



Open Innovation for Earth Observation Programmes

2-4 November 2022 | ESA-ESRIN | Frascati (Rm), Italy

Perspective on creating Open Source user workflows
The Pangeo Community



Outline

- What is Pangeo?
- How does Pangeo foster Open Innovation?
- Successes, Challenges and Lessons learned
- Recommendations to Space Agencies

What is Pangeo?



A global community initiative for Big Geoscience Data that promotes open, reproducible, and scalable science

- Open Community
- Open Platform

Code of conduct

- https://github.com/pangeo-data/governance/blob/master/conduct/code_of_conduct.md

Governance

- <https://github.com/pangeo-data/governance/blob/master/governance.md>



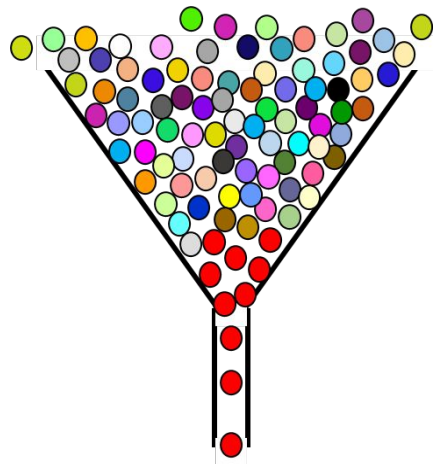
Open Community

Open Innovation for Earth Observation Programmes

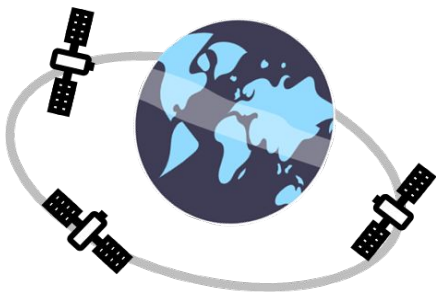
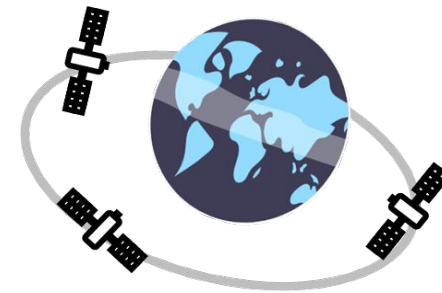
2-4 November 2022 | ESA-ESRIN | Frascati (Rm), Italy



Pangeo vision: lower barriers to entry to stimulate innovation



Facilitate contributions to increase diversity



DEI (Diversity, Equity and Inclusion) is essential

A community of developers, scientists and users

America
Time zone +11-> -2

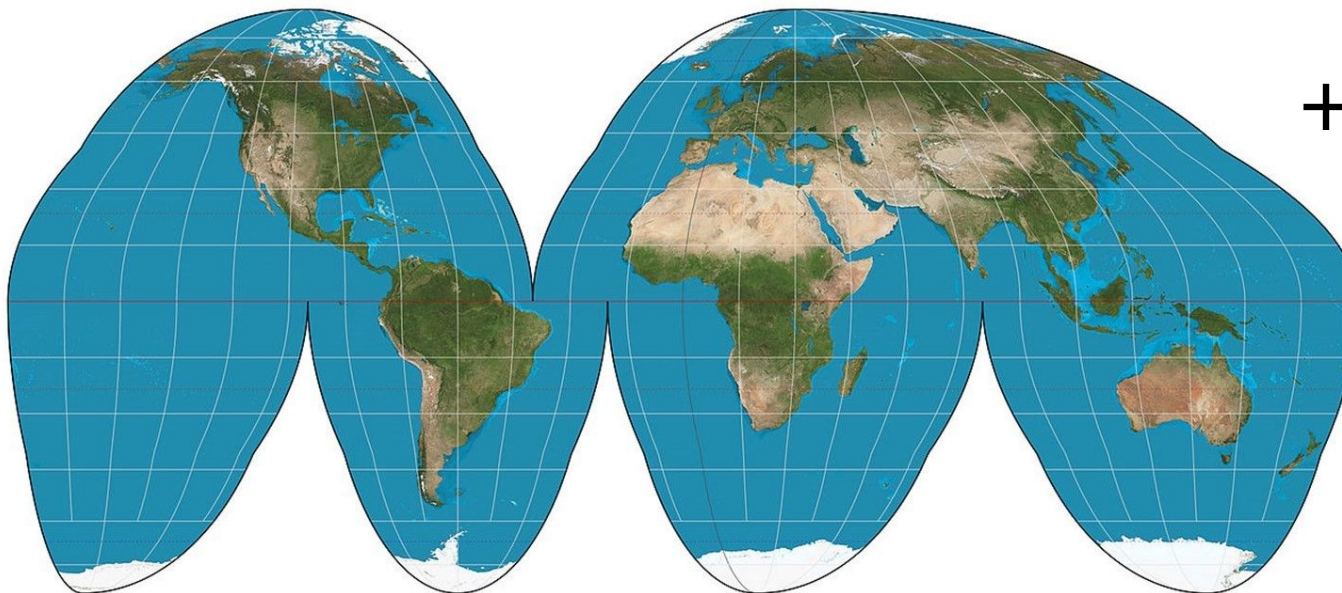
Europe, Africa, West Asia
Time zone -1-> +5

Australia, East Asia
Time zone -1-> +5

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

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

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



+ Working Group Meetings

- *Machine Learning Working Group*
- *Cloud Operations Working Group*
- *Project Pythia (formerly the Education Working Group)*
- *Pangeo Forge (Cloud Data Platform)*
- *Open Science Meeting - discussions to coordinate open science activities*

Pangeo Showcase Webinar Series

- To learn from each other. It is part of “America” community calls. 15 minute talks are recorded, given a DOI and made available on the [Pangeo YouTube Channel](#);

FALL 2022 SHOWCASE

Date	Speaker	Title
2022-09-21 12PM EDT	Peter Marsh, University of Cape Town	Accessing NetCDF and GRIB file collections as cloud-native virtual datasets using Kerchunk DOI 10.5281/zenodo.7140825
2022-10-05 12PM EDT	Leah Wasser, pyOpenSci	PyOpenSci DOI 10.5281/zenodo.7158586
2022-10-12 4PM EDT	Rich Signell, USGS	My ERA5 Journey: From API-to-ARCO DOI 10.5281/zenodo.7226344
2022-10-19 12PM EDT	Matthias Mohr, openEO	openEO: What it is and how it relates to Pangeo DOI 10.5281/zenodo.7229398
2022-10-26 4PM EDT	Hauke Schulz, CICOES/University of Washington	Xbitinfo: Compress datasets based on their information content

Show & Tell

- At the request of anyone who is interested to **go in depth and show what they are working on**. It is expected to be hands-on and is also recorded.
- FAIR Digital Object: a Research Object (<https://reliance.rohub.org/>) is created and aggregate all the material used/showed (slides, recorded video, jupyter notebooks or other codes, input data, etc.) and given a DOI.
- Text mining service to extract metadata and facilitate interdisciplinary re-usage

DOI: 10.24424/TG01-KV33

Created: 04.10.2022 (11:22), last modified: 31.10.2022 (10:31)

OPEN MANUAL SNAPSHOT EXECUTABLE RESEARCH OBJECT PANGEO

APPLIED SCIENCES EARTH OBSERVATION EARTH SCIENCES

DGGS and their potential impact in Geoscience and Geospatial communities

Alexander Kmoch

Contributed by Pangeo Europe
Published by Simula Research Laboratory

Overview Content Completeness Activity Life cycle Relations Impact

A Discrete Global Grid Systems (DGGS) is a unique type of spatial reference system comprising of a hierarchy of uniquely identifiable discrete grid cells that span the globe at multiple resolutions. A DGGS can support efficient management, storage, integration, exploration, mining, and visualisation of large geospatial datasets, and several systems of tessellation and indexing schemes exist. The main topic of this session is to introduce the audience to the theoretical background of Discrete Global Grid Systems (DGGS), current real-world implementations and exemplary use cases. This includes grid generation,...

☆ 0.00 / 5 0 0

9 Downloads 3 Views

Hide more details

Resources	9
Annotations	54
Events	86
Forks	0
Snapshots	0
Archives	0
Size	10064.65 KB

AGENTS

Pangeo Europe
Creator

COMPLETENESS 100%

DISCOVERED METADATA: 0

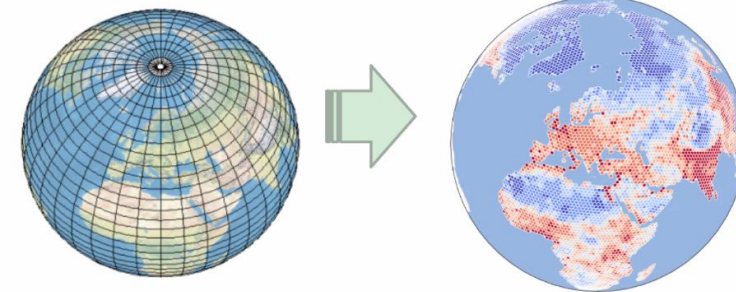
CARTOGRAPHY DATABASE
EARTH SCIENCES ATMOSPHERIC SCIENCE
MEDICAL PROCEDURE-TEST
MATHEMATICAL AND COMPUTER SCIENCE
COMPUTER OPERATIONS AND HARDWARE

TOOLBOX

SHARE

CITE AS

Kmoch, Alexander, and Pangeo Europe.



LOCATION: [dropdown]

CONTENT

- biblio
 - Pangeo discourse on "Discrete Global Grid Systems (DGGS) use..."
 - DGGS and their potential impact in Geoscience and Geospatial... (9025Kb)
 - Pangeo discourse on "Discrete Global Grid Systems (DGGS) use..."
 - DGGS and their potential impact in Geoscience and Geospatial... (9025Kb)
- input
 - Pangeo discourse announcement Show & Tell on "October 6, 20..."

Open Innovation for Earth Observation Programmes

2-4 November 2022 | ESA-ESRIN | Frascati (Rm), Italy



A Culture of Collaboration



Lamont-Doherty Earth Observatory
COLUMBIA UNIVERSITY | EARTH INSTITUTE



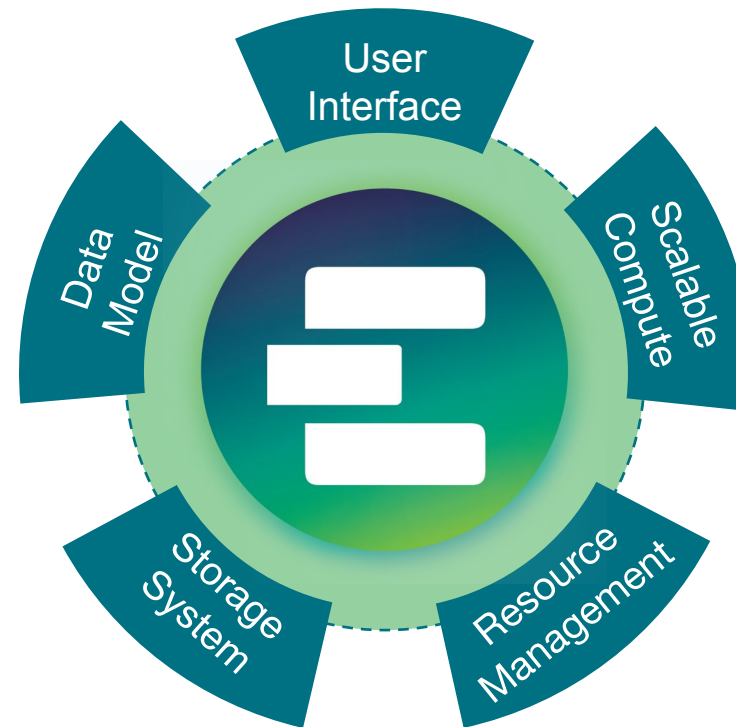
UNIVERSITY OF LEEDS



Based on GitHub and papers affiliations

Open Platform

Pangeo **platform** *is scalable*



from laptop to cloud or HPC

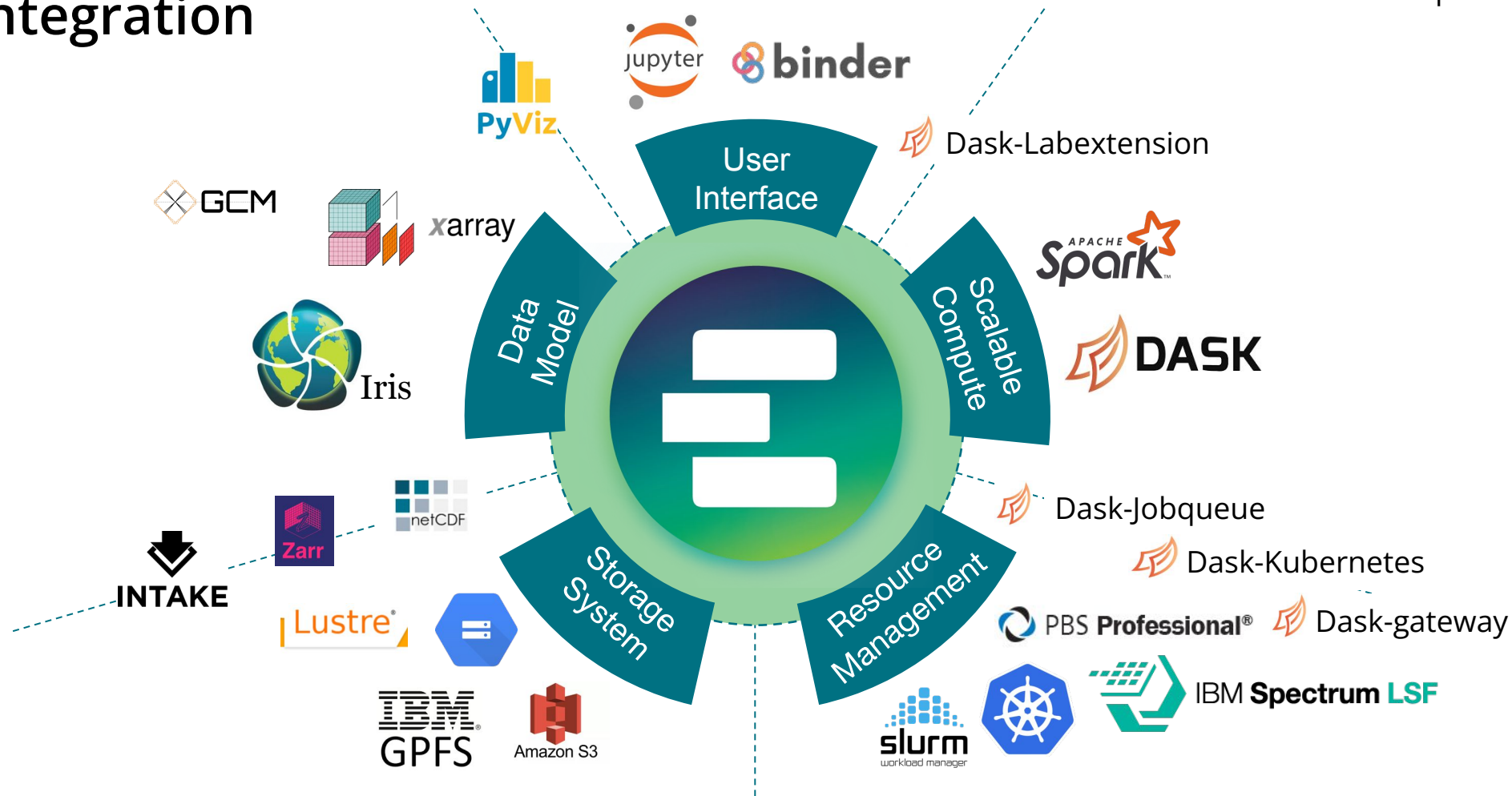
Open Innovation for Earth Observation Programmes

2-4 November 2022 | ESA-ESRIN | Frascati (Rm), Italy



Pangeo is about integration

Conda integrates most platforms



Pangeo makes it possible to explore geoscience data using HPC or cloud in an interactive manner

Pangeo deployments

- HPC
 - CNES, IFREMER, PRACE, Fugaku, ..
- Cloud
 - Public Clouds (EGI-ACE, EOSC)
 - Commercial clouds (AWS, Microsoft, Google)
- Laptop
 - Anywhere for anyone!

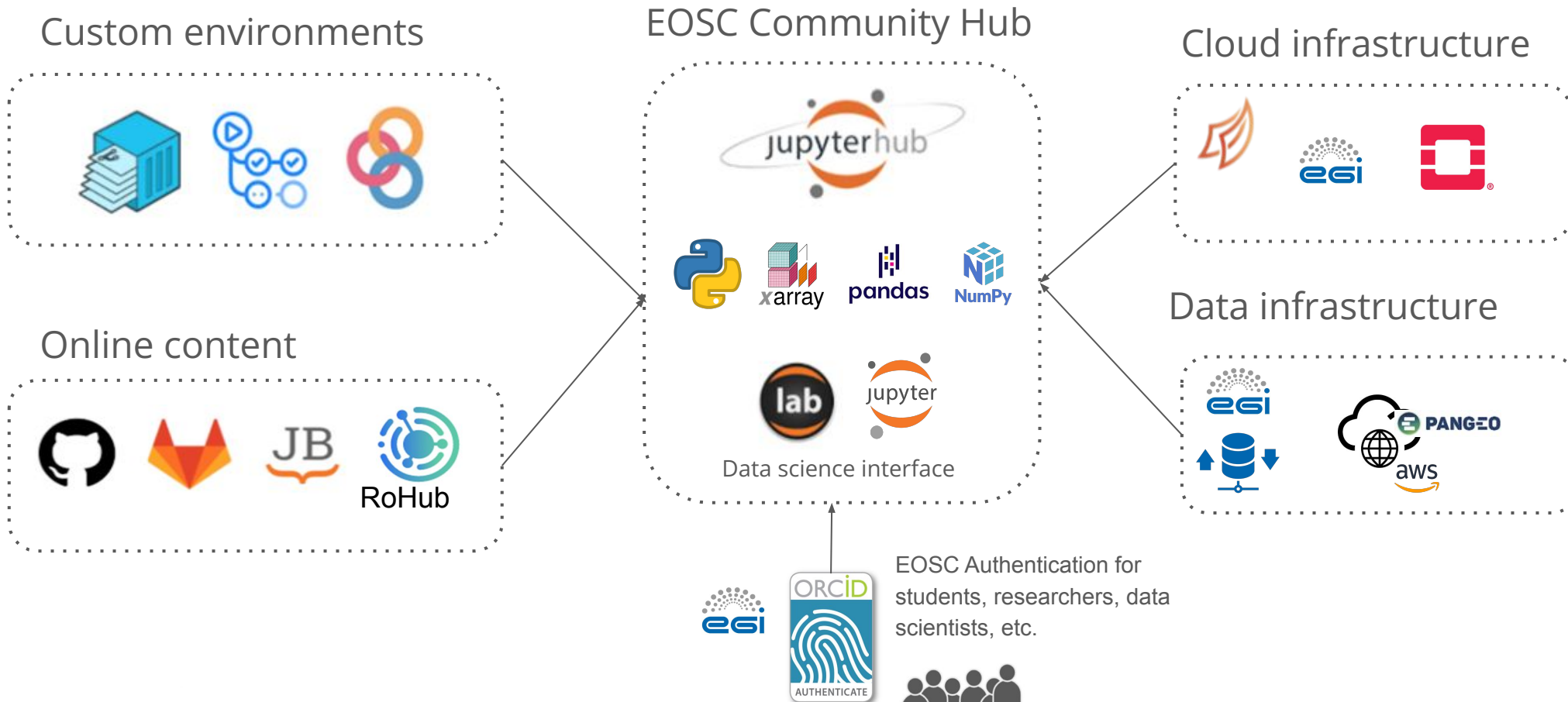
Docker images (pangeo-notebook, ml-notebook, etc.) developed and maintained by the Community to ease usage and deployments.

Open Innovation for Earth Observation Programmes

2-4 November 2022 | ESA-ESRIN | Frascati (Rm), Italy



Pangeo at scale with deployment on the European Open Science Cloud (EOSC)



Mirrored from Zi2c Pangeo deployment (US)



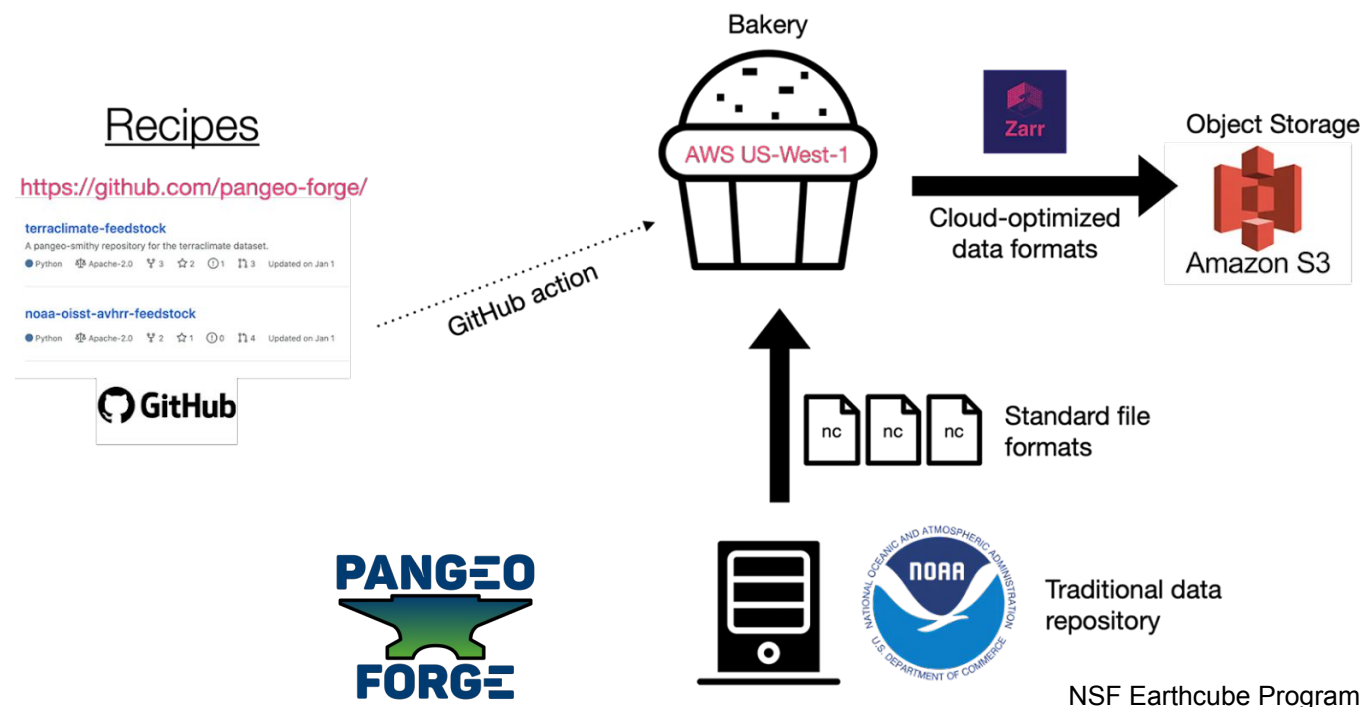
Easy and efficient
access to Geoscience
data

Pangeo-forge for efficient data management within the Pangeo Community

→ Make it easy to extract data from any traditional repository and deposit this data in cloud object storage in an analysis-ready, cloud optimized (ARCO) format;

Two main components:

- Open source Python package for describing and running data pipelines;
- Cloud platform for automatically executing recipes stored in Github repos.



STAC and Pangeo

- Pangeo-forge supports the creation of analysis-ready cloud optimized (ARCO) data in cloud object storage from "classical" data repositories;
- STAC is used to create catalog and goes beyond the Pangeo ecosystem;
- Work is ongoing to figure out the best way to expose Pangeo-Forge-generated data assets via STAC catalogs.

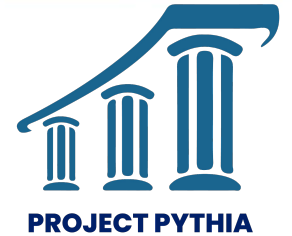


Team-up with other initiatives to
onboard new members and
increase diversity

Team-up with Project Pythia

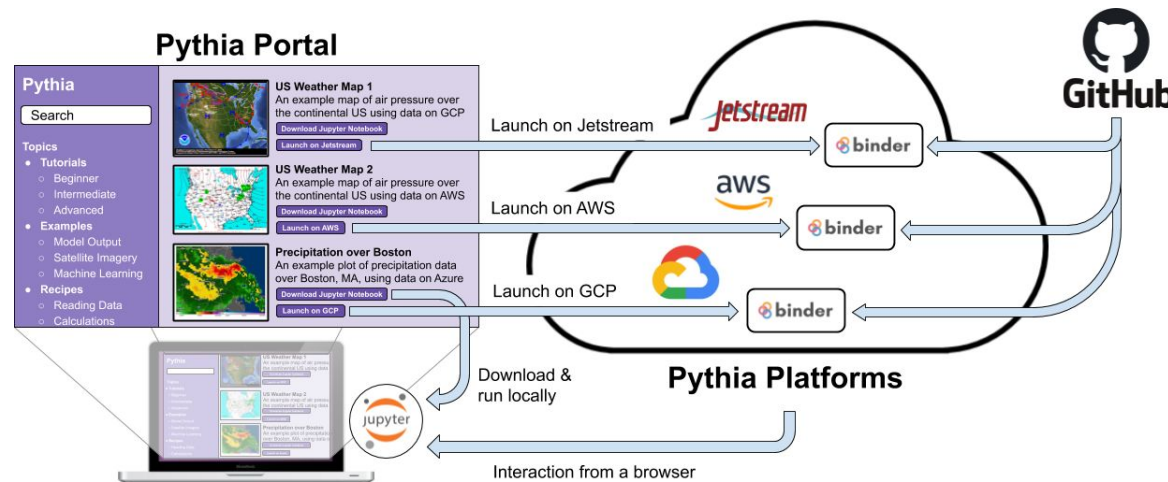
A Community Learning Resource for Geoscientists

<https://projectpythia.org/>



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

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

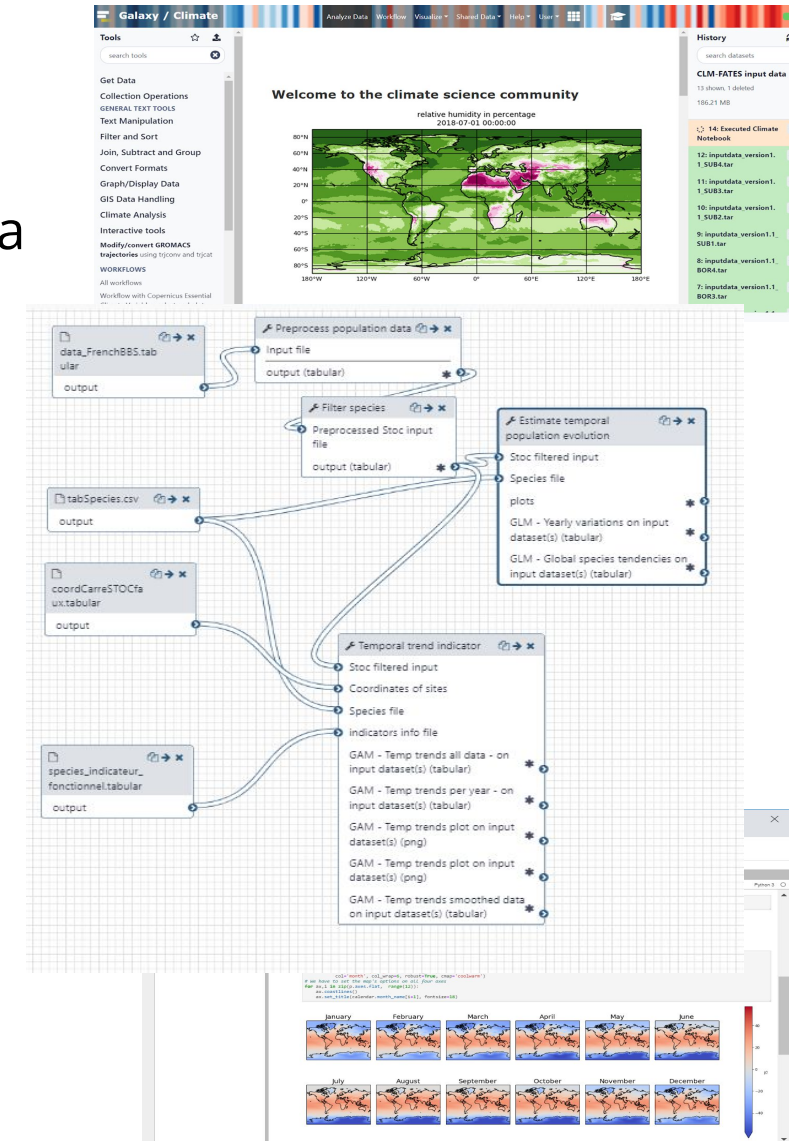


- ★ Geoscience focused
- ★ From beginner to the power user
- ★ Tutorials, videos, examples, on-line courses, and sample data
- ★ Community owned!

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

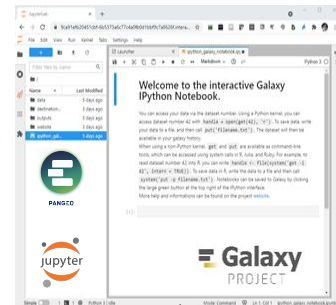


Co-design: from Jupyter Notebook to standalone Galaxy Tools



Create a Jupyter Python/R/Julia Notebook

[Galaxy Europe](#)



One Jupyter Notebook is turned into one script

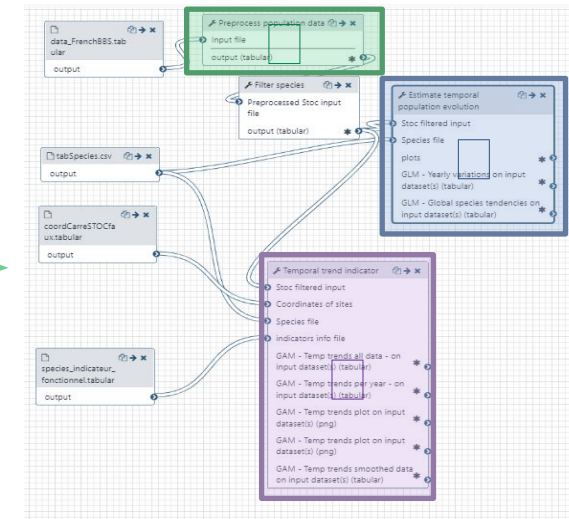
```
tab[is.na(tab)] <- 0
# filename <- "trouverUnNom"
# chemin <- paste(rep(filename, sep="/")
# write.table(tab, chemin)
colnames(tab) <- sub("nombre.", "", colnames(tab))
return(tab)

## sous jeux de données si choix d'espèce d'année ou d'un pourcentage de carrés
makesousstab <- function(tab, vecsp=NULL, echantillon=1,
  methodeechantillon="carre", vecannees=NULL) {
  cat("-- Fabrication du sous jeu de données --\n")
  flush.console()
  ## reduction de la table à certaines espèces
  if(!is.null(vecsp)) {
    cat("selection", length(vecsp), "espèce(s):\n -> ")
    cat("\n")
    tab <- data.frame(carre = tab$carre, annee = tab$annee, tab[,vecsp])
    colnames(tab) <- c("carre", "annee", vecsp)
  }
  ## reduction de la table pour certaines années
  if(!is.null(vecannees)) {
    tab <- subset(tab, annee==vecannees[1] & annee <= vecannees[2])
  }

  if(echantillon != 1) {
    if(echantillon < 1 & echantillon > 0) {
      nbinit <- nrow(tab)
      if(methodeechantillon == "global") {
        nb <- round(nrow(tab)*echantillon)
        cat("echantillonage", echantillon*100,
          "% des données par la methode", methodeechantillon, "\n")
        cat(" -> conservation de", nb, "lignes sur", nbinit, "\n")
      } else {
        if(methodeechantillon == "carre") {
          if(methodeechantillon == "carre") {
            echantillonage, echantillon*100,
            "% des carrés par la methode", methodeechantillon, "\n")
            nbcarrereinit <- length(unique(tab$carre))
            chat=sample(unique(tab$carre),
              length(unique(tab$carre))*echantillon, replace=F)
            cat(" -> conservation de", length(chat), "carrés sur",
              nbcarrereinit)
            tab=subset(tab, subset(carre %in% chat)
            cat(" ", nrow(tab), " lignes sur ", nbinit, "\n", sep="")
          } else {
            stop("Methode d'echantillonage non reconnue")
          }
        }
      }
    }
  }
}
```



Several atomic scripts from which generic and fully annotated Galaxy Tools are created



Team-up with the Environmental Data Science Book

<https://the-environmental-ds-book.netlify.app/>

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

Reproducible, scalable, & shareable
ENVIRONMENTAL DATA SCIENCE



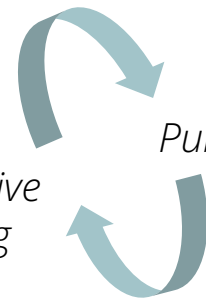
The Alan Turing Institute



Contribution



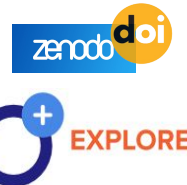
Collaborative Reviewing



Publication



jupyter {book}



Title

Tags (Environment, Theme)

RoHub FAIR Executable Research Object

launch binder

Context
purpose, highlight, contributions

Data

Analysis

Citation

Successes

- **No vendor lock-in;**
- Easy to start and **deployment on laptops, cloud and HPC;**
- “Reference” deployments on different cloud infrastructures;
- **Team up with other initiatives.** It can help to increase DEI (Diversity, Equity and Inclusion):
 - Educational material and deliver trainings (Pythia, Galaxy);
 - Training infrastructure as a Service (Pythia, Galaxy, EOSC);
 - Use Pangeo from GUI (no programming skills required) on Galaxy Europe.
- **Contribute to easy creation of data in analysis-ready, cloud optimized (ARCO) format (pangeo-forge);**
- Promote the work done by the Pangeo Community and other Geosciences initiatives (Pangeo Show & Tell/Showcase);
- Pangeo heavily used in industry;
- **Spin-off** (often from Pangeo community members) and many startups & companies using Pangeo software stack and contributing to Pangeo ecosystem.

Challenges

- **Enthusiastic individuals and volunteered driven community:** Pangeo is growing in many different ways and it is often **difficult to inform everyone** about all the initiatives;
- Initiatives are usually driven by a few dedicated individuals with **short-term funding (or no funding)**;
- Several **spin-off** but not easy to **track them down and/or help them**;
- Difficult to decide when a sub-project **needs to be “professionalized”** e.g. being able to “collect” funding, secure long-term funding (NumFOCUS but US based only);
- Difficult to onboard (train) wide range of contributors;
- **Roadmap is community driven** and can be a bit obscure to newcomers;

Lessons learned

- Important to **create an Identity** early on ;
- Need to **team-up with other initiatives** for instance to get infrastructure, training network, interdisciplinary research, and different types of contributors (**diversity of contributions**), etc.;
- Important to be able to **bootstrap new sub-projects** (new packages);
- Need to fund a **Community Manager**.

Recommendations to Space Agencies

Space Agencies could play a larger role, and in particular:

- **Contribute actively** to Open Source and Open Science projects (prevent re-inventing the wheel);
- Be more involved in Open Science initiatives, and even **lead** some of them (cf. #NASATOPS);
- Provide infrastructure (compute & Storage)
- Sponsor Open Source projects (similar to NumFOCUS) with EO industry compliance;
- Make calls to attract “small-players” and foster innovation;
- Fund training/hackathons events to engage with a wider range of actors, including training of community manager (CSCCE);
- Promote and fund initiative such as OSGeo;
- Support and invest in tech start-ups (incubators).

What about a “Matchmaker program” to connect with existing/upcoming entrepreneurs, deep-tech consultants, researchers, community managers, software developers, governments, citizen, funders, etc.?

Thank you for your attention!

