

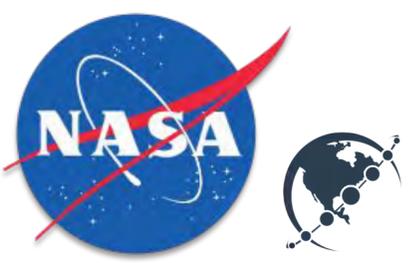
# Empowering Transformational Science

Chelle Gentemann (Farallon Institute) twitter: [@ChelleGentemann](https://twitter.com/ChelleGentemann)

Ryan Abernathy (Columbia / LDEO) twitter: [@rabernat](https://twitter.com/rabernat)

Aimee Barciauskas (Development Seed) twitter: [@aimeeb](https://twitter.com/aimeeb)

(there are lots of links in this presentation! click away!)



**EARTH SCIENCE**  
DATA SYSTEMS

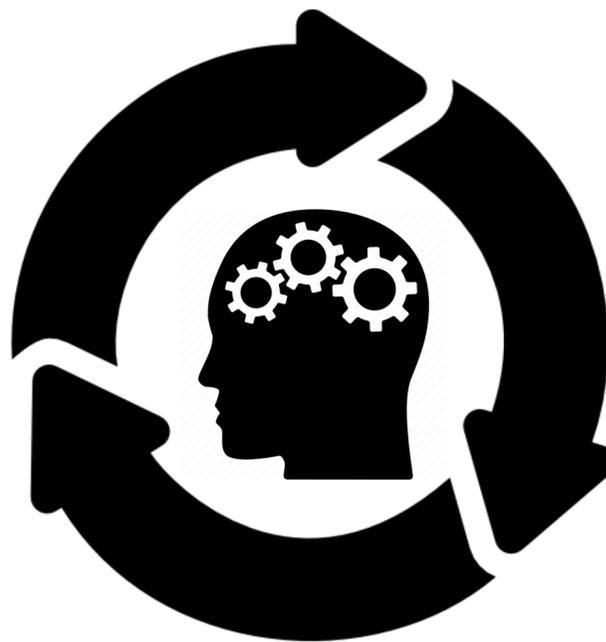
NASA Physical Oceanography Program



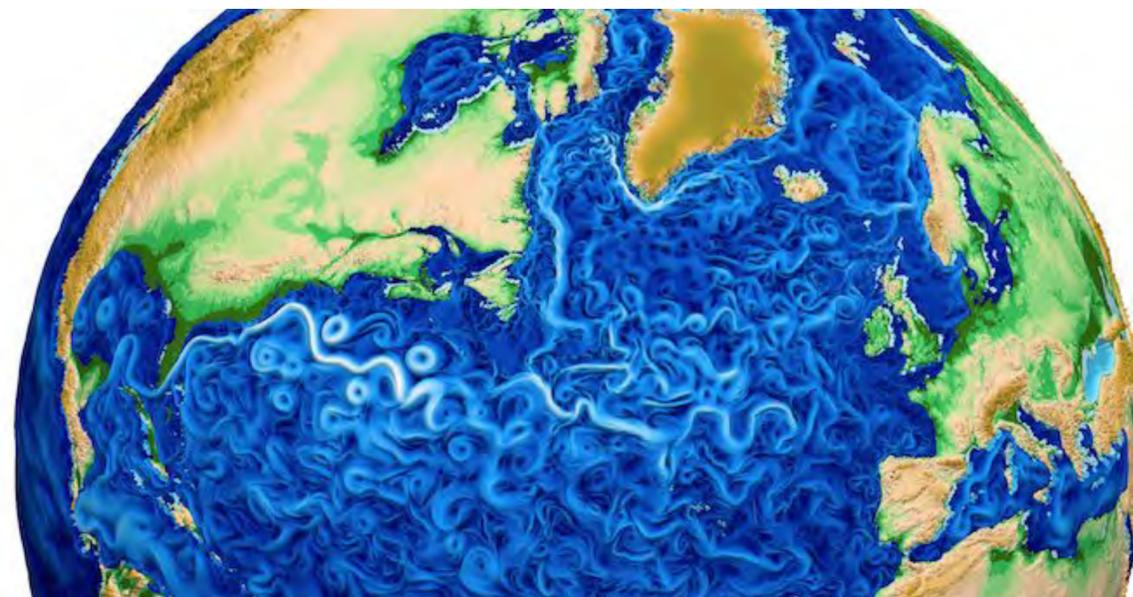
# What Drives Progress in Earth System Science?

**New Ideas / Hypotheses**

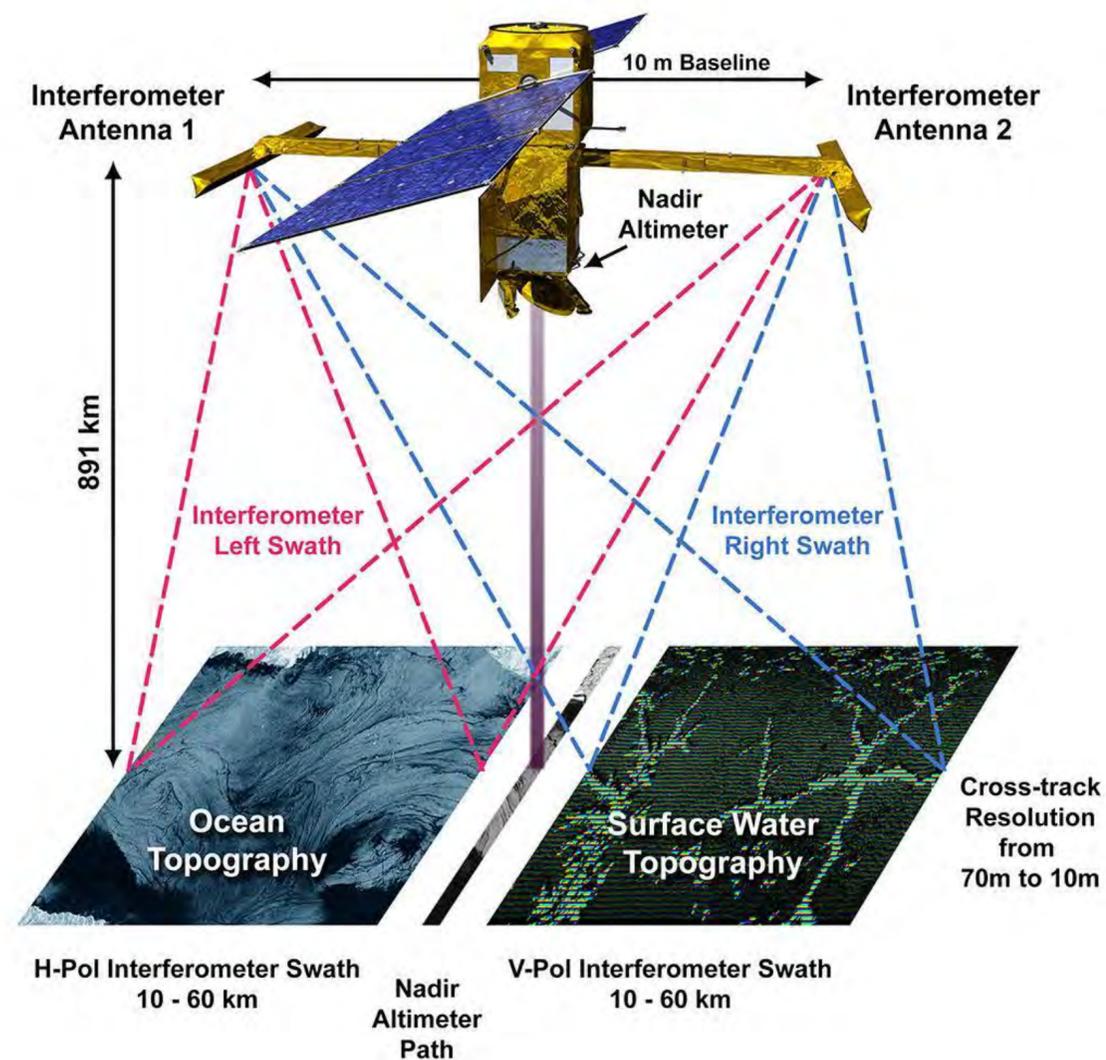
$$\rho \frac{Du}{Dt} = -\nabla p + \nabla \cdot \tau + \rho g$$



**New Simulations**



**New Observations**



# What impacts the **velocity** of progress? Data, Software, & Compute

**Data:** time to find, access, clean, & format data for analysis

**Software:** what tools are easily available?

**Compute:** access to compute == speed of results

80%  
Data Preparation  
(download, clean, & organize files)

10%  
Batch  
Processing

10%  
Think about  
science

# Traditional methods of data access cannot leverage large volumes of data

US\$1.5 billion

With a total cost estimated at US\$1.5 billion, NISAR is likely to be the world's most expensive Earth-imaging satellite.



en.wikipedia.org › wiki › NISAR\_(satellite) ▾

[NISAR \(satellite\) - Wikipedia](#)

140 petabytes

As Dobson notes, **NISAR** is expected to generate a tremendous **volume** of data over its scheduled three-year **mission** – as much as 140 petabytes (PB).

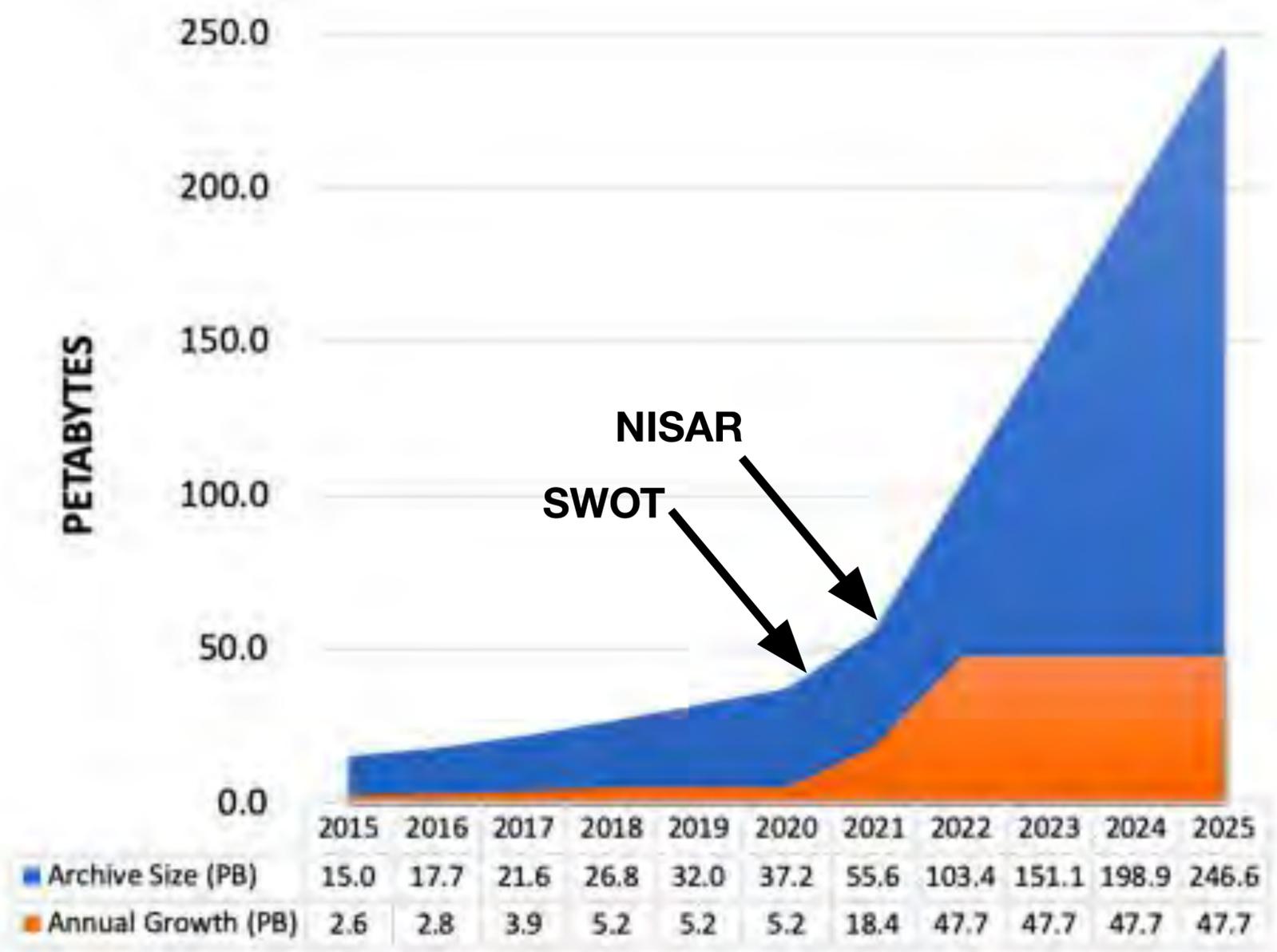
asf.alaska.edu › wp-content › uploads › 2019/06 › 201... ▾ PDF

[Getting Ready for NISAR - Alaska Satellite Facility](#)

# Data, Software, Compute

DATA ANALYTICS

New climate model data now in Google Public Datasets



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



Registry of Open Data on AWS

## Multi-Scale Ultra High Resolution (MUR) Sea Surface Temperature (SST)

climate earth observation environmental natural resource oceans satellite imagery sustainability water weather

### Description

A global, gap-free, gridded, daily 1 km Sea Surface Temperature (SST) dataset created by 2 satellite SST datasets. Those input datasets include the NASA Advanced Microwave Scanning Radiometer 2 (AMSR-2) on GCOM-W and the JAXA Advanced Microwave Scanning Radiometer 2 (AMSR-2) on GCOM-W Resolution Imaging Spectroradiometers (MODIS) on the NASA Aqua and Terra platforms, WindSat radiometer, the Advanced Very High Resolution Radiometer (AVHRR) on several situ SST observations from the NOAA iQuam project. Data are available from 2002 to present. The original source of the MUR data is the NASA JPL Physical Oceanography DAAC.

### Update Frequency

Daily

### License

There are no restrictions on the use of these data.

### Documentation

<https://podaac.jpl.nasa.gov/dataset/MUR-JPL-L4-GLOB-v4.1>

### Managed By

<https://faralloninstitute.org>

See all datasets managed by <https://faralloninstitute.org>.

### Contact

[podaac@podaac.jpl.nasa.gov](mailto:podaac@podaac.jpl.nasa.gov)

### Usage Examples

#### Tutorials

- Python Jupyter Notebooks by Chelle Gentemann, Rich Signell
- Python Reader Software by PO DAAC

Microsoft Azure | Open Datasets

### GOES-16

SatelliteImagery EarthObservation AlforEarth NOAA

Overview Data access

Weather imagery from the GOES-16 satellite.

The GOES-R (Geostationary Operational Environmental Satellite) program images weather phenomena from a set of satellites in geostationary orbits. The GOES-16 satellite is the first of four planned GOES-R satellites; GOES-16's orbit provides a view of the Americas.

This dataset currently includes the ABI-L2-MCMIPF product (Advanced Baseline Imager, Level 2, Cloud and Moisture Imagery, Full-disk). We may on-board other GOES-16 and GOES-17 products on request; please contact [goearthdatasets@microsoft.com](mailto:goearthdatasets@microsoft.com) if you are interested in using additional GOES data on Azure.

This dataset is available on Azure thanks to the [NOAA Big Data Program](#).

#### Storage resources

Data are stored in blobs in [NetCDF](#) format (one blob per image) in the East US data center, in the following blob container:

<https://goes.blob.core.windows.net/noaa-goes16>

Within that container, data are named as:

[product]/[year]/[day]/[hour]/[filename]

# Analytics Optimized Data Store (AODS)

a few examples of AODS formats



## Current method -

**NetCDF files** - organized into 'reasonable' data sizes per file, usually by orbit, granule, or day. Filename has information about date, sensor, version. Reading usually involved calculating the filename, opening, reading, processing, closing.

## Analytics Optimized Data Store (one example of many different formats)

**Zarr** - makes large datasets easily accessible to distributed computing. Original data is stored in directories each having chunked data corresponding to dataset dimensions. Metadata is read by zarr libraries to read only the chunks necessary to complete a subsetting request.

## Technology advances -

**Lazy loading** - also known as asynchronous loading - defer initialization of an object until the point at which it is needed. Developed for webpages. Delays reading data until needed for compute.

## Advanced OSS libraries:

**Xarray** - library for analyzing multi-dimensional arrays, lazy loading.

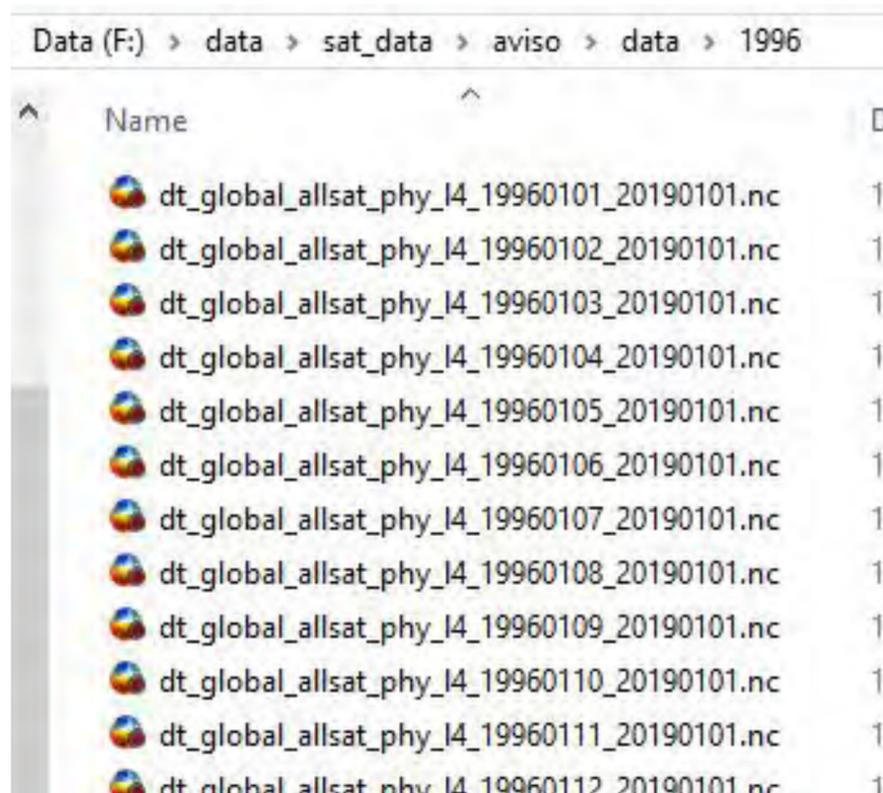
**Dask** - able to break a large computational problems into a network of smaller problems for distribution across multiple processors

**Intake** - lightweight set of tools for loading and sharing data in data science projects

# What does a data store look like?

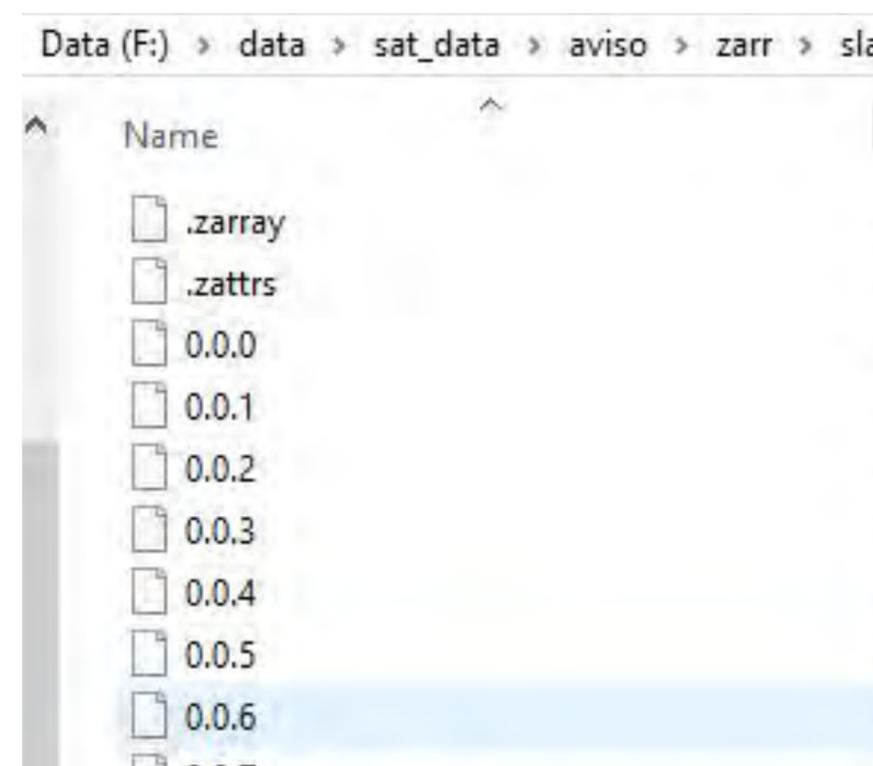
## NetCDF

Organized so that each file can fit into RAM, usually by day, orbit, or granules



## Zarr

organization and format invisible to user, data accessed by metadata



# How to access data?



My version of lazy loading before I knew python - on bedrest, pregnant with twins

```
jupyter Test_CCMP_zarr_simple_version Last Checkpoint: 4 minutes ago (autosaved)
File Edit View Insert Cell Kernel Widgets Help
+ -> Run Code
Import libraries
In [1]: import xarray as xr
        from glob import glob
        import numpy as np

In [2]: pattern_netcdf = 'F:/data/sat_data/ccmp/v02.0/*/*_v02.0_L3.0_RSS.nc'
        pattern_zarr = 'F:/data/sat_data/ccmp/zarr/'

Reading netCDF files
In [3]: %%time
        #list files
        files = [x for x in glob(pattern_netcdf)]
        #open dataset
        ds=xr.open_mfdataset(files,combine='nested',concat_dim='time')
        Wall time: 12min 28s

Reading Zarr files
In [4]: %%time
        ds= xr.open_zarr(pattern_zarr)
        Wall time: 557 ms

Advantage of using these tools
• Create a 32 year timeseries at a single point in < 1 min, because only the data chunks that are needed are loaded

In [5]: %%time
        ts = ds.uwnd.sel(latitude=slice(-10,0),longitude=slice(170,180)).mean({'latitude','longitude'}).plot()
        Wall time: 33.6 s
```

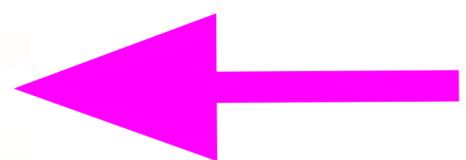
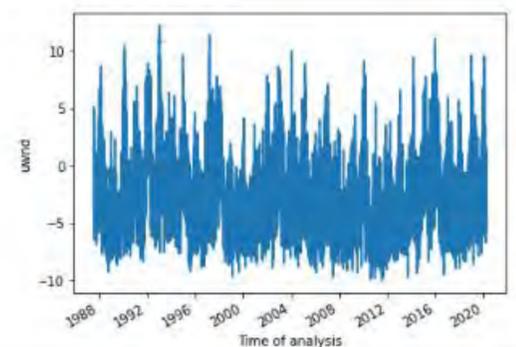
Modern software tools use lazy loading to access large datasets

Reading in netCDF data: **13 minutes (depends on computer)**

- 1 - user creates list of filenames
- 2 - access dataset by reading the metadata distributed through files

Reading in Zarr data: **0.1 seconds (metadata consolidated)**

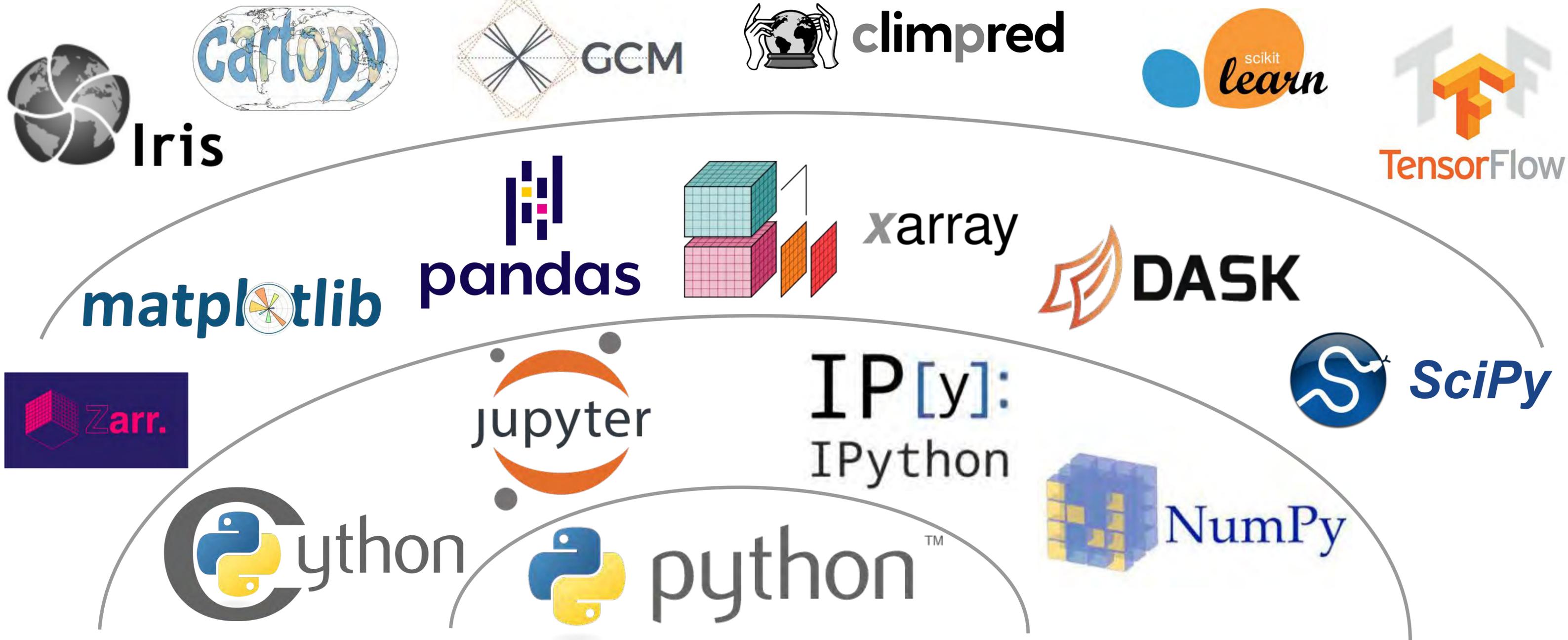
- 1 - access dataset by reading the consolidated metadata



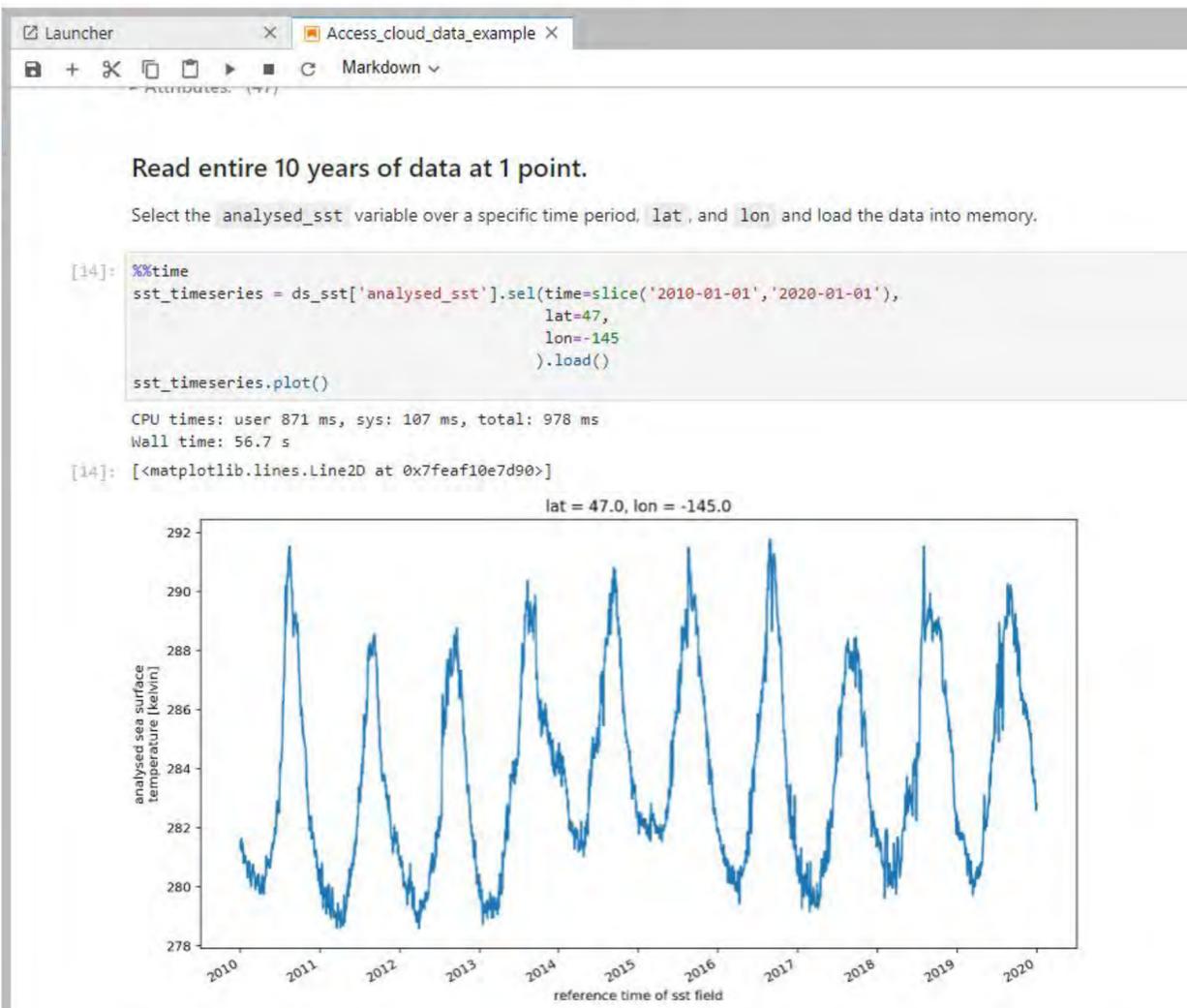
**STOP ----- THIS IS DIFFERENT -----**

- 1 line of code to access a 32-year, global, 25km dataset
  - 1 line of code to select a region, calculate mean, & plot time series
- in LESS than 1 minute**

# Data, Software, Compute







Coding with coordinates == CLARITY

Make a timeseries:

```
data.sel(latitude=10,longitude=0,method='nearest').plot()
```

Regrid data for multivariate analysis:

```
data_regrid = data.interp_like(other_data)
```

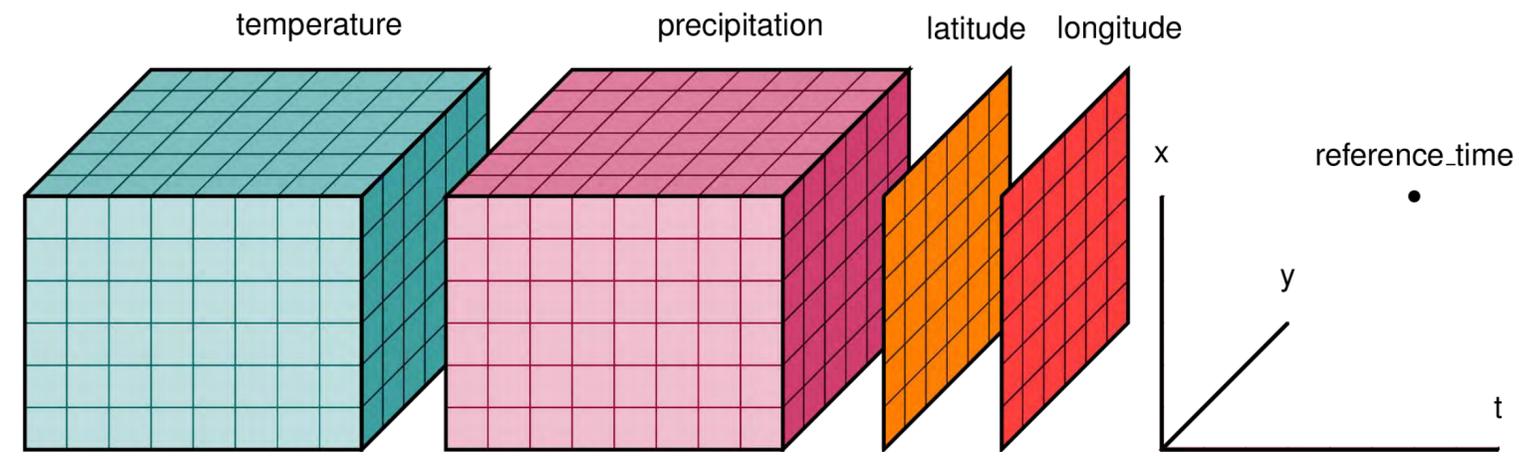
# Xarray

@xarray\_dev

NUMFOCUS  
OPEN CODE = BETTER SCIENCE



Xarray introduces labels in the form of dimensions, coordinates and attributes on top of raw NumPy-like multidimensional arrays, which allows for a more intuitive, more concise, and **less error-prone** developer experience.

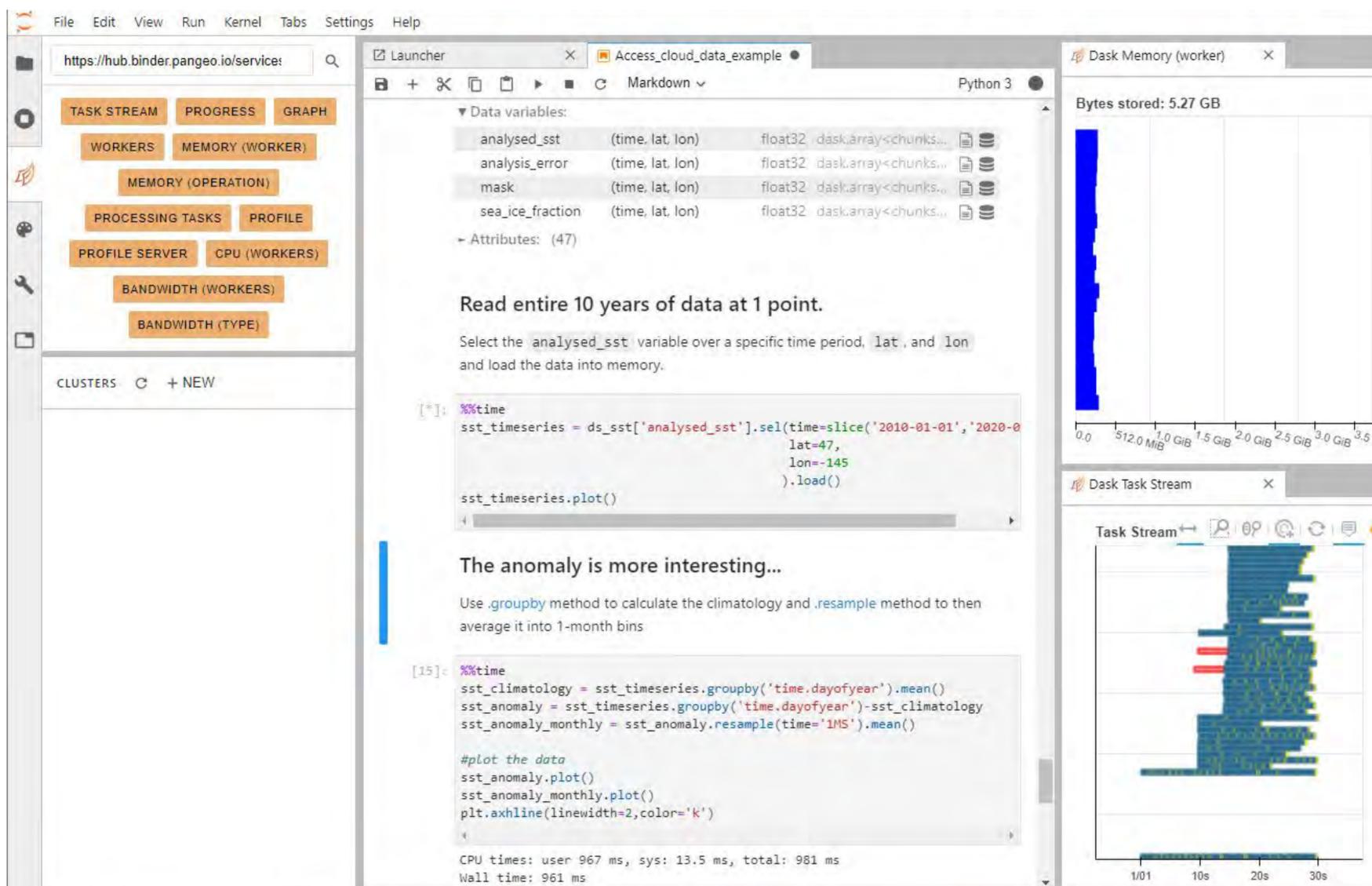


<https://github.com/pangeo-gallery/osm2020tutorial>



# Dask

[@dask\\_dev](#)



Launcher Access\_cloud\_data\_example Python 3

Data variables:

analysed_sst	(time, lat, lon)	float32	dask.array<chunks...
analysis_error	(time, lat, lon)	float32	dask.array<chunks...
mask	(time, lat, lon)	float32	dask.array<chunks...
sea_ice_fraction	(time, lat, lon)	float32	dask.array<chunks...

Attributes: (47)

Read entire 10 years of data at 1 point.

Select the `analysed_sst` variable over a specific time period, `lat`, and `lon` and load the data into memory.

```
[*]: %%time
sst_timeseries = ds_sst['analysed_sst'].sel(time=slice('2010-01-01', '2020-01-01'),
                                             lat=47,
                                             lon=-145)
sst_timeseries.load()

sst_timeseries.plot()
```

Dask Memory (worker)

Bytes stored: 5.27 GB

Dask Task Stream

Task Stream

```
[15]: %%time
sst_climatology = sst_timeseries.groupby('time.dayofyear').mean()
sst_anomaly = sst_timeseries.groupby('time.dayofyear') - sst_climatology
sst_anomaly_monthly = sst_anomaly.resample(time='1MS').mean()

#plot the data
sst_anomaly.plot()
sst_anomaly_monthly.plot()
plt.axhline(linewidth=2, color='k')
```

CPU times: user 967 ms, sys: 13.5 ms, total: 981 ms  
Wall time: 961 ms

Dask is a flexible library for parallel computing in Python.

Xarray integrates with Dask to support parallel computations and streaming computation on datasets that don't fit into memory. When you are using Xarray, you are likely using Dask without even realizing it.

The Pangeo binder and jupyterhub use [Dask Gateway](#) to manage access to the Dask clusters & kubernetes.

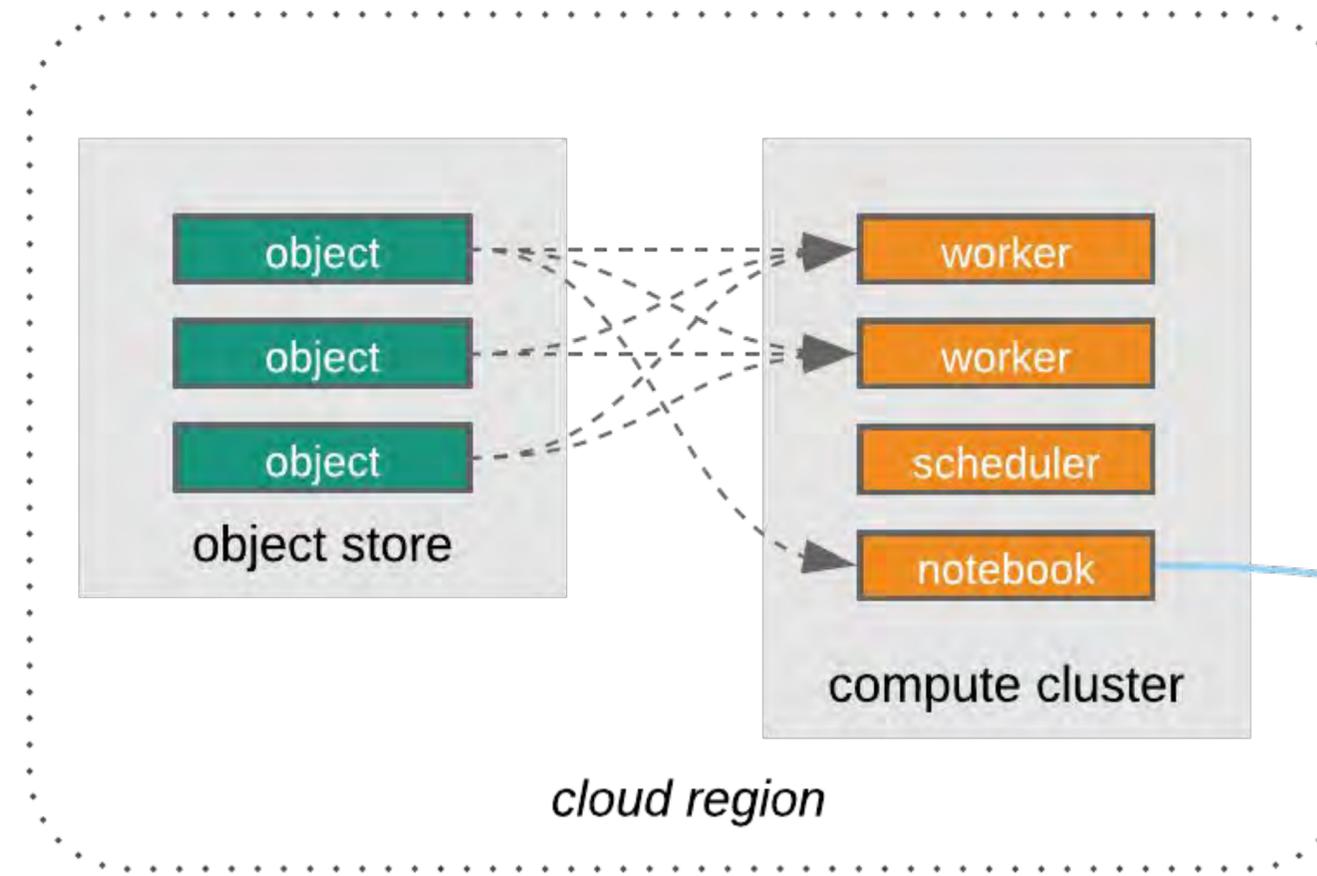
Cluster performance visualization using the Dask viewers: memory use, profile, CPU, etc.

# Data, Software, Compute

Analytics Optimized Data Store (AODS)



Scalable Parallel Computing Frameworks



Data Provider's \$

Data Consumer's \$



# Agency driven solutions

Earth Science Data Systems (ESDS) Program **Multi-Mission Algorithm and Analysis Platform (MAAP)**

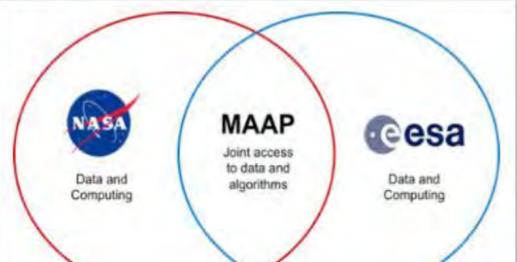
**ESDS Program**

- Earth Science Data Systems Program
- Program Components
- IMPACT
- Competitive Programs
- Commercial Smallsat Data Acquisition Program
- Multi-Mission Algorithm and Analysis Platform (MAAP)
- ESDS Geographic Information Systems Team (EGIST)
- Continuous Evolution
- Adding New Data to EOSDIS
- Open Data, Services and Software Policies

## Multi-Mission Algorithm and Analysis Platform (MAAP)

**About MAAP**

The Multi-Mission Algorithm and Analysis Platform (MAAP) a collaborative project between NASA and ESA, designed to support aboveground biomass research. MAAP will bring together relevant data, algorithms, and computing capabilities in a common environment in order to address the challenges of sharing and processing data from field, airborne, and satellite measurements related to ESA and NASA missions.



**Challenge and Solution**

New missions such as NASA's GEDI, ESA's B... exponentially higher than any currently op... storage, processing, and sharing challenge... heterogeneous nature of the data, which a... resolutions, coverages, and processing lev... and immediate need for improved data sh...

MAAP is addressing these community need...

- Enabling researchers to easily discover, NASA missions and validation/calibratio
- Harmonizing satellite, airborne, and gro data generation.
- Developing tools for repeatable and sha

SENTINEL Hub by SINERGISE

EXPLORE ▾ DEVELOP ▾ ABOUT ▾ PRICING BLOG SIGN



## CLOUD API SATELLITE IMAGERY

Use. Pick. Enhance. Expose.

Import imagery

[Button]

geohazards tep Site ▾ Page ▾ Source

## ESA CloudToolbox

The ESA CloudToolbox is a Virtual Machine (VM) that offers a flexible amount of CPUs, RAM and dedicated storage, tailored to user needs and type of machine required. When needed, users can request upgrades of the configuration (for example, asking more processing power) at any time, compatibly with the Cloud infrastructure constraints. A pre-built VM template offers ready-to-use machines for SAR Interferometric processing or generic EO data processing. However, besides free and licensed software tools (e.g. Sentinel-1 toolbox, NEST, GAMMA, Matlab, etc) that can be installed on the machines, users may request installation of additional tools.

### Create a CloudToolbox

To create your own CloudToolbox:

- Access the cloud dashboard (see Cloud Dashboard)
- Click on  to create a new Virtual Machine
- Set the Virtual Machine name (e.g 'my esa toolbox')
- Select the **ESA CloudToolbox** template
- Click on **Create**
- Wait for the VM to be deployed
- Get the <ESA CloudToolbox IP>.

### Create Virtual Machine

my esa toolbox

 **ESA Cloud Toolbox**

# Grass-Roots Solutions


**Erin Dougherty**  
 @edougherty\_

My #dayofscience: paper revisions 🤖 and using @pangeo\_data to analyze massive amounts of high-res climate data to understand floods in a current and future climate over the U.S. When not doing this, I ❤️ observing #wx directly via field work and watching storms. 🌩️



2:38 PM · Oct 15, 2019 · [Twitter for iPhone](#)

5 Retweets 51 Likes


**Patrick Gray**  
 @clifgray

One of the highlights of #AGU19 for me was a workshop on @pangeo\_data. If you're interested in earth sci + geospatial analysis at scale I can't recommend their tutorial enough. Find it at [github.com/pangeo-data/pa...](https://github.com/pangeo-data/pangeo-tutorial) Clear and concise intro to #xarray, #dask, #geopandas, and #intake.



pangeo-data/pangeo-tutorial  
 Interactive jupyter notebooks for pangeo tutorial events - pangeo-data/pangeo-tutorial  
[github.com](https://github.com)

2:03 PM · Dec 16, 2019 · [Twitter Web App](#)

15 Retweets 64 Likes


**Andrew Williams**  
 @AndrewWilliams

It's taken a while, but I think that the whole xarray/Dask/cloud thing has finally clicked !

I've been working with CESM Large Ensemble data for a few months now - moved onto a Pangeo server and managed to speed up my workflow massively !

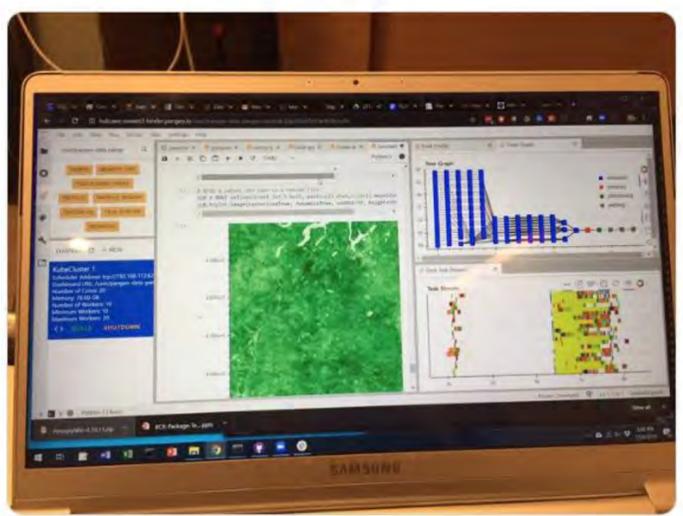
[@xarray\\_dev](#), [@dask\\_dev](#), [@pangeo\\_data](#)

10:23 AM · Jan 12, 2020 · [Twitter Web App](#)

6 Retweets 43 Likes


**Chelle Gentemann**  
 @ChelleGentemann

Rocking 70GB Landsat data at the @pangeo\_data #AGU2019 tutorial. Such a powerful #OpenSource software stack.



10:45 PM · Dec 8, 2019 · [Twitter for iPhone](#)

17 Retweets 92 Likes


**Andrew Pauling**  
 @andrewp109

#cmip6hack is just wrapping up, and has changed the way I will think about, and hopefully do, climate model analysis in the future. The @pangeo\_data infrastructure makes it all so easy.

5:29 PM · Oct 18, 2019 · [Twitter Web App](#)

6 Retweets 24 Likes


**Scott Collis**  
 @Cyclogenesis\_au · Nov 11, 2019

Teaching a course on #@ONScience #OpenScience with @Shobenase at @MonashUni with 20 very diverse attendees @environmentca @argonne @armnewsteam



1 3 13

pangeo Retweeted


**Scott Collis**  
 @Cyclogenesis\_au

Replying to @Cyclogenesis\_au @Shobenase and 4 others

Forgot to mention all done on @pangeo\_data ! Making future Pangeans!

9:15 PM · Nov 11, 2019 · [Twitter Web App](#)

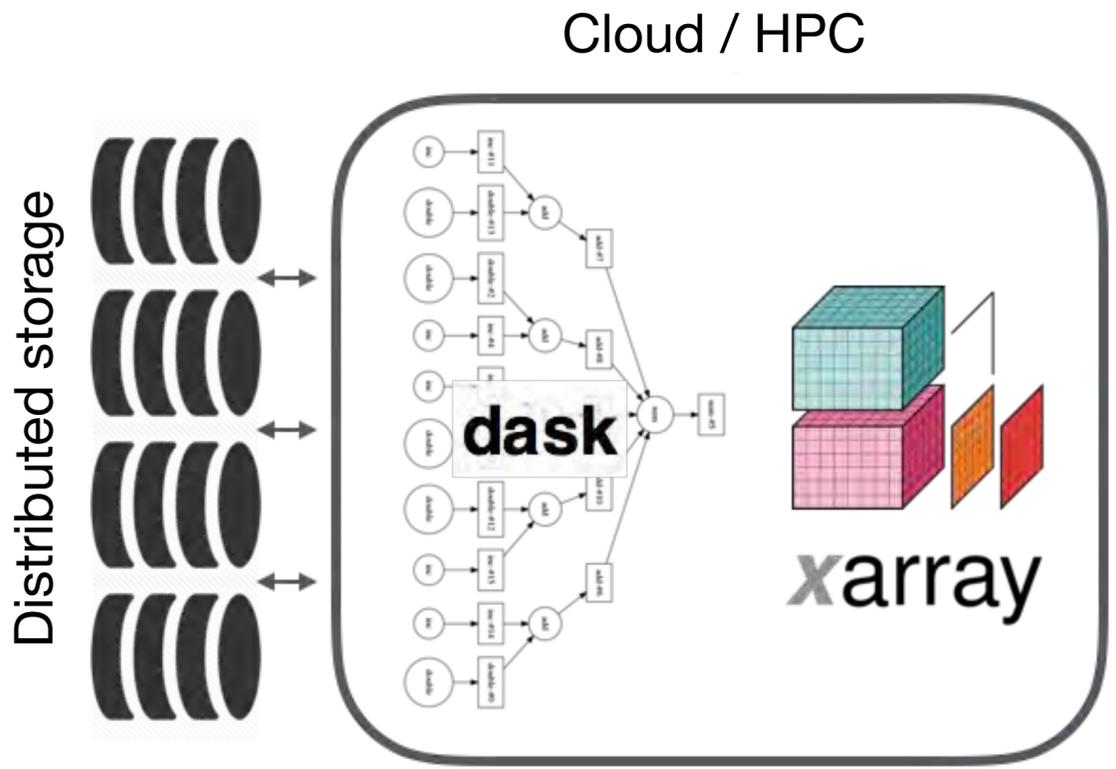
1 Retweet 5 Likes

# Pangeo Architecture



@pangeo\_data

“Analytics Optimized Data Stores” stored on globally-available distributed storage.



Jupyter for interactive data analysis on remote systems

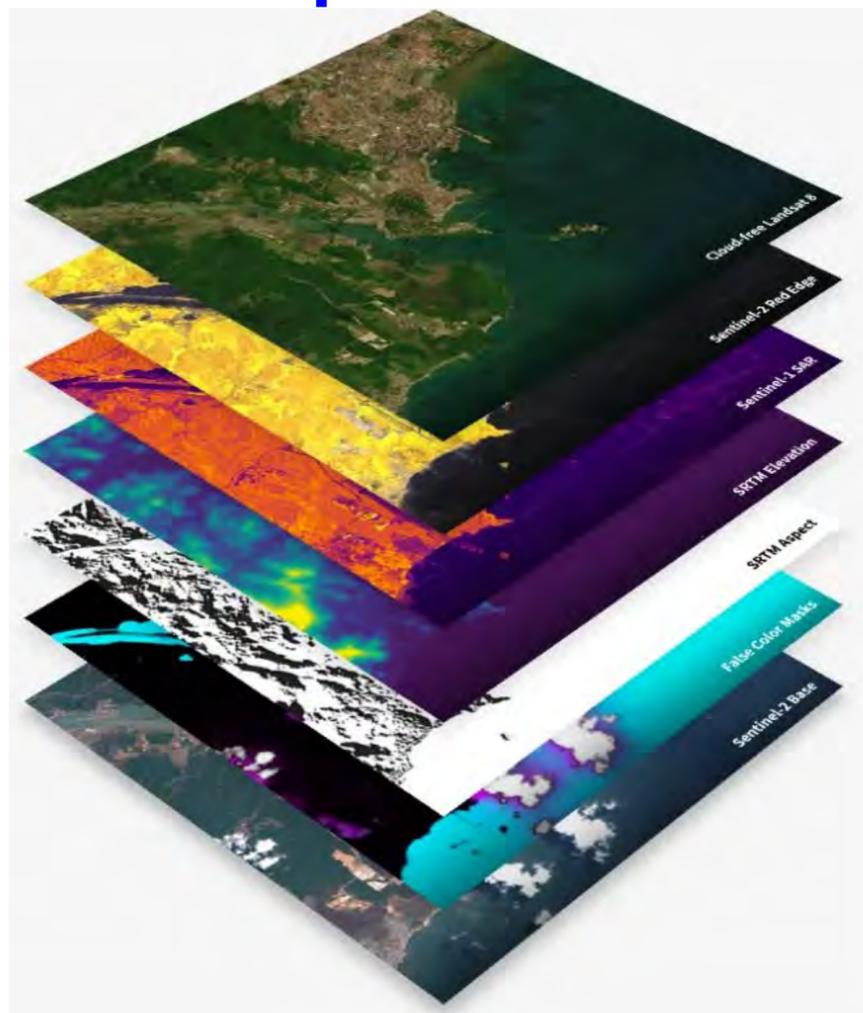
Xarray provides data structures and intuitive interface for interacting with datasets

Parallel computing system allows users deploy clusters of compute nodes for data processing.  
Dask tells the nodes what to do.

# How can data providers reduce barriers?

Reimagine how cloud data access and tools can enable transformational science

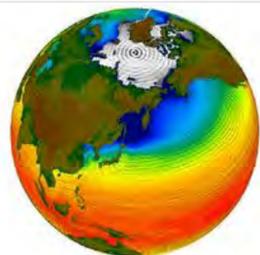
## Publish cloud-optimized data



## Interactive tutorials

### PANGEO GALLERY

Welcome to the Pangeo Gallery website. This site containing one or more notebooks. Each gallery i see the [Contributor Guide](#).



### GALLERY FOR CESM LENS ON AWS

A gallery of notebooks developed to demonstrate analysis of CESM LENS data publicly available on Amazon S3 (us-west-2 region) using xarray and dask

license [BSD-3-Clause](#) last commit [may](#)  
[Binderbot](#) [failing](#) [launch](#) [binder](#)

## Increase user interactions/feedback



## Contribute to OSS tools

Anaconda Cloud 0.1.4 conda-forge 557 pypi v0.1.4 build passing build passing codecov 45% License MIT DOI 10.5281/zenodo.3743397

### cmip6\_preprocessing

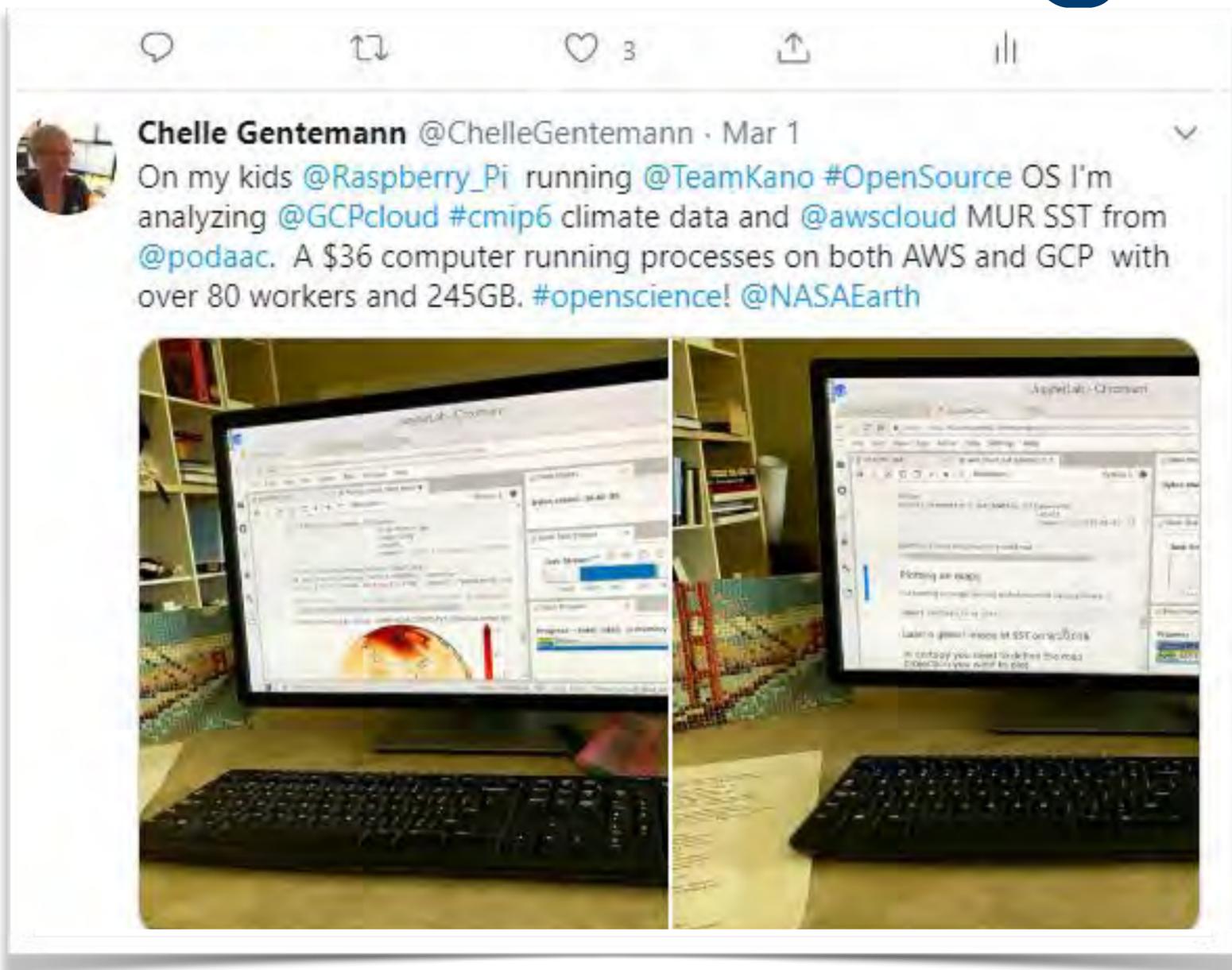
Frustrated with how 'dirty' CMIP6 data still is? Do you just want to run a simple (or complicated) analysis on various models and end up having to write logic for each separate case? Then this package is for you.

Developed during the [cmip6-hackathon](#) this package provides utility functions that play nicely with [intake-esm](#).



**Julius Busecke**  
 @JuliusBusecke Follows you

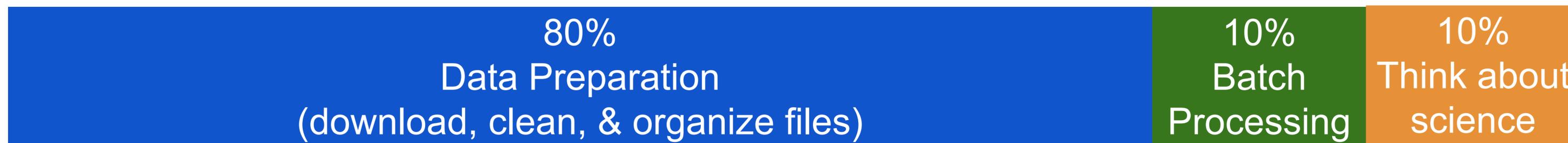
# How does minimizing barriers to data change science?



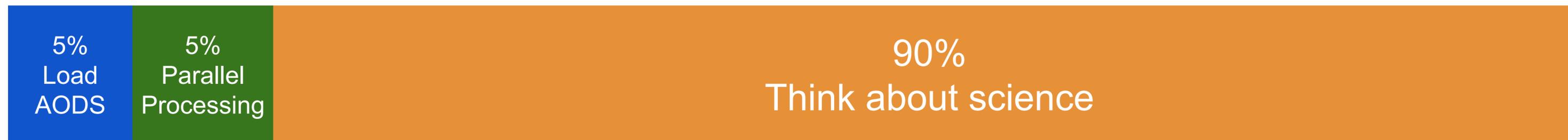
Levels the playing field for all who want to contribute

# Impacts: Reduce Time to Science

## Traditional Project Timeline



## Cloud-based Project Timeline



# Impacts: Reproducibility

*Reproducibility in data-driven science requires more than just code!*

## Traditional Project Code

```
# step 1: open data (stored on local hard drive)
>>> data = open_data("/path/to/private/files")
Error: files not found
```

## Cloud-based Project Code

```
# step 1: open data (globally accessible)
>>> data = open_data("http://catalog.pangeo.io/path/to/dataset")
# step 2: process data
>>> process(data)
```

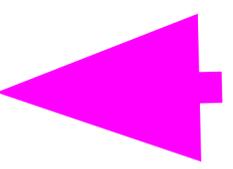
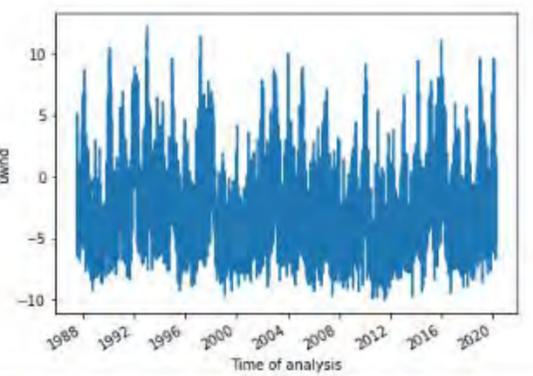
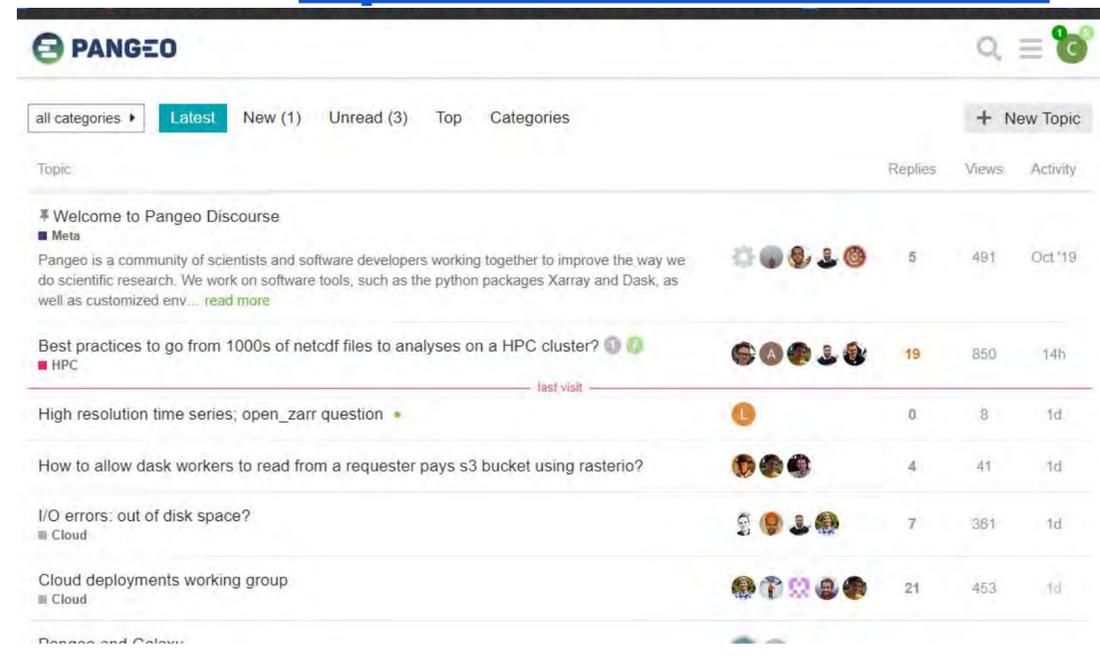
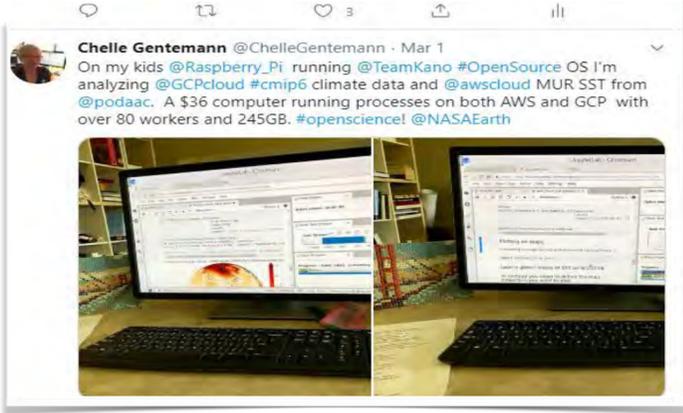
# What impacts the **velocity** of progress? Data, Software, & Compute

Open source science

## Thank you!



**NUMFOCUS**  
OPEN CODE = BETTER SCIENCE



**STOP ----- THIS IS DIFFERENT -----**  
 1 line of code to read in entire 32 year global 25km dataset  
 1 line of code to select a region, calculate & plot a mean time series  
in LESS than 1 minute

