

The Pangeo Platform

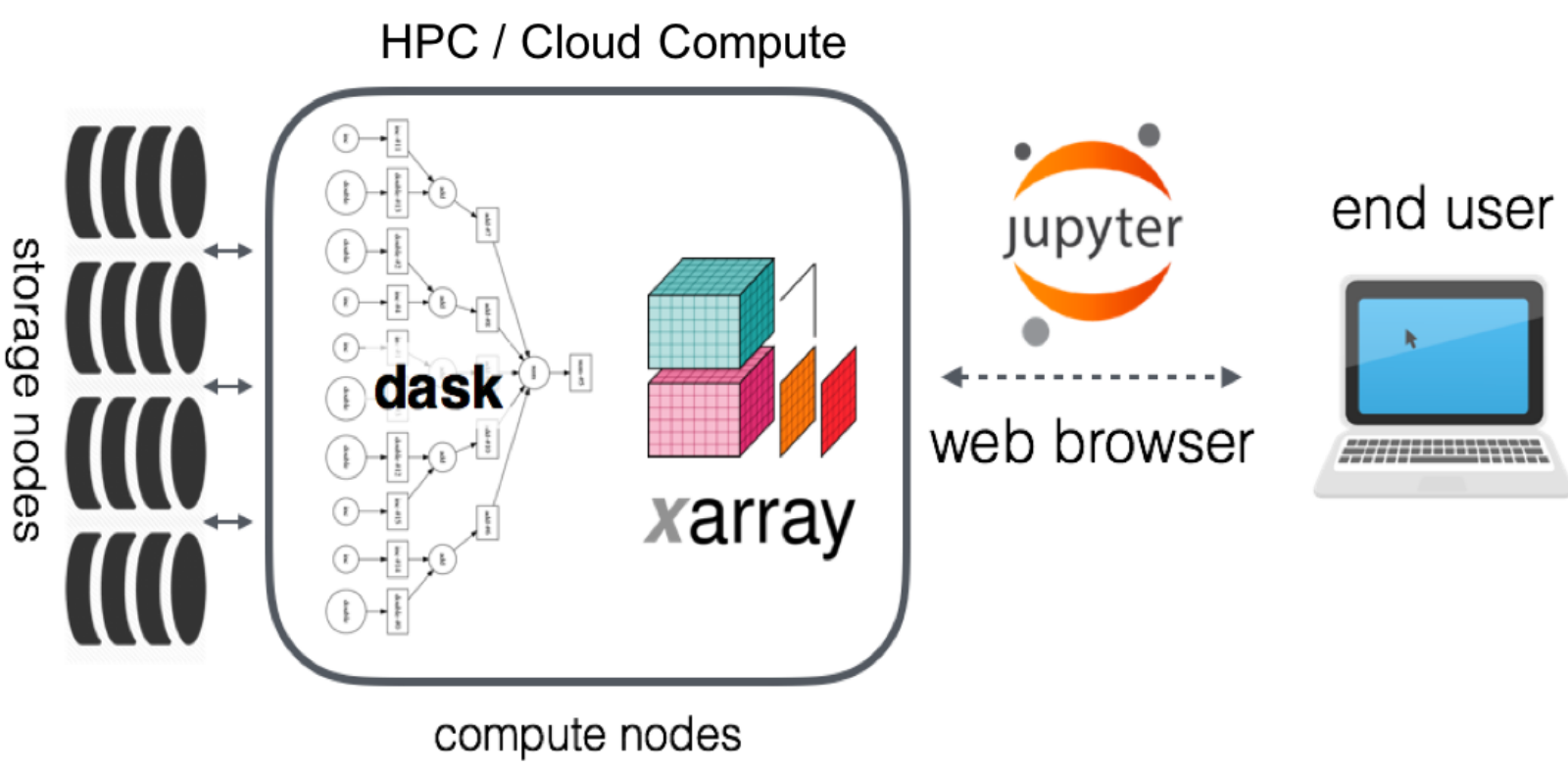
Joseph J. Hamman, NCAR
& The Pangeo Project

AGU Poster #IN11D-0693

The Pangeo Project

- Pangeo is a **community** promoting open, reproducible, and scalable science
- Pangeo is an integrated **ecosystem** of open source software tools
- Pangeo is a community **platform** for Big Data Geoscience

Architecture



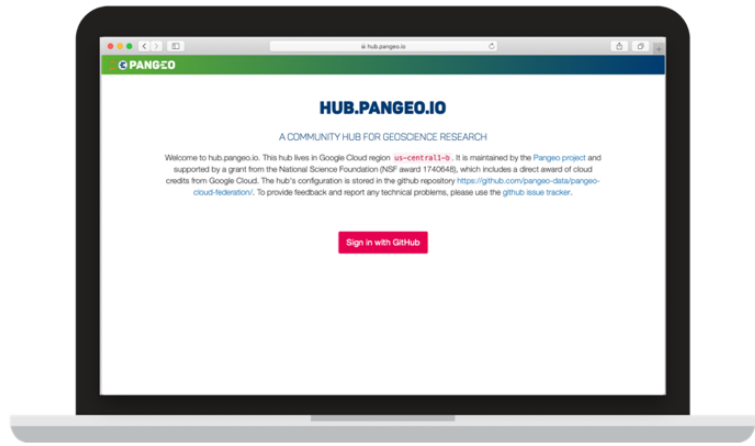
The Pangeo Platform assembles a collection of unitary components:

- User Interface (e.g. Jupyter Notebooks)
- Data Model (e.g. Xarray, Iris, or Pandas)
- Parallel Job Distribution (e.g. Dask)
- Resource Management System (e.g. Kubernetes)
- Raw Storage System (e.g. Object Storage)
- High-level Data Broker (e.g. Intake)

Read more about the design principles behind Pangeo in our recent paper: <https://arxiv.org/abs/1908.03356>.

Community Deployments

The Pangeo Community supports multiple Pangeo deployments on Google Cloud Platform and Amazon Web Services. Log in with a **GitHub ID**:



- hub.pangeo.io & ocean.pangeo.io
- aws-uswest2.pangeo.io & aws-useast1.pangeo.io
- ooi.pangeo.io
- jupyterhub.ucar.edu

Deploy Your Own Pangeo

Pangeo can be easily deployed on any Kubernetes or HPC Cluster. Learn more about the setup and deployment process: pangeo.io/setup_guides.



The Pangeo Project is supported, in part, by the National Science Foundation and the National Aeronautics and Space Administration.

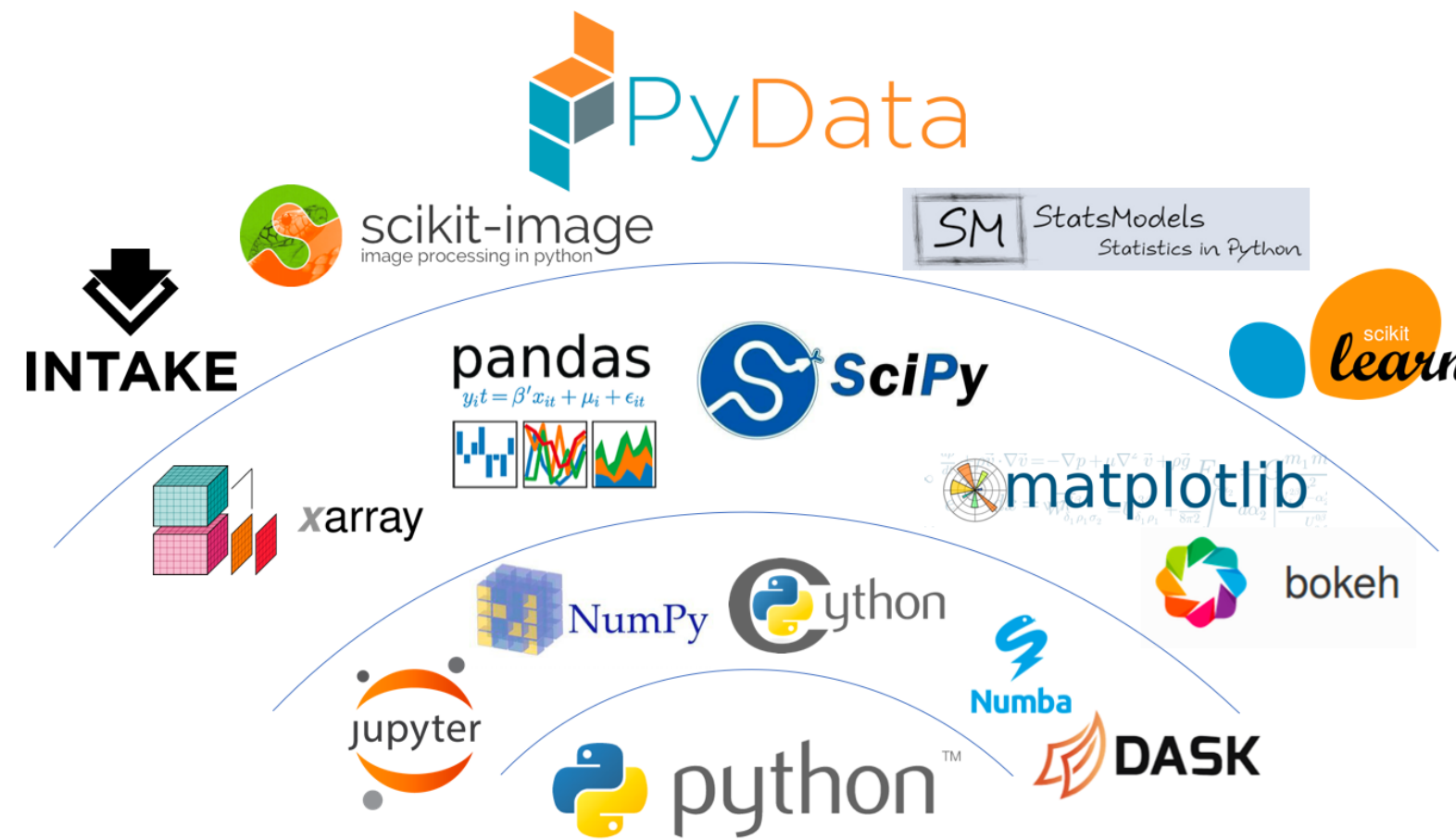
PANGEO

A community-driven,
open source,
big data platform
for the geosciences.



Take a picture to
learn more about the Pangeo Project

Scientific Python Ecosystem



Pangeo integrates across the thriving open source scientific Python ecosystem. Below, we highlight a few of the core software libraries Pangeo uses.

Interactive computing

- Multi-user gateway to single-user Jupyter Notebooks, jupyter.org/hub
- Web-based user interface for Jupyter Notebooks, jupyterlab.readthedocs.io

Data Search and Discovery

- General interface for cataloging and loading data in Python, intake.readthedocs.io
- Intake driver for loading collections of Earth Observation data, intake-stac.readthedocs.io
- Intake driver for loading catalogs of climate model data, intake-esm.readthedocs.io

```
[2]: cat = Intake.Catalog("catalog.yaml") # Load an Intake Catalog
ds = cat["gmet_v1"].to_dask()[["pcp"]] # Load an xarray dataset from catalog
display(ds)

xarray.Dataset

> Dimensions: (ensemble: 100, lat: 224, lon: 464, time: 12054)
> Coordinates: (4)
> Data variables:
pcp (ensemble, time, lat, lon) float64 dask.array<chunksize=(1, 366, 224, 464)...
```

Data Analysis

- Flexible parallel computing library for analytics, dask.org
- N-D labeled arrays and datasets in Python, xarray.pydata.org

```
[3]: # Calculate the monthly climatology
climatology = ds.groupby('time.month').mean()
climatology

[3]: xarray.Dataset

> Dimensions: (ensemble: 100, lat: 224, lon: 464, month: 12)
> Coordinates: (4)
> Data variables:
pcp (month, ensemble, lat, lon) float64 dask.array<chunksize=(1, 1, 224, 464)...
```

Data Storage

- N-D array-oriented scientific data, unidata.ucar.edu/software/netcdf
- Cloud friendly, chunked, compressed, N-D arrays, zarr.readthedocs.io
- File-system-like abstractions for remote data, filesystem-spec.readthedocs.io