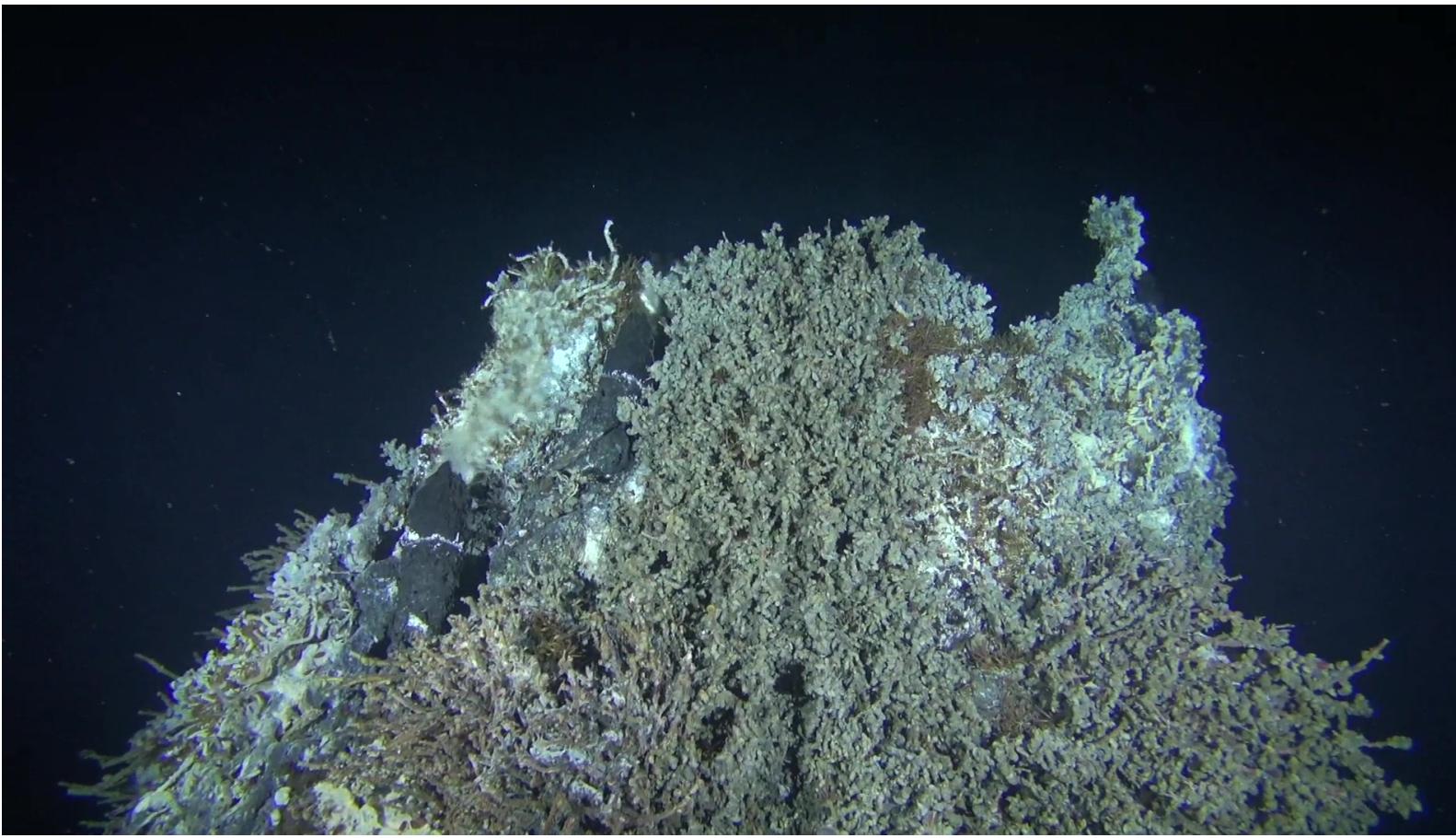
<https://github.com/ooicloud>



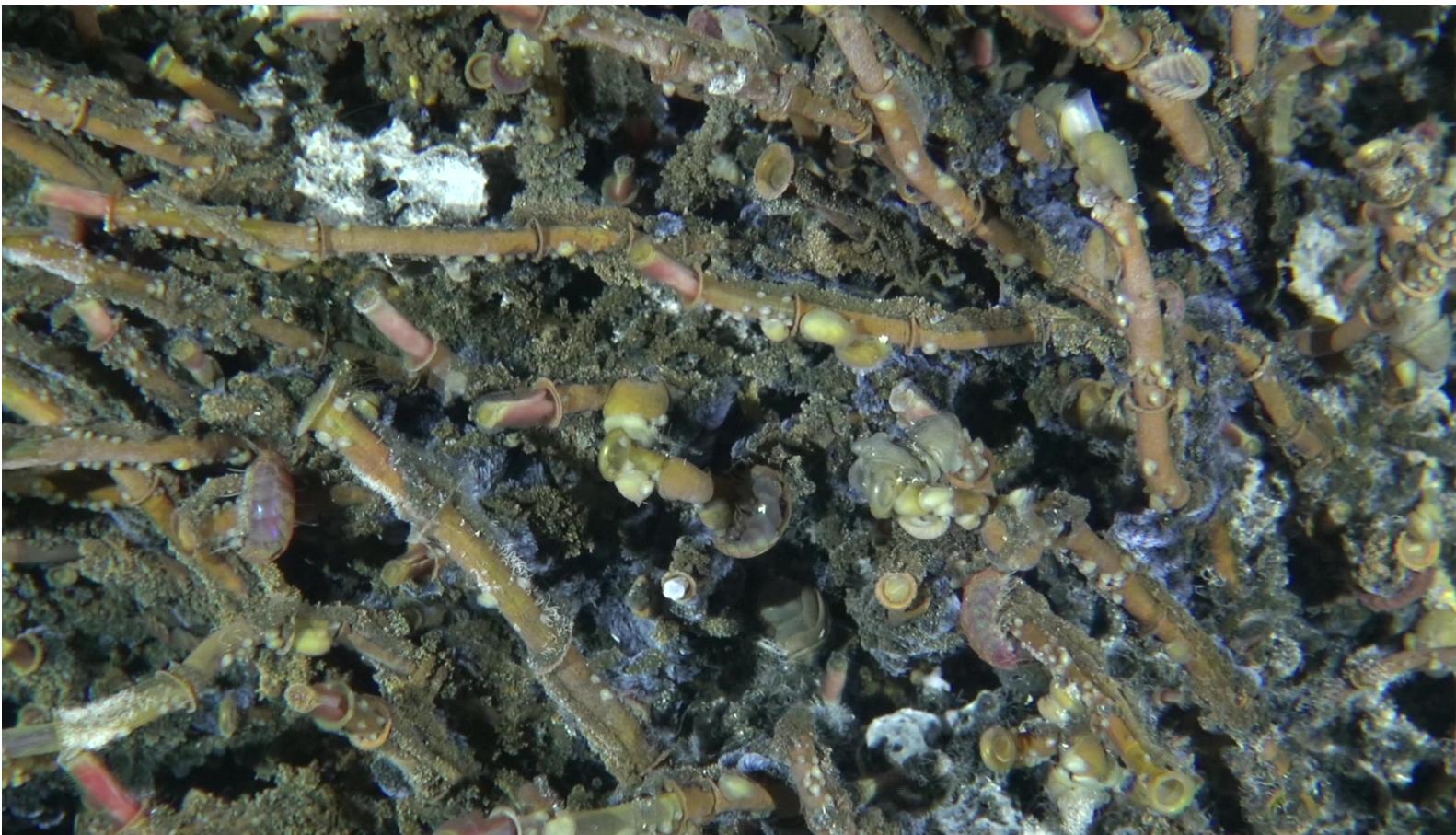


Science!

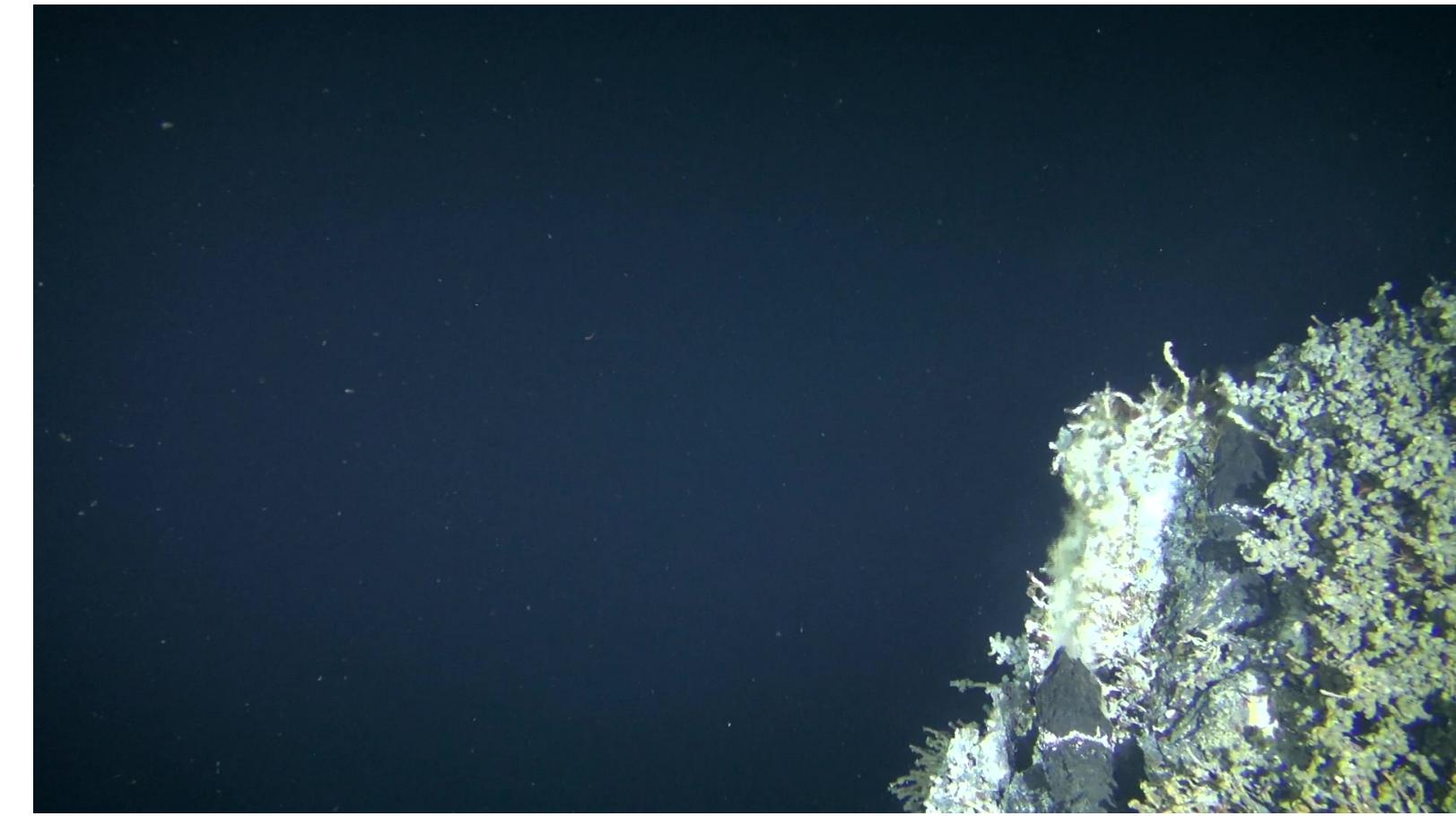
Geology



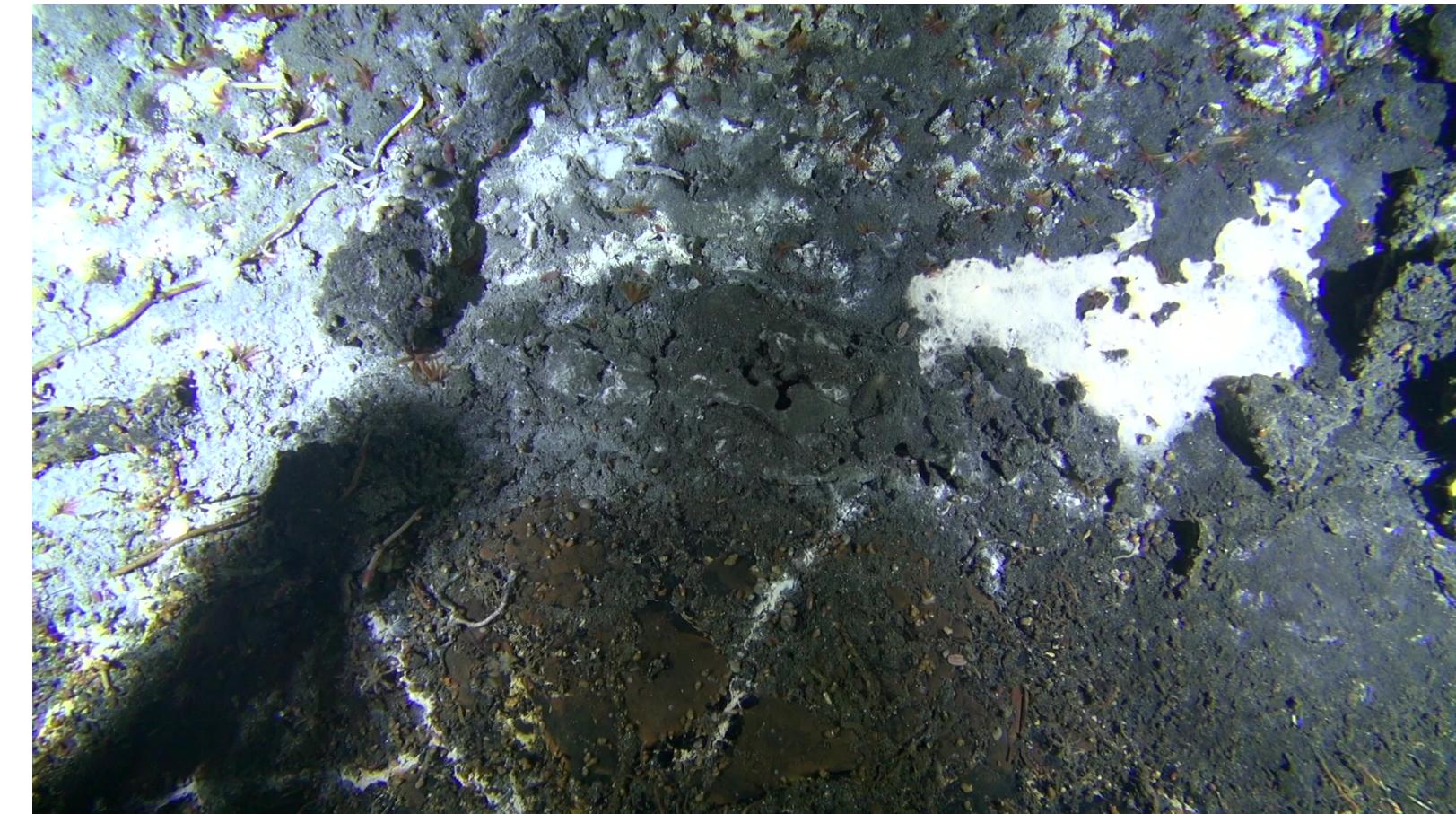
Macrofauna



Oceanography and Hydrology



Bacterial Mats



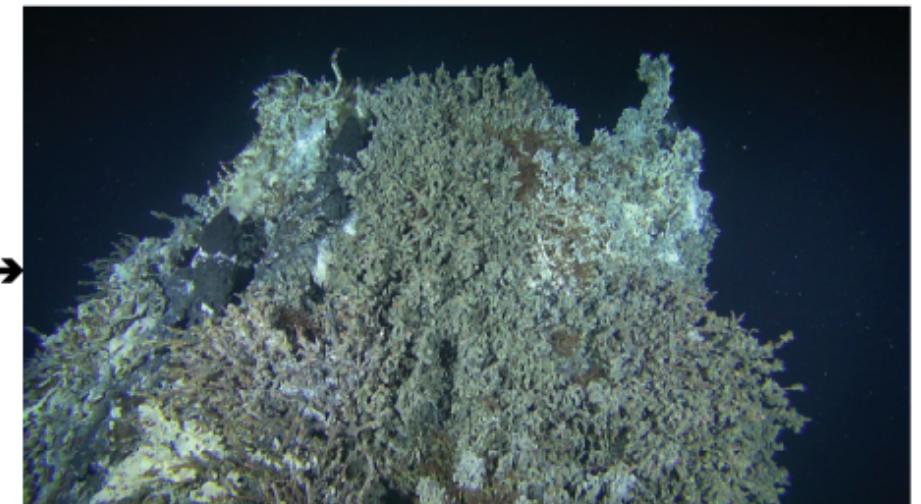
PyCamHD

<https://github.com/tjcrone/pycamhd>

- Python library for working with the CamHD raw data archive
- Image extraction functions
- File information functions
- Archive information functions

Archive Stats

```
camhd.get_stats() Return the total number of MOV files and the total size of the MOV files (in TB) in the data archive. Returns an integer and a float.  
camhd.get_file_list() Return a list of all MOV files in the data archive as fully-qualified URLs. Returns a list of strings.
```



Study: vent structure growth, dynamics, fluid flow rates, heat and chemical fluxes.

File Information

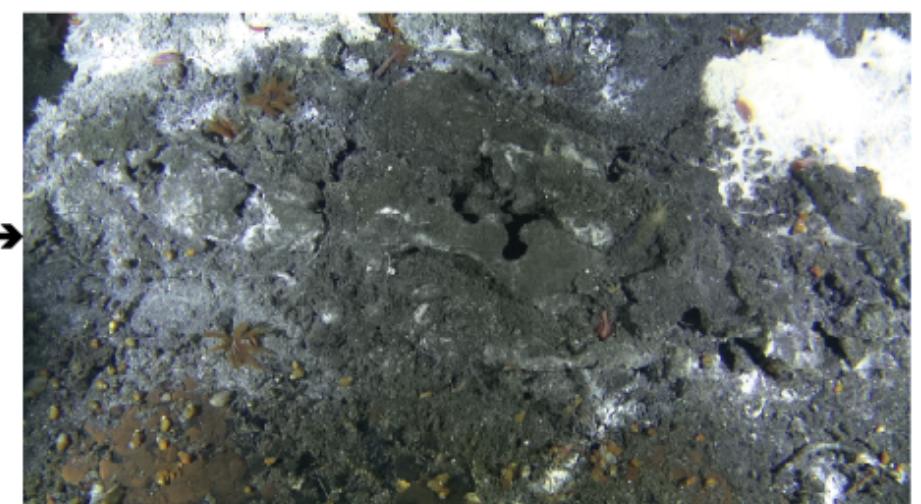
```
camhd.get_atom_sizes(filename) Return the sizes of the three top-level atoms in a remote file. Returns three integers.  
camhd.get_frame_count(filename[, moov_atom]) Return the number of frames in a remote file. Returns an integer.  
camhd.get_frame_sizes(filename[, moov_atom]) Return the sizes of all frames in a remote file. Returns a list of integers.  
camhd.get_frame_offsets(filename[, moov_atom]) Return the offsets of all frames in a remote file. Returns a list of integers.
```



Study: macrofaunal population dynamics and interactions.

File Components

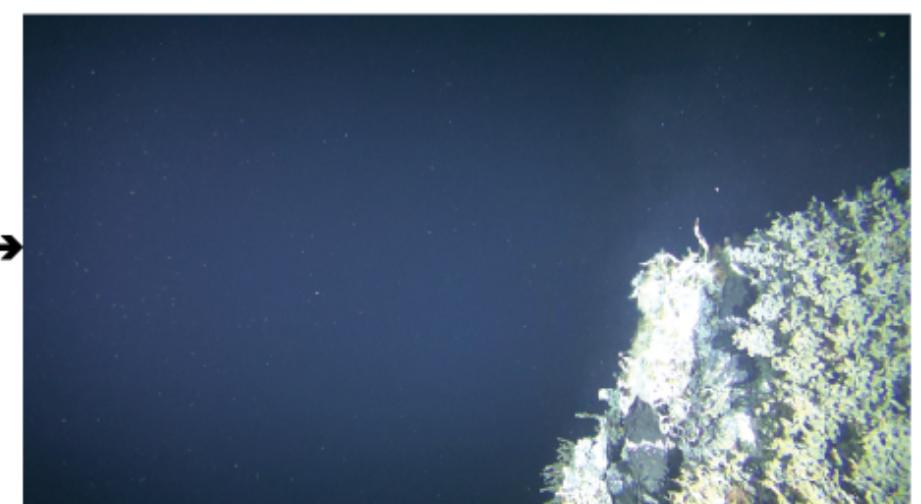
```
camhd.get_moov_atom(filename) Retrieve the moov atom from a remote file. Returns a string containing raw packed binary data.  
camhd.get_frame_data(filename, frame_number[, moov_atom]) Retrieve the raw ProRes encoded frame data from a frame in a remote file. Returns a string containing raw packed binary data.  
camhd.get_avi_file(frame_data) Adds an appropriately structured AVI header to frame_data. frame_data should be a string containing raw packed binary data as returned by get_frame_data(). Returns a string containing raw packed binary data.
```



Study: sulfide alteration, colonization, chemosynthesizing bacterial mats, things that snails do.

Write Output

```
camhd.write_frame(filename, frame_number[, moov_atom]) Writes a single-frame AVI file. The resulting AVI file can be converted to a TIFF, PNG, YUV, or another image or movie format using ffmpeg. YUV conversions are lossless, as would be conversions to any valid container format using a video stream copy. All CamHD ProRes encoded video frames are key frames
```



Study: suspended bacterial floc concentration, water column turbidity.

Interactive Example

```
>>> import camhd  
>>> filename = 'https://rawdata.oceanobservatories.org/files/RS03ASHS/PN03B/06-CAMHDA301/2016/07/29/CAMHDA301-20160729T210000Z.mov'  
>>> moov_atom = camhd.get_moov_atom(filename)  
>>> frame_set = [1980, 7845, 11460, 15930]  
>>> for frame_number in frame_set:  
...     camhd.write_frame(filename, frame_number, moov_atom)
```