



PANGEO

Enabling Open, Reproducible, and scalable big data geoscience for everyone

Anne Fouilloux, *Simula Research Laboratory, Oslo, Norway*

Tina Erica Odaka, *LOPS UMR 6523, CNRS-IFREMER-IRD-Univ.Brest-IUEM, Brest, France*

Benjamin Ragan-Kelley, *Simula Research Laboratory, Oslo, Norway*

Francesco Nattino, *Netherlands eScience Center, Netherlands*

Meiert W. Grootes, *Netherlands eScience Center, Netherlands*

Ou Ku, *Netherlands eScience Center, Netherlands*

Raymond Oonk, *SURF, Netherlands*

What is Pangeo?

A global community initiative for Big Geoscience Data

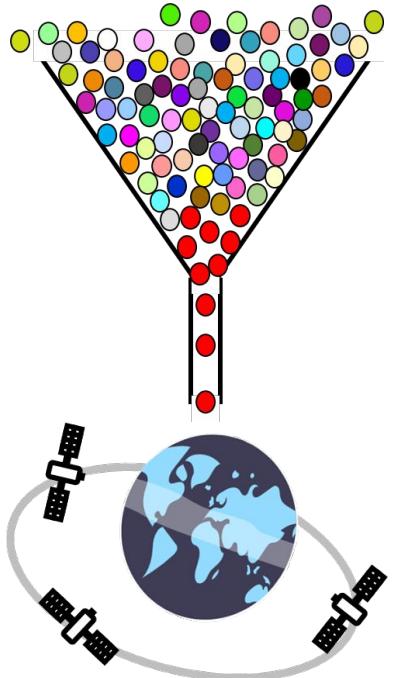


- Based on
 - Open Community
 - Open Platform
 - Open deployments
- Try it yourself!

Open Community

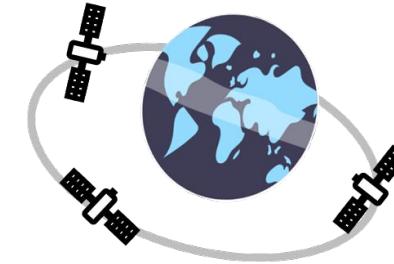
The Pangeo community

Vision: more equitable access to data, software and hardware to facilitate innovation



Earth Observation market

Open Earth Observation data



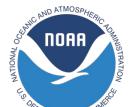
Facilitate take-up of the technology in the real world



Facilitate take-up of the technology in the real world

A culture of collaboration

Lamont-Doherty Earth Observatory
COLUMBIA UNIVERSITY | EARTH INSTITUTE



simula



National Centre for
Atmospheric Science
NATIONAL ENVIRONMENT RESEARCH COUNCIL



Inria



UNIVERSITY OF LEEDS



UNIVERSITY OF TORONTO



UNIVERSITY OF SASKATCHEWAN



The Alan Turing Institute



Massachusetts Institute of Technology ®



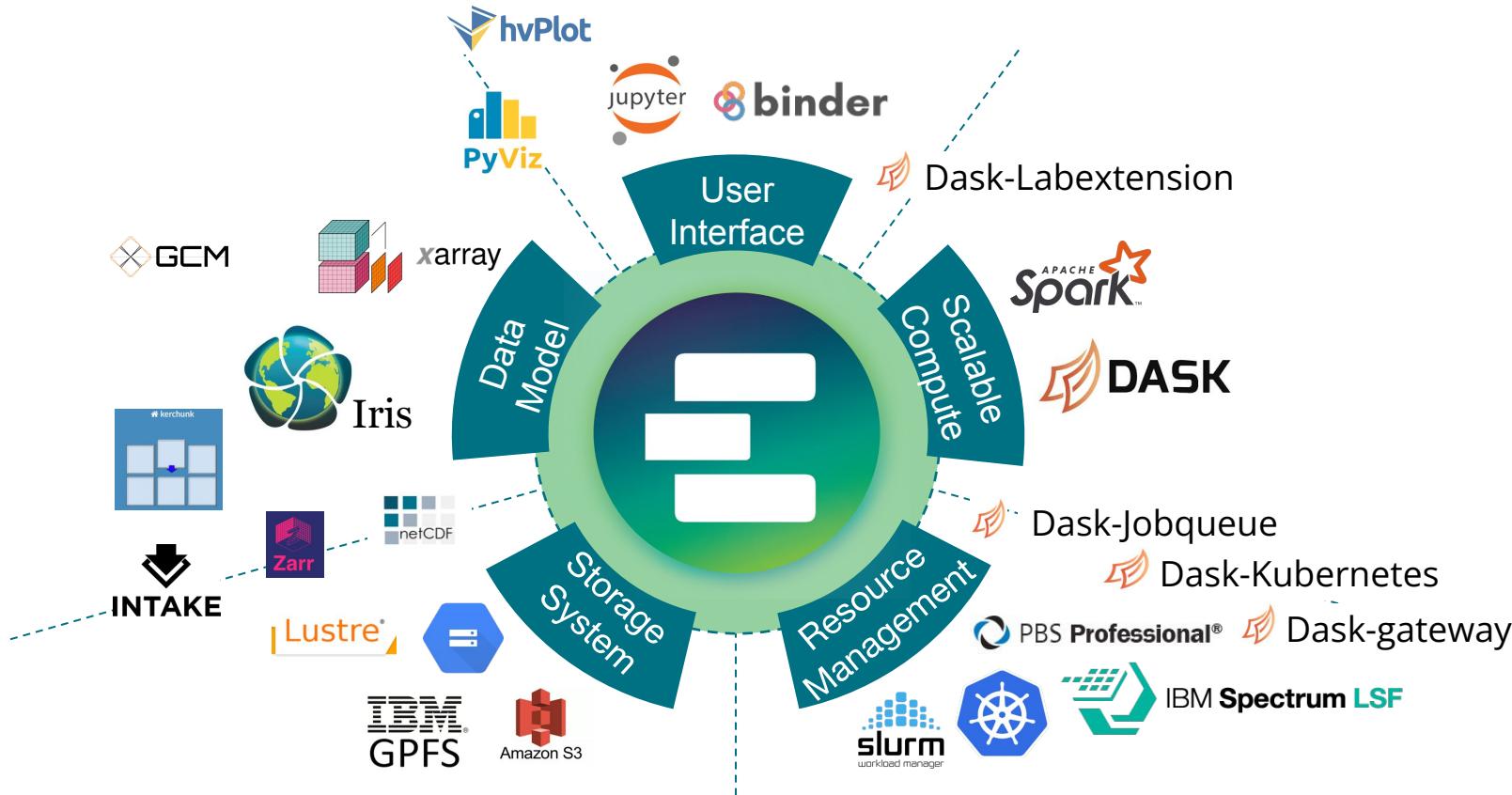
MEMORIAL UNIVERSITY



Based on GitHub and papers affiliations

Open Platform

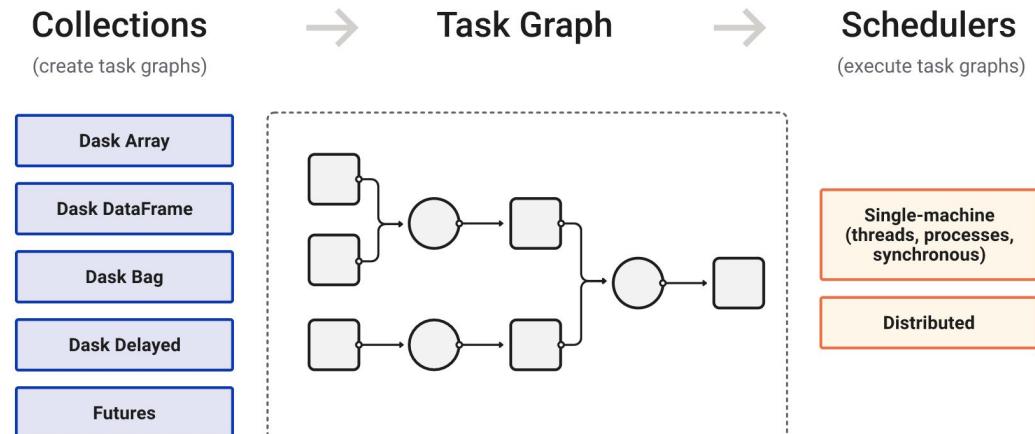
Modular, open and scalable platform



from laptop to cloud and HPC

Flexible library for parallel and distributed computing on "big (larger than memory) data".

- Dynamic task scheduling optimized for computation
- Extends pandas and NumPy APIs to larger-than-memory data (chunks)
- Fully customizable
- Scales from laptop to cluster (and back)
- Designed for interactive use
- Native python
- Almost like writing serial code (some idiosyncrasies)

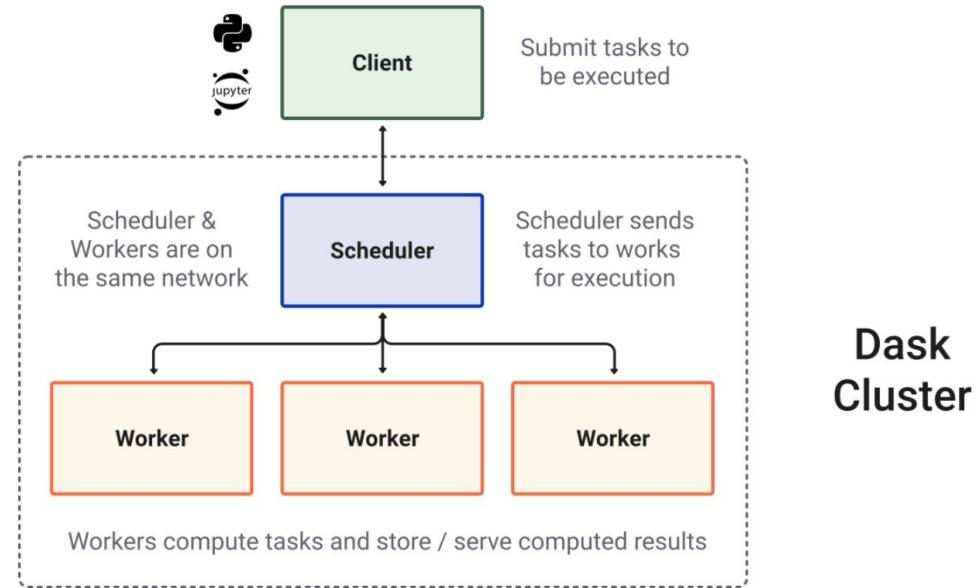


High level collections are used to generate task graphs which can be executed by schedulers on a single machine or a cluster.

Dask distributed scheduler

Supports fully asynchronous workflows and provides sophisticated data locality.

Seamlessly integrates with many/most HPC and cloud systems



Deploy anywhere
and for everyone

Pangeo deployments

- Laptop, e.g.
<https://github.com/pangeo-data/pangeo-docker-images>
- Cloud
 - Public Clouds (EOSC), e.g. Pangeo@EOSC
<https://pangeo-data.github.io/pangeo-eosc/>
 - Commercial clouds (AWS, Microsoft, Google)
- EO Platforms
 - Formerly DIAS (CREODIAS ...), e.g.
<https://dataspace.copernicus.eu/>
- HPC
 - CNES, IFREMER, SURF, ..

Pangeo@EOSC: Why?

- Pangeo is **infrastructure-agnostic**, making it a great use case for EOSC (and other federated clouds).
- When you use Xarray, Dask and Jupyter Notebooks, you are essentially using Pangeo.
- Pangeo serves as a community and **facilitator** for deployment on various infrastructures.
- It is **open source**, eliminating vendor lock-in risks.
- At times, end-users are limited to using EU infrastructure.
- A **federation** of infrastructure **can't be** a single **centralized** infrastructure.

Pangeo eosc: How?



Custom environments



EOSC Community Hub



Online content



Reuse existing [2i2c Pangeo deployment configuration](#)



Check-in

Cloud infrastructure



EOSC Authentication for
students, researchers, data
scientists, etc.



Pangeo Training Infrastructure as a Service

FOSS4G workshop

Overview

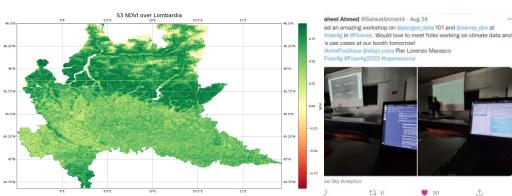


Questions

- Why do chunking matter?
- How can I read datasets by chunks to optimize memory usage?

Objectives

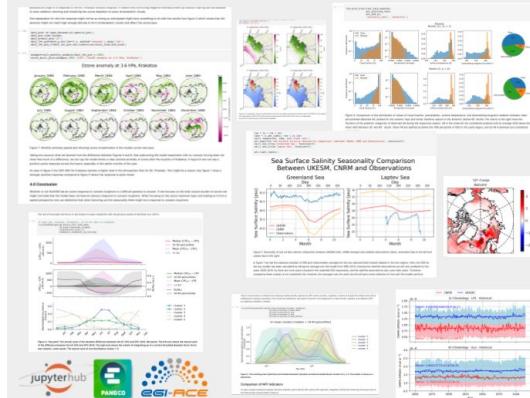
- Learn about chunking
- Learn about zarr
- Use kerchunk to consolidate chunk metadata and prepare single ensemble datasets for parallel computing



Pangeo 101 workshop at the FOSS4G conference



eScience Course



Integrate EOSC into master students' courses

CLIVAR Arctic Bootcamp

Supporting researchers in becoming competent practitioners in their scientific domain

Pangeo & OpenEO at BiDS



cesa

Request PTIaaS to train and teach!

Fill the form

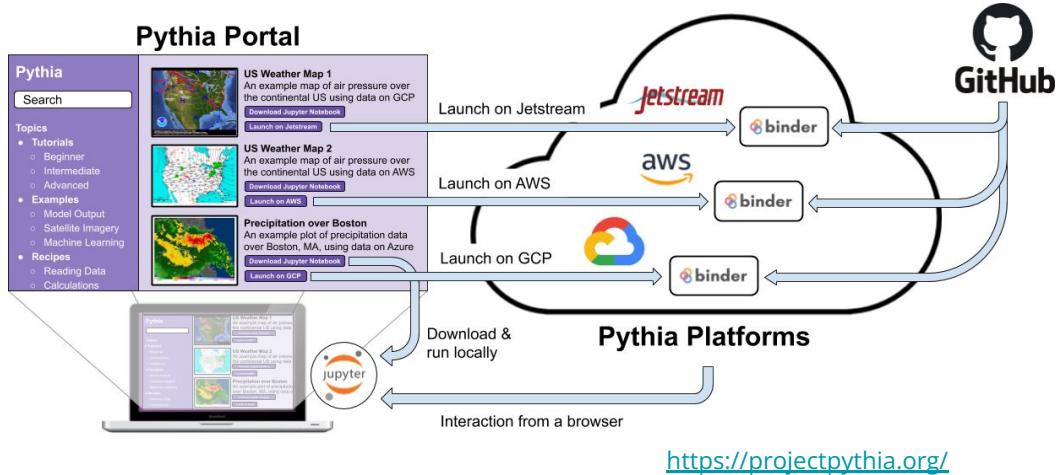


Team-up with other initiatives to
onboard new members and
increase diversity

Team-up with Project Pythia



A Community Learning Resource for Geoscientists



This work supported through the National Science Foundation Award #2026899.

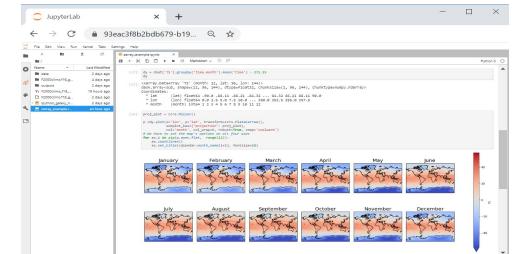
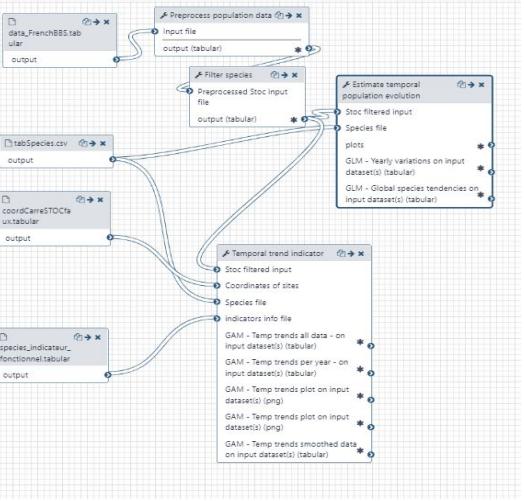
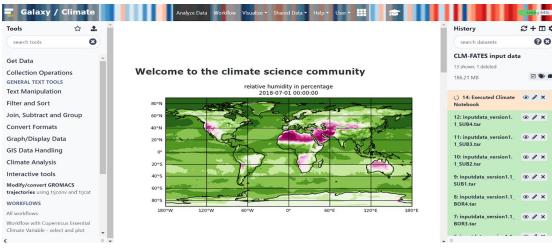
Aspiration goal: Be the goto resource for learning the *Scientific Python Ecosystem*

- ★ Geoscience focused
- ★ From beginner to the power user
- ★ Tutorials, videos, examples, on-line courses, and sample data
- ★ Community owned!

Team-up with Galaxy Europe

Galaxy is an open-source community and platform for FAIR data analysis. It offers:

- **Graphical User Interface (GUI) for users with no programming skills**
- Workflow editor to create and run fully reproducible data analysis
- Compute & Storage to everyone (free registration)
- [Self Paced Learning material](#) with the Galaxy Training Network
- [Training Infrastructure as a Service](#) (TiaaS) free and ready to use with private queues where only training's jobs run



Team-up with the Environmental Data Science Book

<https://edsbook.org/welcome.html>

Reproducible, scalable, & shareable
ENVIRONMENTAL DATA SCIENCE



Scriberia ↗

The
Alan Turing
Institute



Contribution



Collaborative
Reviewing

jupyter {book}



Publication



EXPLORE

CC-BY image by Turing Way and Scriberia, CC-BY 4.0, EDS Book community, @EnvDSBook

Living, open and community-driven online resource to **showcase and support the publication** of data, research and open-source tools.

Reliance

RoHub

The screenshot shows the RoHub interface for a FAIR Executable Research Object. It includes sections for Title, Tags (Environment, Theme), launch (binder), Context (purpose, highlight, contributions), Data, Analysis, and Citation. Logos for Reliance, RoHub, and jupyter are visible at the bottom right.

GeoPython community

- Geopandas
- GDAL
- Pysal

Digital open books, building on Jupyter notebooks etc.

- “Geocomputation with Python” (or R) by Michael Dorman, Anita Graser, Jakub Nowosad, Robin Lovelace: <https://py.geocompx.org/>
- “Geographic Data Science with Python”, by Sergio J. Rey, Dani Arribas-Bel and Levi J. Wolf: <https://geographicdata.science>

→ JupyterGIS: <https://blog.jupyter.org/jupytergis-d63b7adf9d0c> (new ESA Digital Innovation project)

Pangeo@SURF: scaling on HPC/HTC

Pangeo@SURF

RS-DAT Components

<https://github.com/RS-DAT/JupyterDaskOnSLURM>

JupyterDask-Examples (Public)
Collection of examples for RS-DAT JupyterDask deployments
• Jupyter Notebook ⭐ 0 Apache-2.0 0 0 0 Updated 3 hours ago

JupyterDaskOnSLURM (Public)
• Shell ⭐ 0 Apache-2.0 0 0 0 Updated 11 hours ago

JupyterDaskOnSRC (Public)
Deploy JupyterHub and Dask on SURF Research Cloud
• Jinja ⭐ 0 Apache-2.0 0 4 0 Updated 15 hours ago

Scalable analysis with Jupyter

Using legacy (Docker) containers for HPC

DockerToSingularity (Public)
Examples of converting a docker image to singularity, and execute the singularity image
• Dockerfile ⭐ 0 0 0 0 Updated on 12 Aug

<https://github.com/RS-DAT/DockerToSingularity>

<https://github.com/RS-DAT/dcachefs>

dcachefs (Public)
Python file-system interface for dCache
• Python

Python interface to dCache

<https://github.com/RS-DAT/stac2dcache>

stac2dcache (Public)
Python tool to create and manipulate STAC catalogs on a dCache storage system
• Jupyter Notebook

utility functions to manage STAC catalogs (and the underlying data) on dCache.

Jupyter Dask on Slurm

SURF H*C Infrastructure: Spider (HTC), Snellius (HPC); SLURM scheduler

JupyterDaskOnSLURM Public



Python 3 Apache-2.0 0 4 0 Updated on Feb 2

- JupyterLab instance with Dask and Git extensions, scalable Dask cluster, running on SLURM managed HTC/HPC system
- Basic idea
 - Launch JupyterLab server and Dask scheduler as long-running batch job
 - SSH port-forwarding to connect to JupyterLab server
 - Launch workers as short-lived batch jobs (fast thru queue)

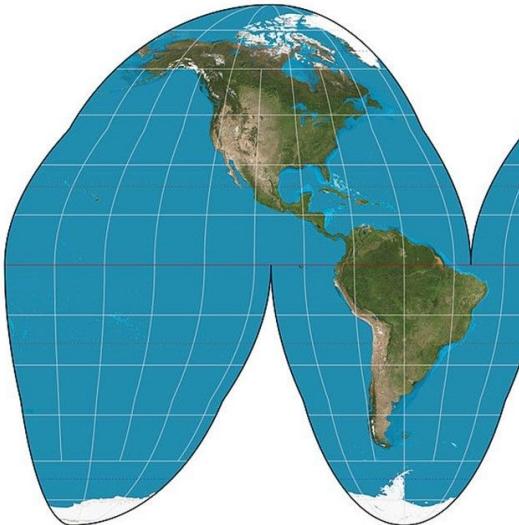
Ensure ease of use by abstracting details from user

A community of developers, scientists and users

America

Time zone +11-> -2

Every Wednesday, alternating
between 12pm ET and 4pm ET

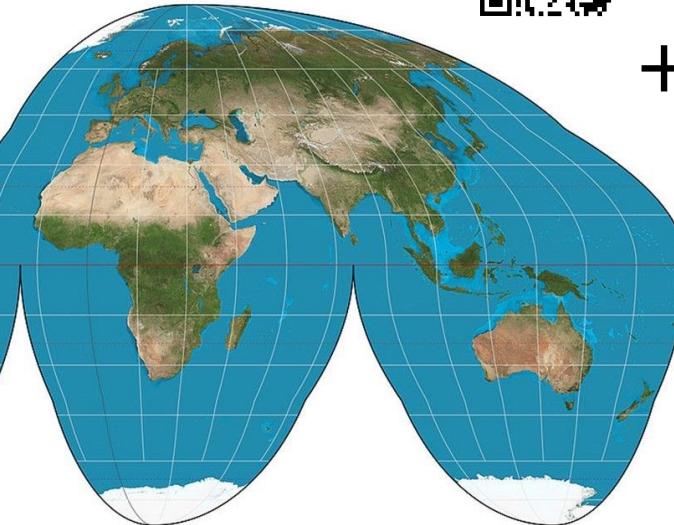


Europe, Africa, West Asia

Time zone -1-> +5

**Every Tuesday at 9.30 a.m.
CET/CEST here**

Join us here



Australia, East Asia

Time zone -1-> +5

3rd Friday of every month at 1pm
Australian Eastern Time



Check out for different discussions :

- [software-production-visualization](#)
- [data-management](#)
- [software-geospatial-bridge](#)
- [cloud-partner-managed-infrastructure](#)
- [cloud-pangeo-managed-infrastructure](#)
- [democratization-science-community](#)
- [scientific-publishing](#)
- [software-massive-scale](#)
- [software-regridding](#)

More info & links here - <https://pangeo.io/meeting-notes.html>

Pangeo: where can I find information?



- <http://pangeo.io/index.html#what-is-pangeo>
 - (<https://github.com/pangeo-data/pangeo>)
- Try some examples in the Pangeo Gallery
 - <http://gallery.pangeo.io>
- Subscribe in discourse, start discussions
 - <https://discourse.pangeo.io/>
- Follow pangeo-data from Github
- Discord channel
 - <https://discord.gg/ex5qqEyyTz>
- Interesting informations exist also in blogs
 - <https://medium.com/pangeo>
- Code of conduct.
 - https://github.com/pangeo-data/governance/blob/master/conduct/code_of_conduct.md

Hands on

<https://pangeo-data.github.io/geo-open-hack-2024/intro.html>

Thank you for your attention!

