

# Pangeo for Plasma

*Lessons for plasma software from the climate data analytics community*

Thomas Nicholas  
(Columbia University / Lamont-Doherty Earth Observatory)

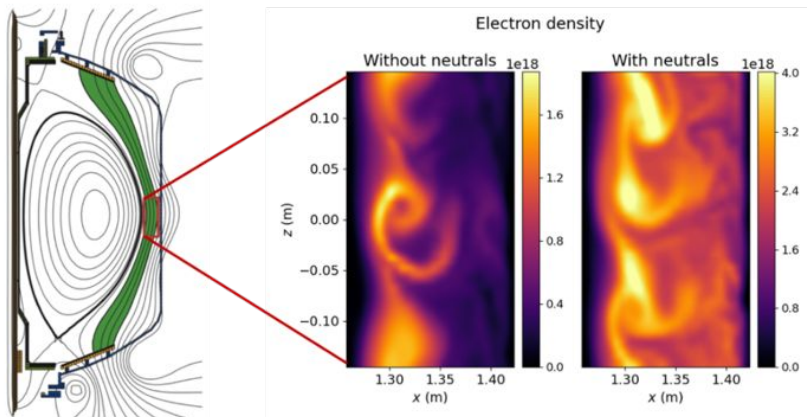
[thomas.nicholas@columbia.edu](mailto:thomas.nicholas@columbia.edu)

# Who am I?

# Who am I?



UNIVERSITY  
*of York*

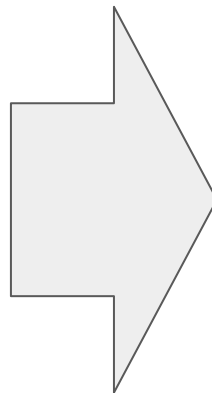
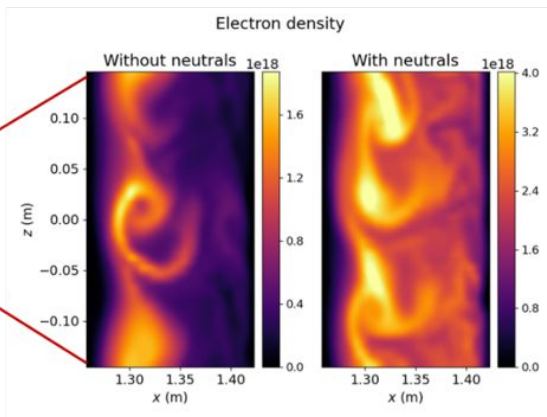
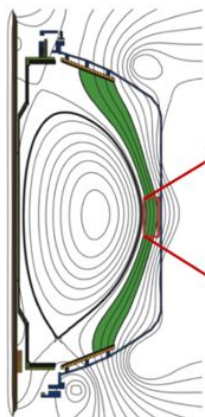


PhD with Ben Dudson, Fulvio Militello, BOUT++

# Who am I?



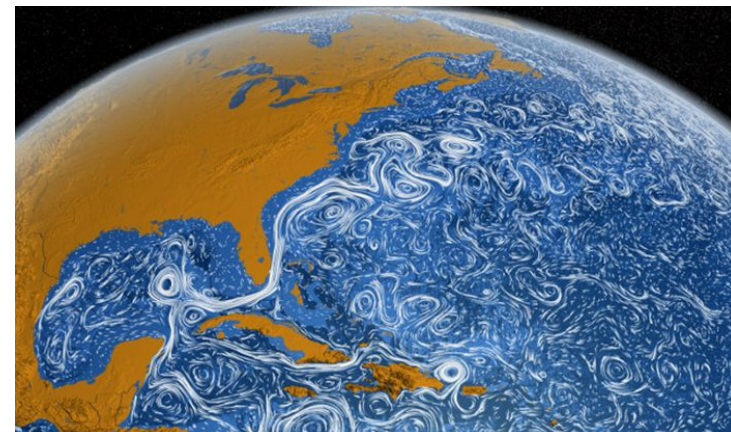
UNIVERSITY  
*of York*



Lamont-Doherty Earth Observatory  
COLUMBIA UNIVERSITY | EARTH INSTITUTE



COLUMBIA UNIVERSITY  
IN THE CITY OF NEW YORK



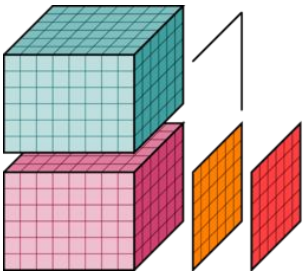
PhD with Ben Dudson, Fulvio Militello, BOUT++

RSE with Ryan Abernathey, various projects

# What do I do now?



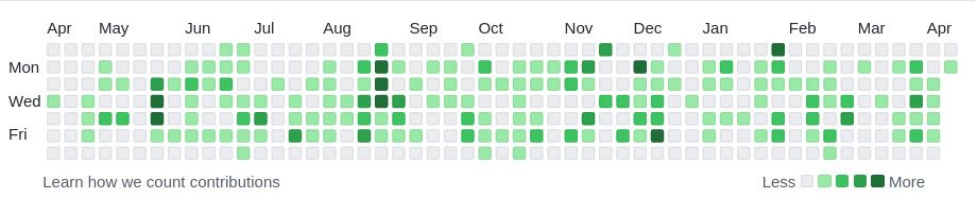
# PANGEO



# xarray

541 contributions in the last year

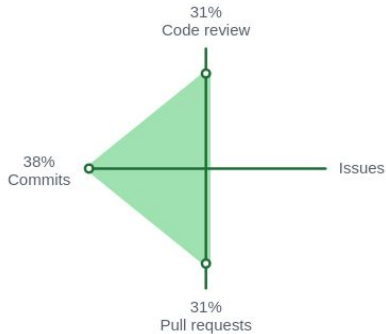
Contribution settings



- @xarray-contrib
- @pydata
- @xgcm
- More

Activity overview

Contributed to [xarray-contrib/datatree](#), [pydata/xarray](#), [xgcm/xgcm](#) and 5 other repositories



# What I hope to convince you of

- Our computational infrastructure **needs to change a lot**
- Can use **solutions from climate science** community
- **Modular approach** makes everyone's work easier
- **Opportunities exist** for plasma coders...

# The White House announces The Federal Year of Open Science



NASA ♦ NSF ♦ NOAA ♦ DOA ♦ DOC ♦ DOE ♦ GSA ♦ NEH ♦ NIH ♦ NIST ♦ USDA ♦ US

Along with other organizations, including CENDI group,  
voluntary collaboration among Federal managers, and  
HELIOS, a coalition of 80+ universities

A multi-agency initiative across the federal  
government to spark change and inspire open  
science engagement through events and activities  
that will advance adoption of open science.

Website: <https://open.science.gov/>

WH: <https://www.whitehouse.gov/ostp/news-updates/>

Nature: <https://doi.org/10.1038/d41586-023-00019-y>



# Climate Science == Plasma Physics

- **Multidimensional** (often fluid turbulent)
- **Large** (bigger than local RAM)
- On regular but warped **grids**
- Often pulled from **central** servers
- From multiple sources but with **common structure** (e.g. experimental and simulation data for same device).



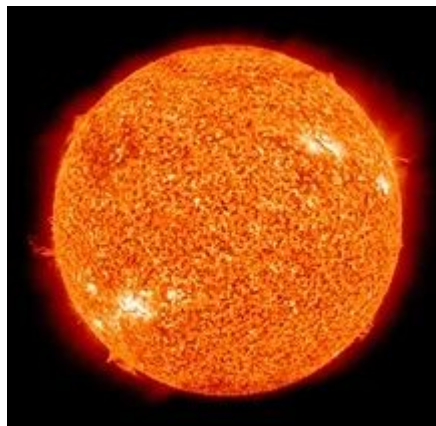
# Climate Science == Plasma Physics

- **Multidimensional** (often fluid turbulent)

- **Large** (bigger than local RAM)

- On regular but warped **grids**

- Often pulled from **central** servers



=



- From multiple sources but with **common structure** (e.g. experimental and simulation data for same device).

# Typical scientific workflow

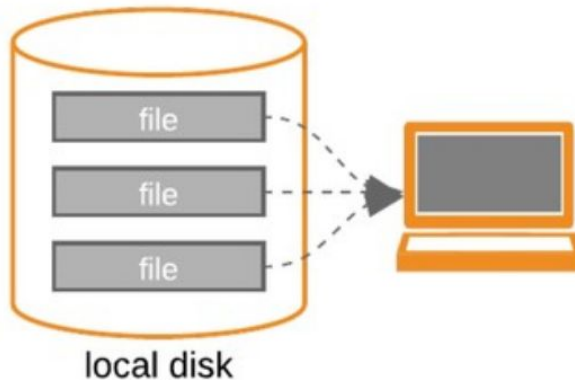
# Typical scientific workflow

## DOWNLOAD

step 1: download



step 2: analyze

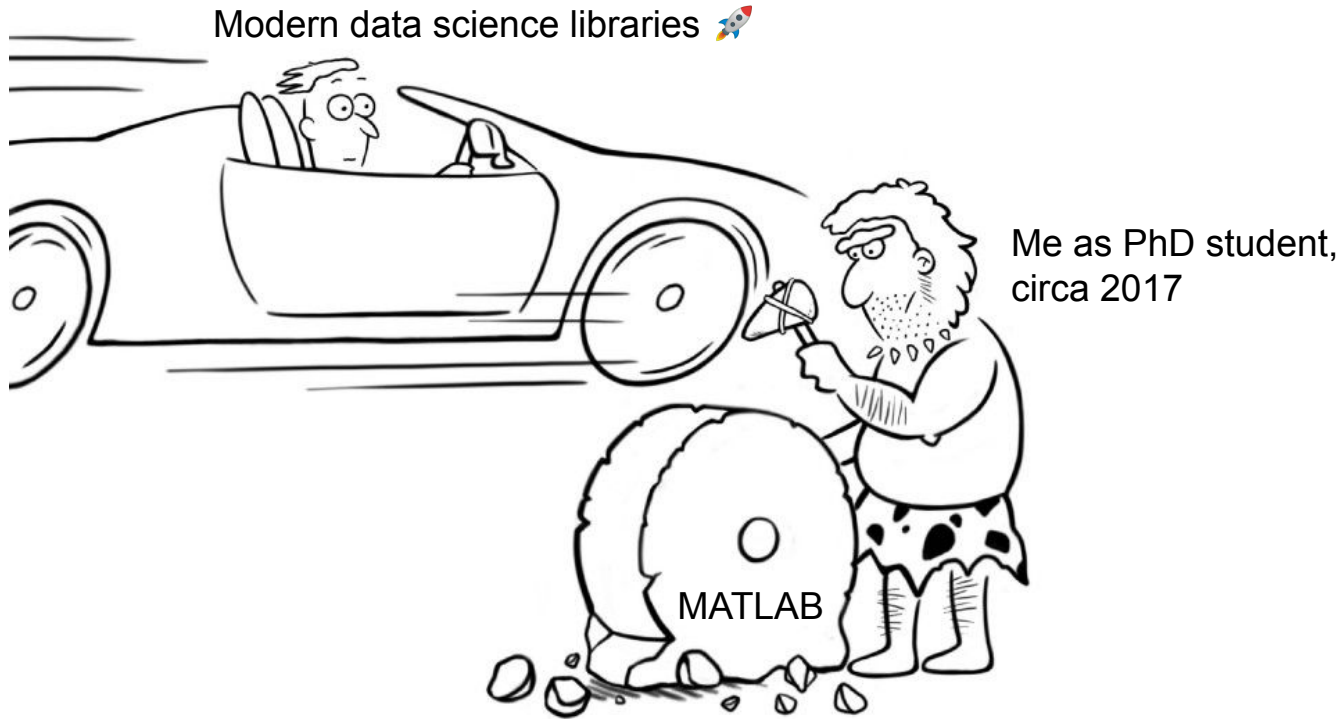


step 3: debug

**Because you likely  
rolled-your-own code...**

# Problem 1: Code not reused

# Problem 1: Code not reused



# Problem 2: Data accessibility

# Problem 2: Data accessibility

PRIVILEGED INSTITUTIONS CREATE  
"DATA FORTRESSES"\*



# Problem 2: Data accessibility

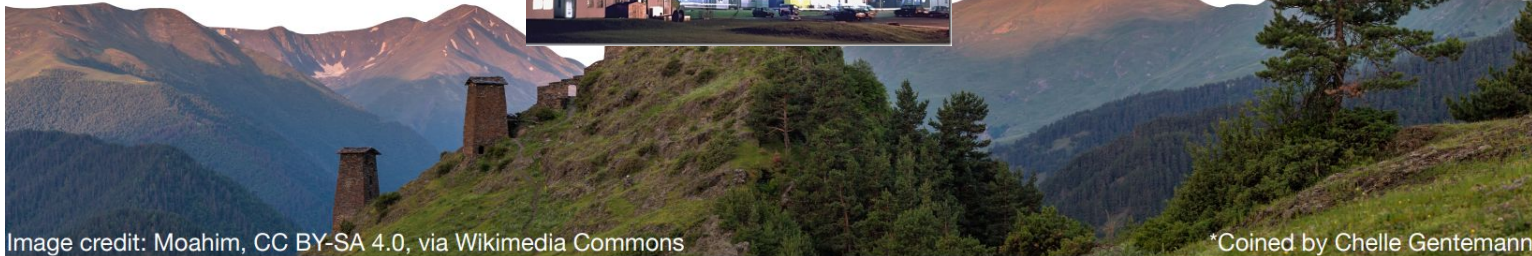
PRIVILEGED INSTITUTIONS CREATE  
"DATA FORTRESSES"\*





# Problem 2: Data accessibility

PRIVILEGED INSTITUTIONS CREATE  
"DATA FORTRESSES"\*



# Problem 3: Scale

**“Brb, let me just go download the data to my laptop...”**

# Problem 3: Scale

“Brb, let me just go download the data to my laptop...”

**MB**



# Problem 3: Scale

**“Brb, let me just go download the data to my laptop...”**

**GB**



# Problem 3: Scale

**“Brb, let me just go download the data to my laptop...”**

**TB**



# Problem 3: Scale

“Brb, let me just go download the data to my laptop...”

**PB**



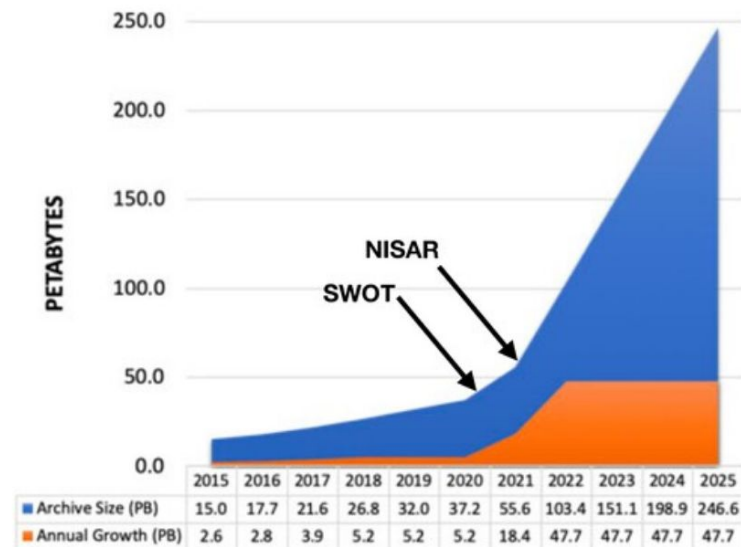
# Problem 3: Scale

## ITER NEWSLINE -

30 NOV, 2020

[Print](#) | [Read the latest published articles](#)

ITER Scientific Data Centre  
 HOW TO MANAGE 2 PETABYTES OF NEW  
 DATA EVERY DAY



Geoscientists' solution:



Geoscientists' solution:

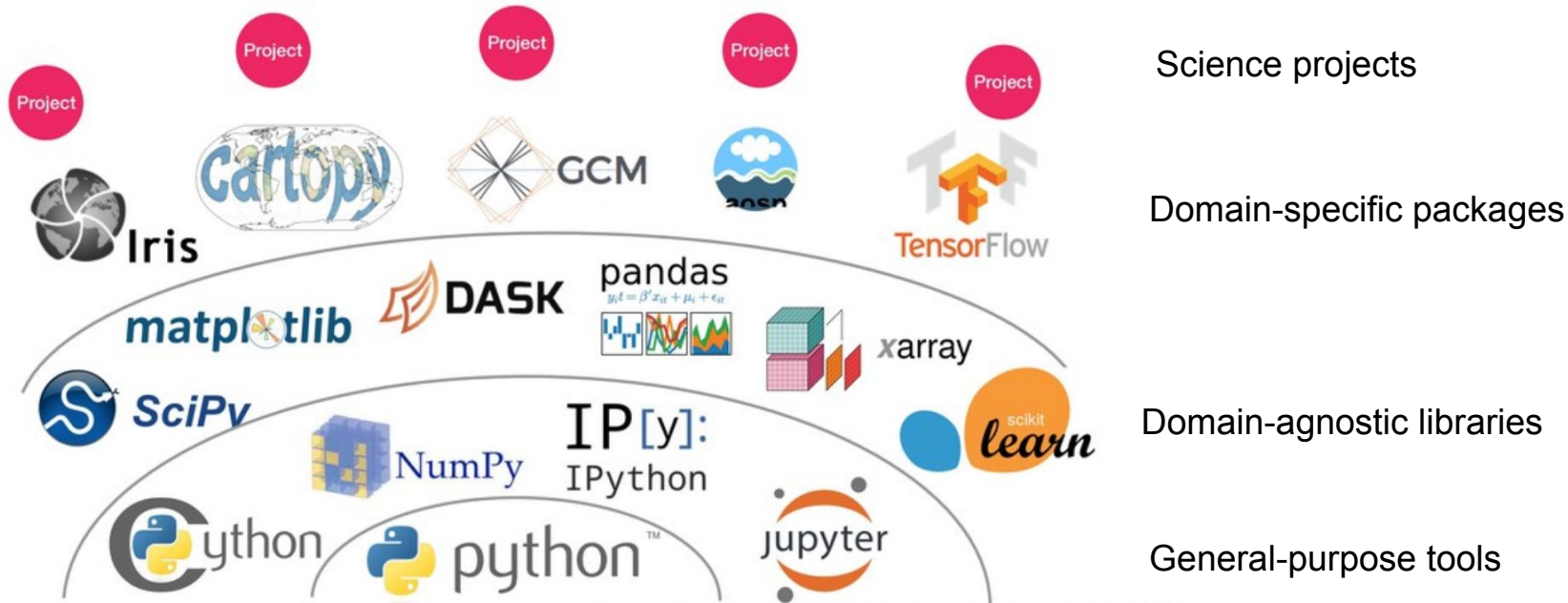


A community platform for Big Data geoscience

# Solution 1: Modular, open ecosystem

# Solution 1: Modular, open ecosystem

## ECOSYSTEM



Credit: Stephan Hoyer, Jake Vanderplas (SciPy 2015)

# Solution 1: Modular, open ecosystem

## PANGEO COMMUNITY

Lamont-Doherty Earth Observatory  
COLUMBIA UNIVERSITY | EARTH INSTITUTE



EARTH CUBI



NATIONAL CENTER FOR ATMOSPHERIC RESEARCH



Met Office



developmentSEED



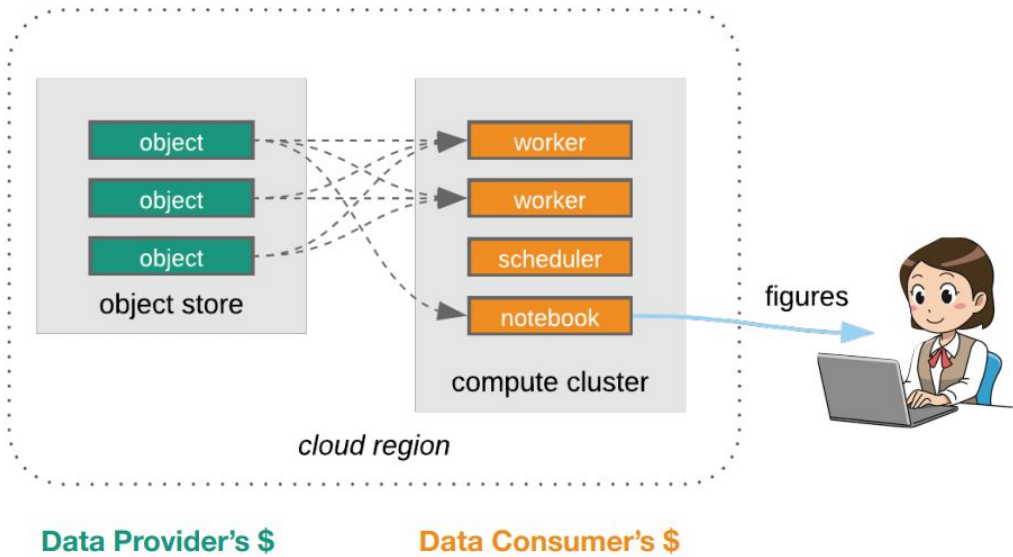
RHODIUM GROUP



CLIMACELL  
Weather. Revealed.

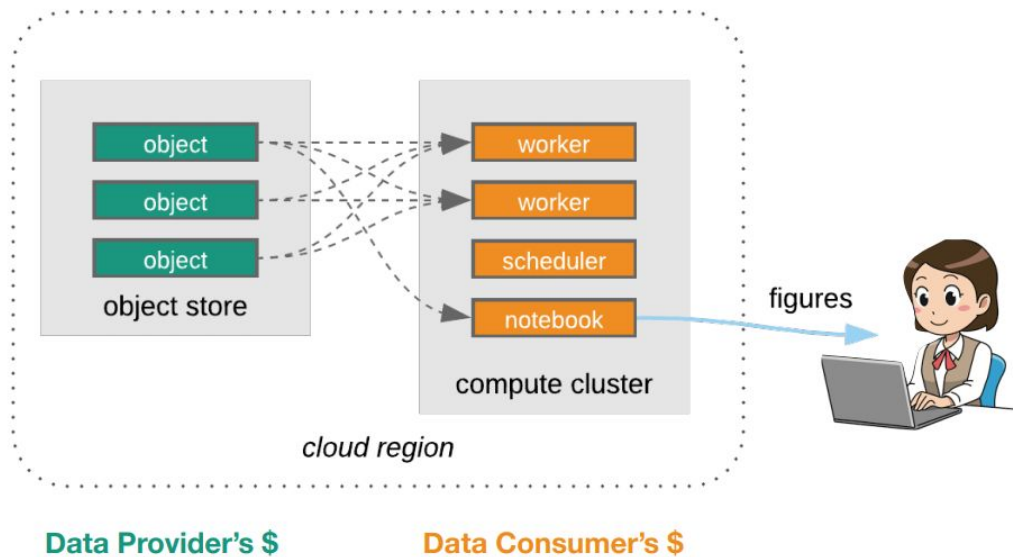
[HTTP://PANGEO.IO](http://pangeo.io)

# Solution 2: Cloud Computing



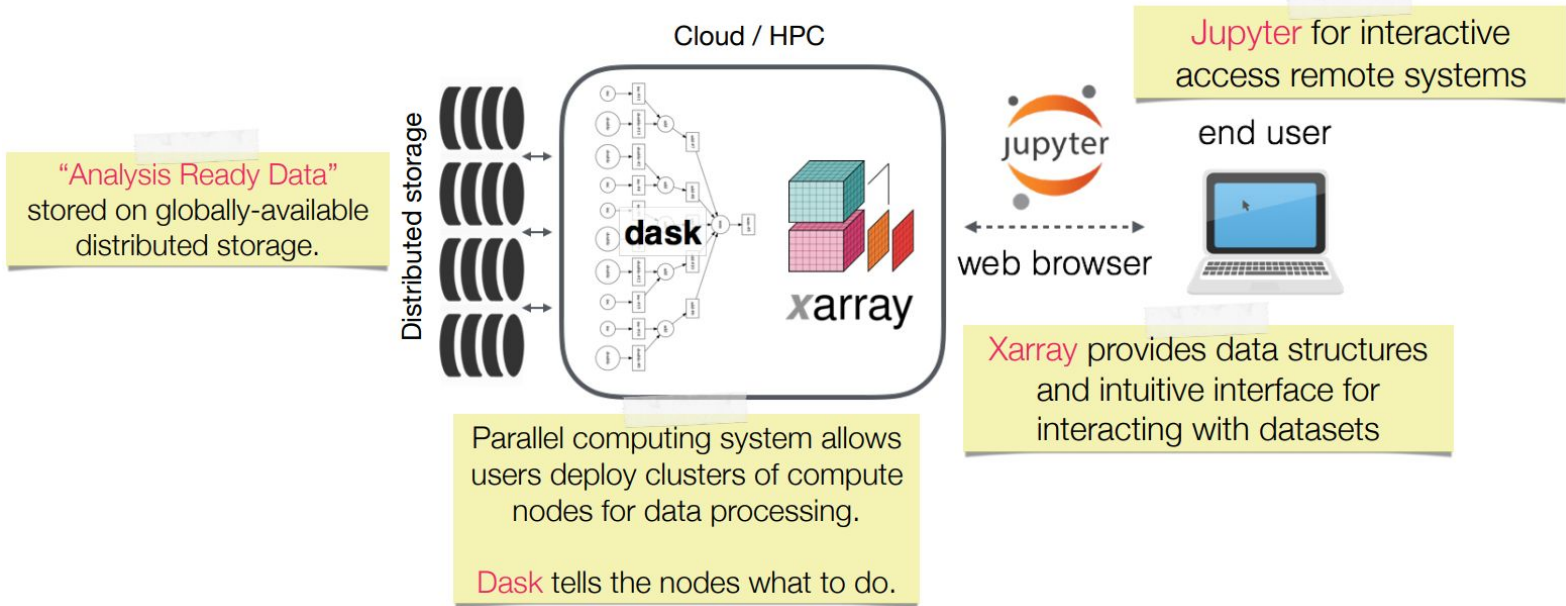
# Solution 2: Cloud Computing

Analysis Ready Data  
Cloud Optimized Formats



# Solution 2: Cloud Computing

## PANGEO ARCHITECTURE



# Solution 2: Cloud Computing

## PANGEO DEPLOYMENTS



[OCEAN.PANGEO.IO](https://ocean.pangeo.io)



Google Cloud Platform

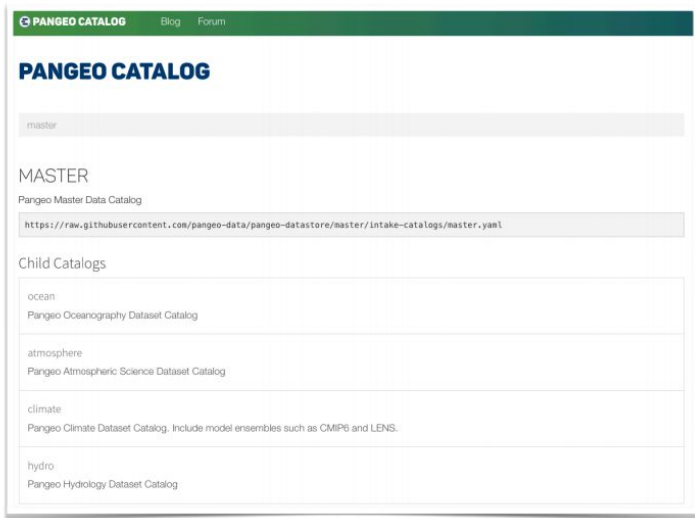




# Solution 2: Cloud Computing

## PANGEO CLOUD DATA CATALOG

[CATALOG.PANGEO.IO](https://catalog.pangeo.io)



**PANGEO CATALOG** Blog Forum

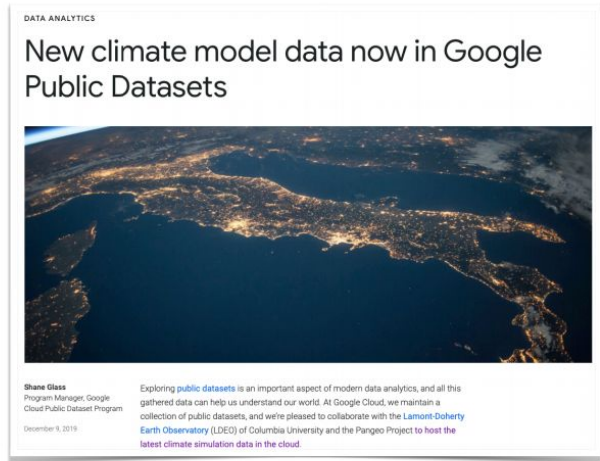
### PANGEO CATALOG

master

**MASTER**  
Pangeo Master Data Catalog  
<https://raw.githubusercontent.com/pangeo-data/pangeo-datastore/master/intake-catalogs/master.yaml>


Child Catalogs

- ocean  
Pangeo Oceanography Dataset Catalog
- atmosphere  
Pangeo Atmospheric Science Dataset Catalog
- climate  
Pangeo Climate Dataset Catalog. Include model ensembles such as CMIP6 and LENS.
- hydro  
Pangeo Hydrology Dataset Catalog



DATA ANALYTICS

### New climate model data now in Google Public Datasets

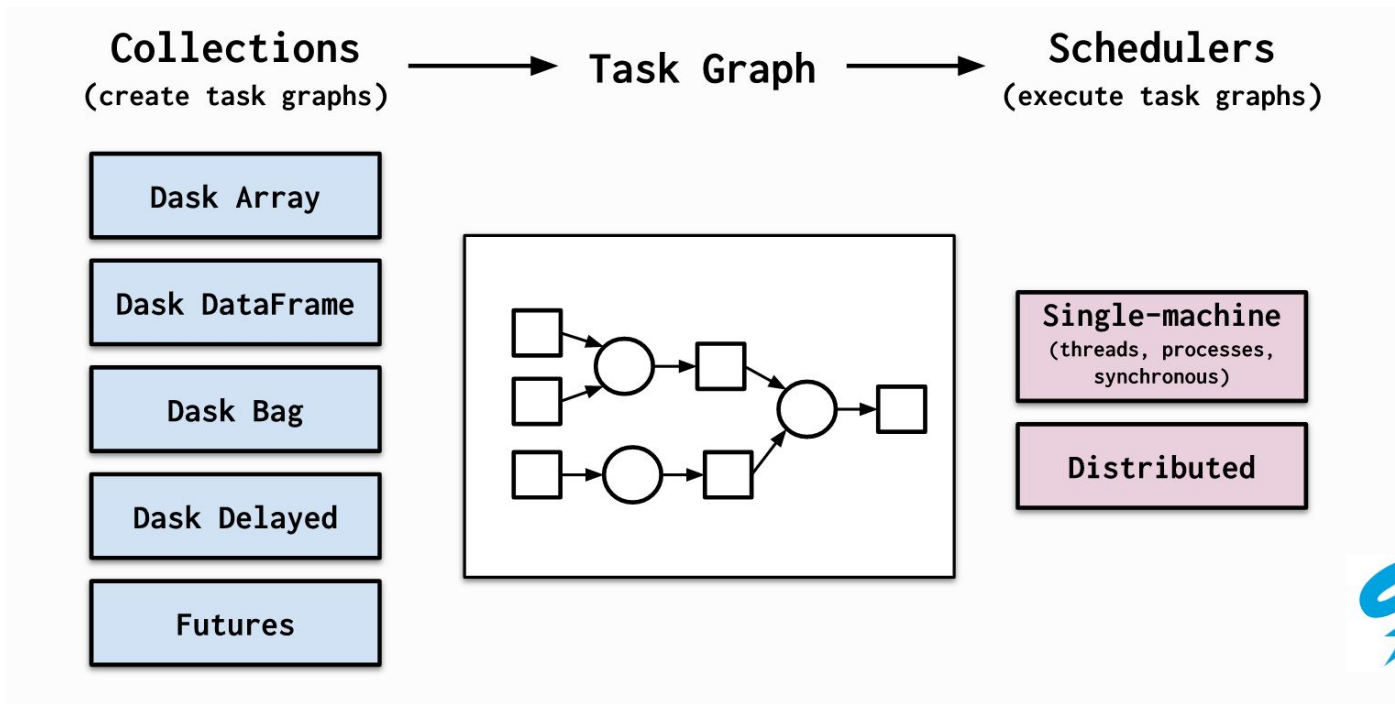


**Shane Glass**  
Program Manager, Google  
Cloud Public Dataset Program

Exploring **public datasets** is an important aspect of modern data analytics, and all this gathered data can help us understand our world. At Google Cloud, we maintain a collection of public datasets, and we're pleased to collaborate with the **Lamont-Doherty Earth Observatory (LDEO)** of Columbia University and the Pangeo Project to host the **latest climate simulation data in the cloud.**

December 4, 2019

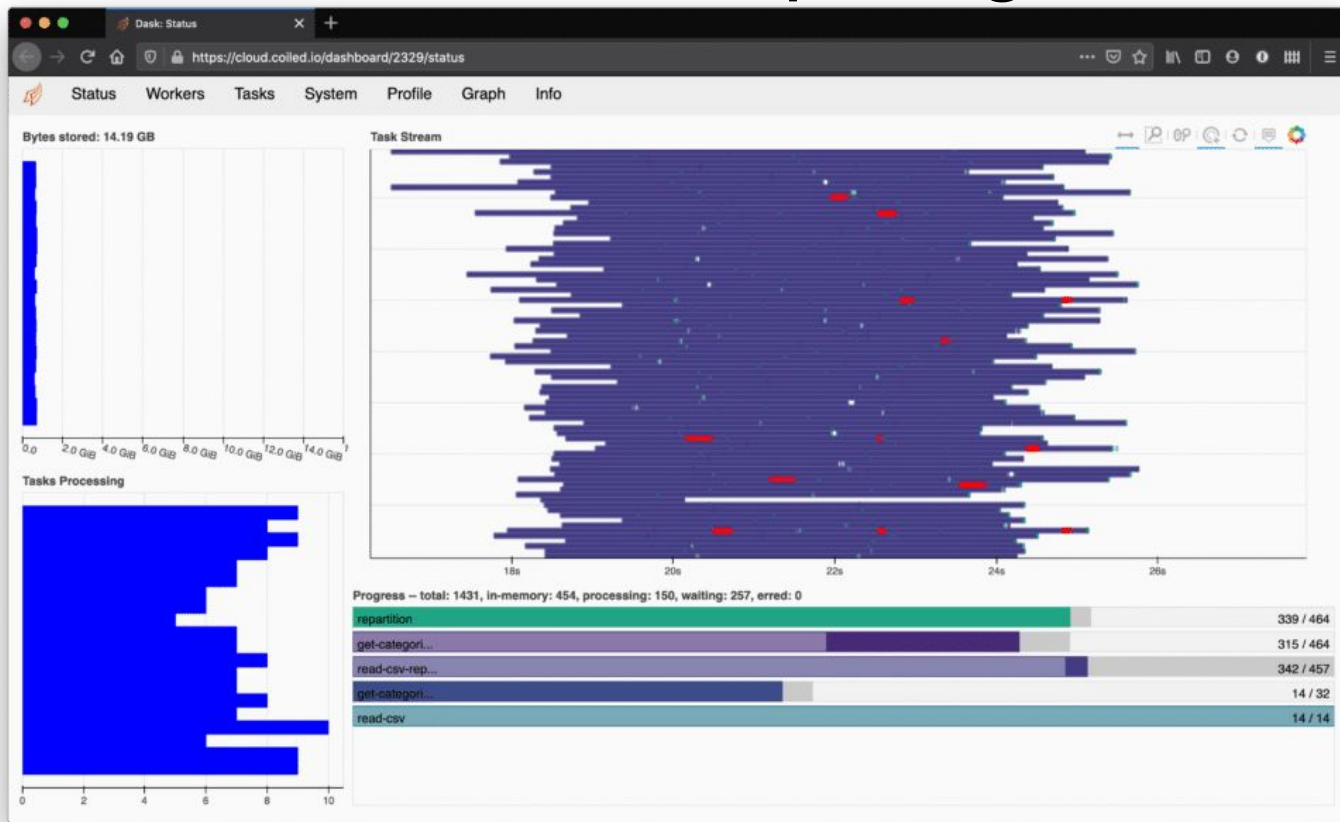
# Solution 3: Parallel computing frameworks



**DASK**



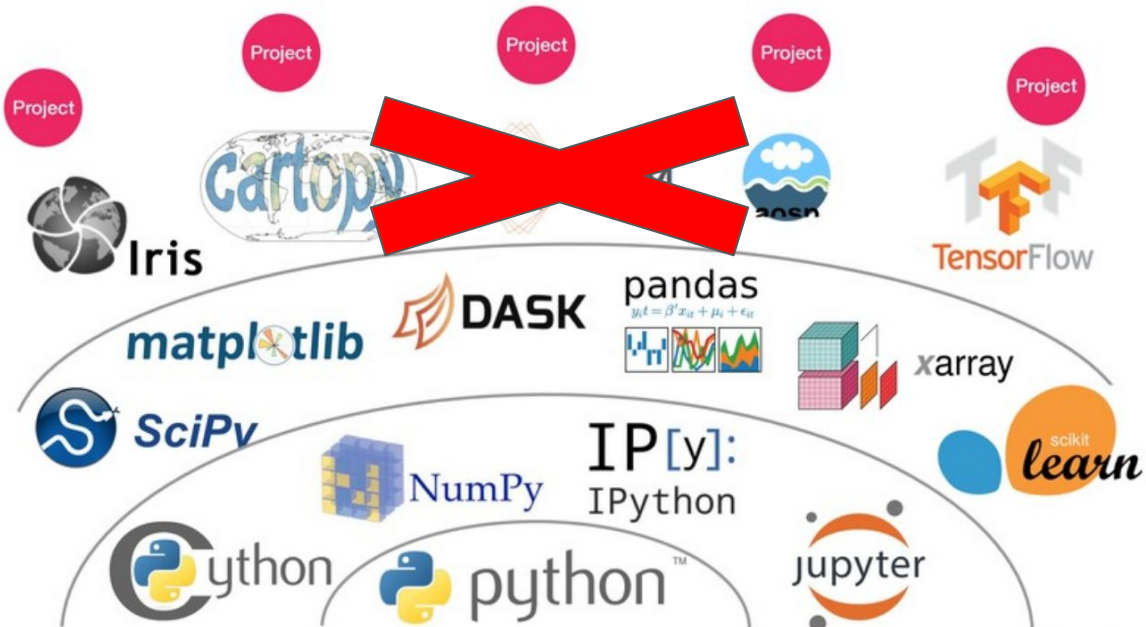
# Solution 3: Parallel computing frameworks



How might this work for plasma?

# How might this work for plasma?

## ECOSYSTEM



Fusion plasma projects

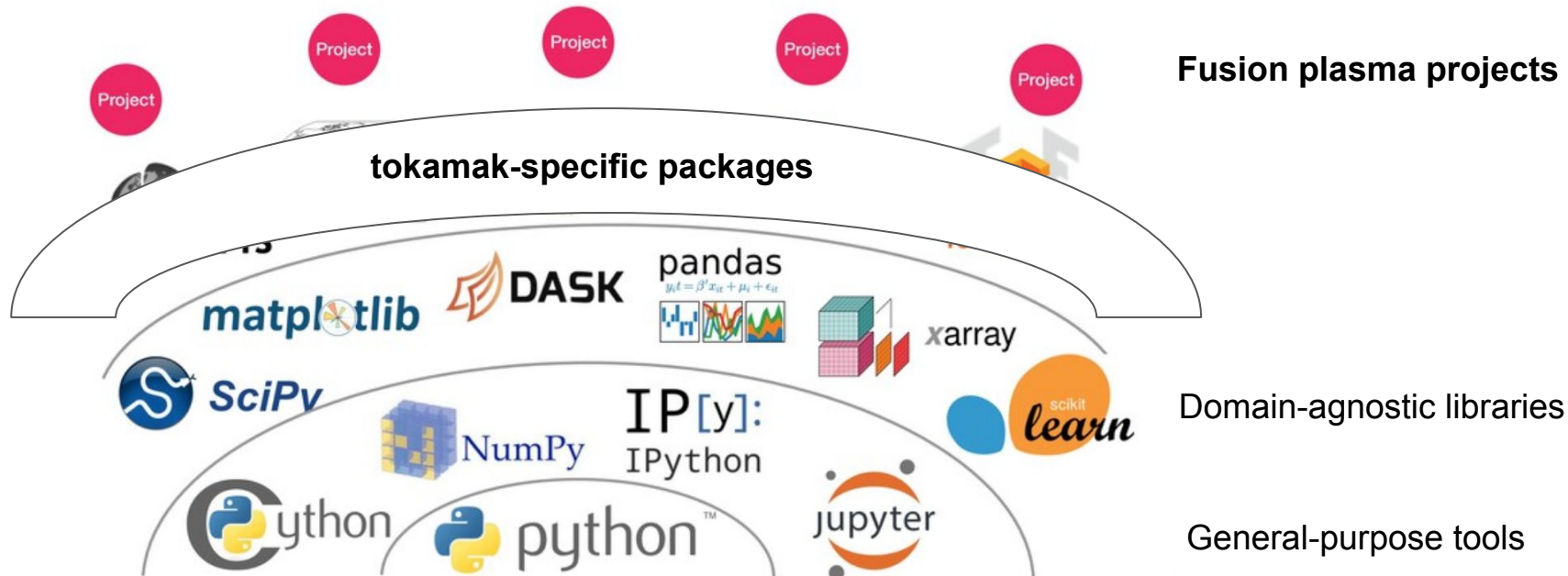
Domain-agnostic libraries

General-purpose tools

Credit: Stephan Hoyer, Jake Vanderplas (SciPy 2015)

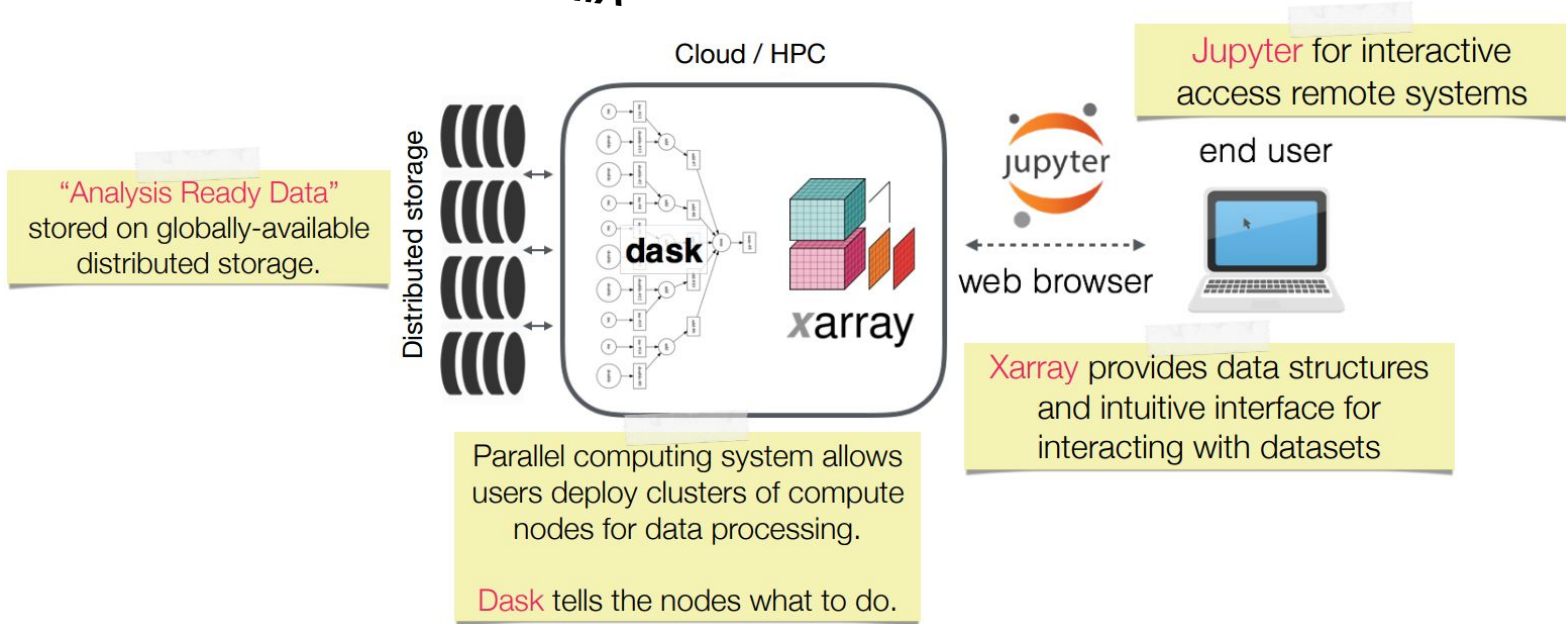
# How might this work for plasma?

## ECOSYSTEM

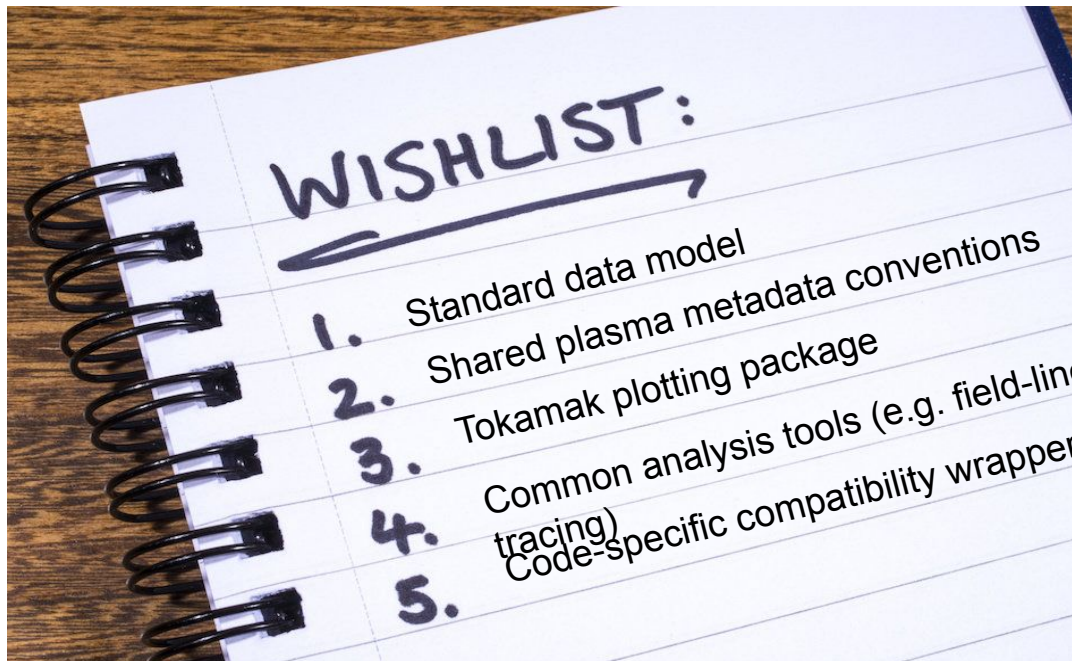


# How might this work for plasma?

## PANGEO ~~ARCHITECTURE~~ -PLASMA



# How might this work for plasma?



Blog post: <https://hackmd.io/@TomNicholas/rkyERwcoO#>



# Other bonuses of joining this ecosystem

- Parallel and out-of-core analysis
- Labelled dimensions
- Unit-aware arithmetic
- Easier reproducibility
- Plotting flexibility
- Machine Learning integration



# Summary

- Geoscience has same problems as plasma physics 🌍🤝☀️
- Being solved using:
  - Modular community software ecosystem 🛠️
  - Cloud computing ☁️
  - Parallel execution frameworks 🚀
- It's working for them - it could work for us! 🔬

# LEARN MORE



<http://pangeo.io>



<https://github.com/pangeo-data/>



<https://medium.com/pangeo>



[@pangeo\\_data](https://twitter.com/pangeo_data)