# PANGEO
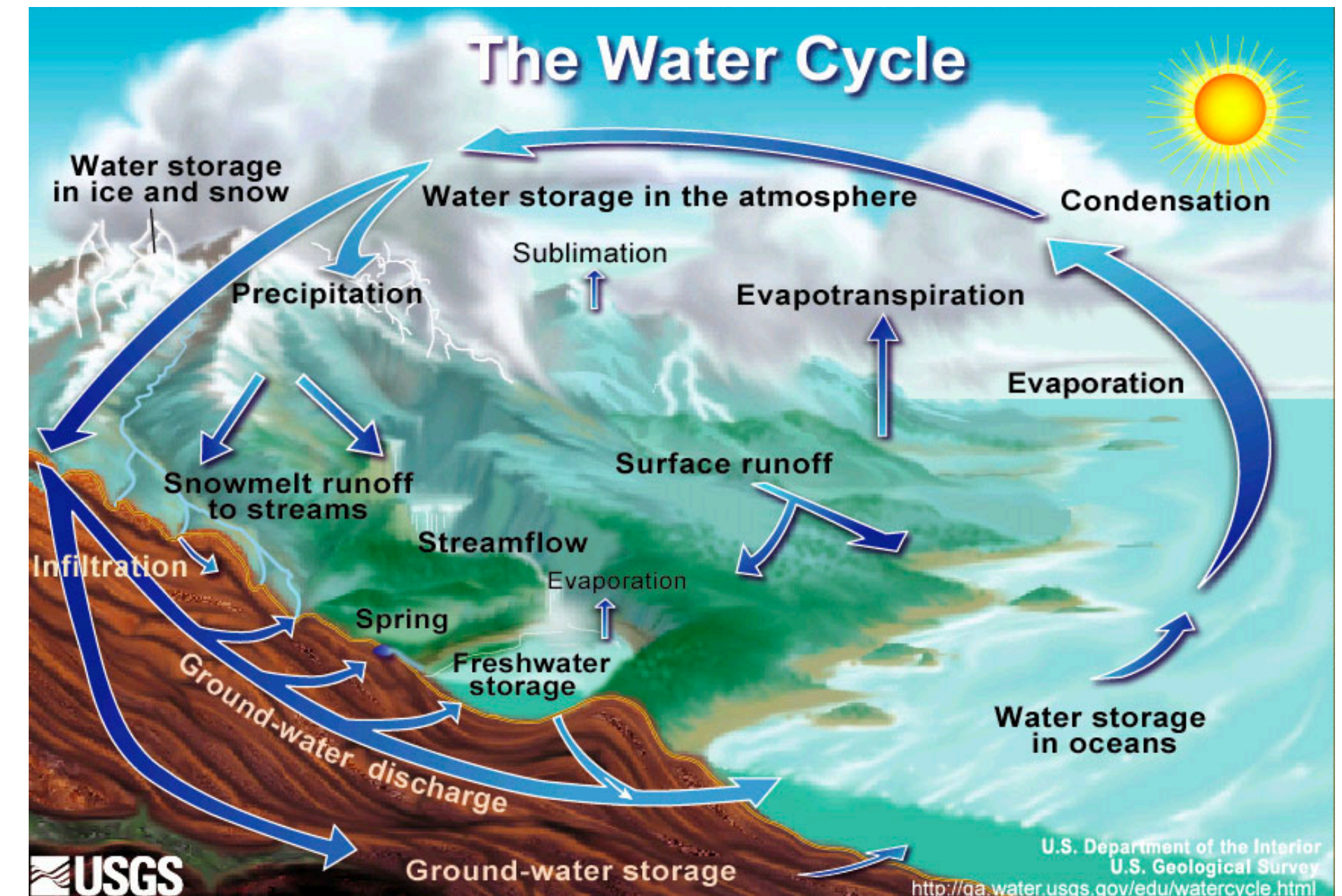
A COMMUNITY-DRIVEN EFFORT FOR
BIG DATA GEOSCIENCE

# HELLO!

- Who am I?

  ▸ Joe Hamman, Ph.D., P.E.

  ▸ I am a scientist at the National Center for Atmospheric Research (RAL & CGD)

  ▸ I study the impacts of climate change on the water cycle.

  ▸ I am a core developer of Xarray

  ▸ I am a founding member of the Pangeo project
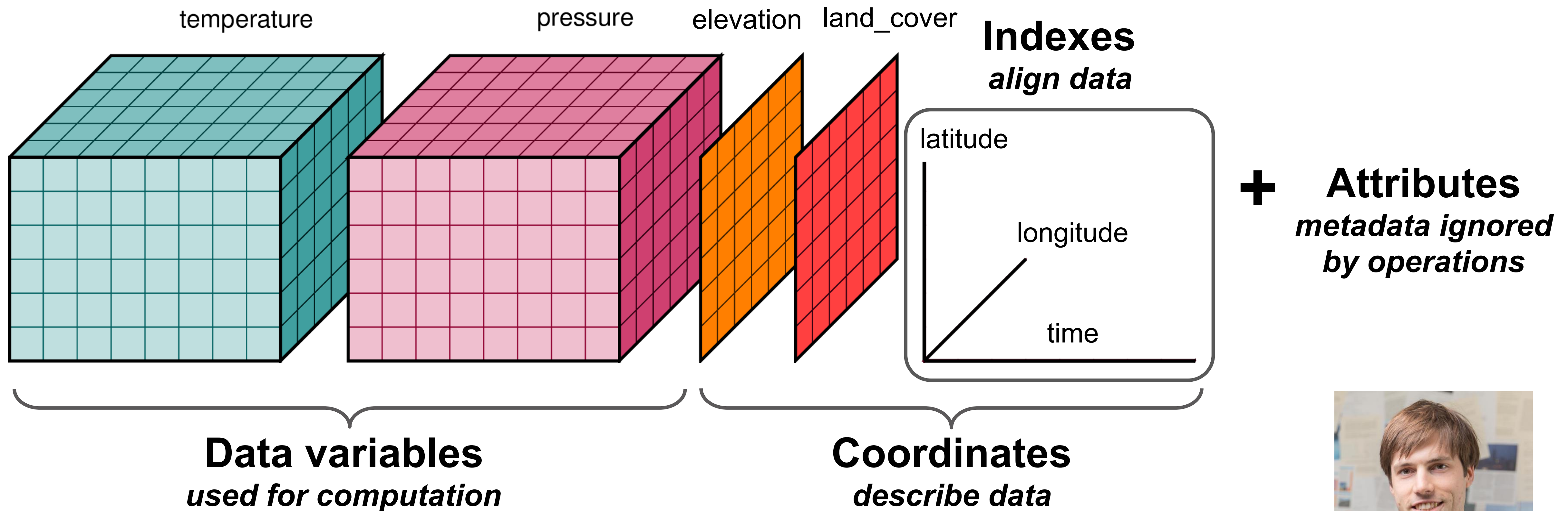


**Github:** @jhamman
**Twitter:** @HammanHydro
**Web:** joehamman.com

# SCIENTIFIC PYTHON FOR DATA SCIENCE



**Credit: Stephan Hoyer, Jake Vanderplas (SciPy 2015)**

# XARRAY DATASET: MULTIDIMENSIONAL VARIABLES WITH COORDINATES AND METADATA



temperature

pressure

elevation

land_cover

**Indexes**
*align data*

latitude

longitude

time

**+ Attributes**
*metadata ignored by operations*

**Data variables**
*used for computation*

**Coordinates**
*describe data*

*"netCDF meets pandas.DataFrame"*
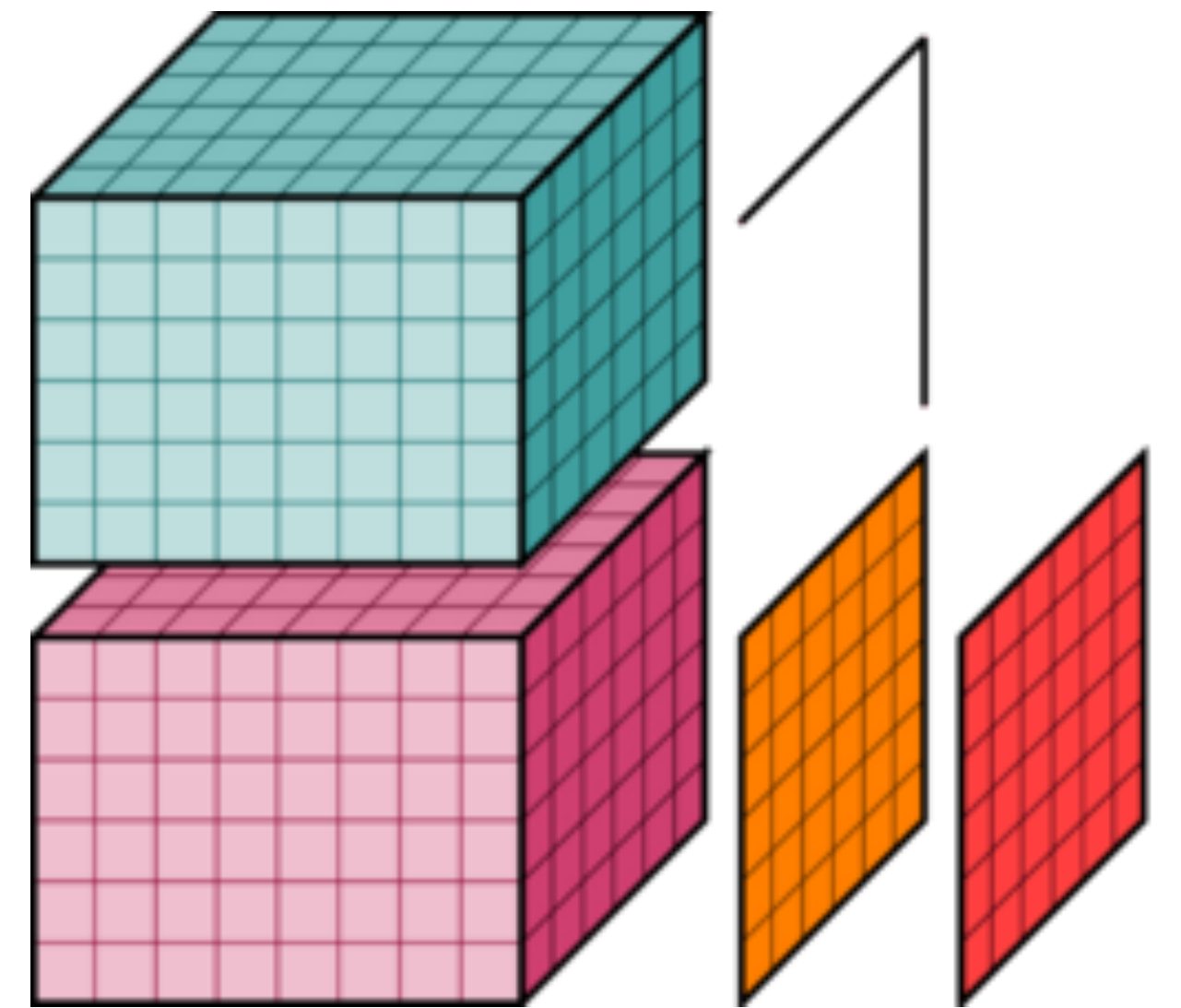
**Credit: Stephan Hoyer**

# X A R R A Y

## http://xarray.pydata.org

- label-based indexing and arithmetic

- interoperability with the core scientific Python packages (e.g., pandas, NumPy, Matplotlib)

- out-of-core computation on datasets that don't fit into memory (thanks dask!)

- wide range of input/output (I/O) options: netCDF, HDF, geoTIFF, zarr

- advanced multi-dimensional data manipulation tools such as group-by and resampling
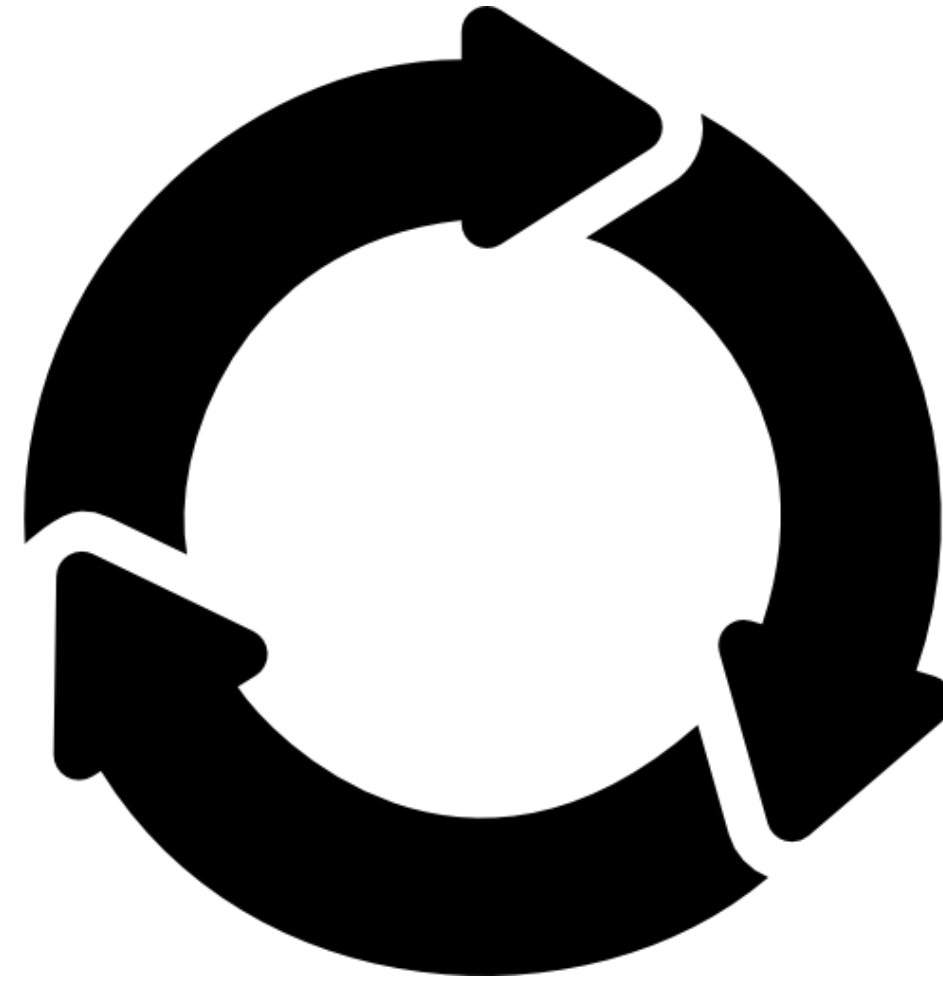
# XARRAY UPDATES

- **Development Roadmap** (http://xarray.pydata.org/en/stable/roadmap.html)

  ▸ More flexible grids/indexing

  ▸ More flexible arrays/computing

  ▸ More flexible storage backends

- **NumFOCUS Sponsorship**

  ▸ https://numfocus.org/project/xarray

- **New Contributors**

  ▸ New core devs: Spencer Clark and Deepak Cherian

# WHAT DRIVES PROGRESS IN GEOSCIENCE?



**New Ideas**

$$q_{liq,z}^{soil} = \begin{cases} q_{rain} - q_{ix} - q_{sx} & z=0 \\ -K^{soil}\dfrac{\partial \psi}{\partial z} + K^{soil} & z > 0 \end{cases}$$
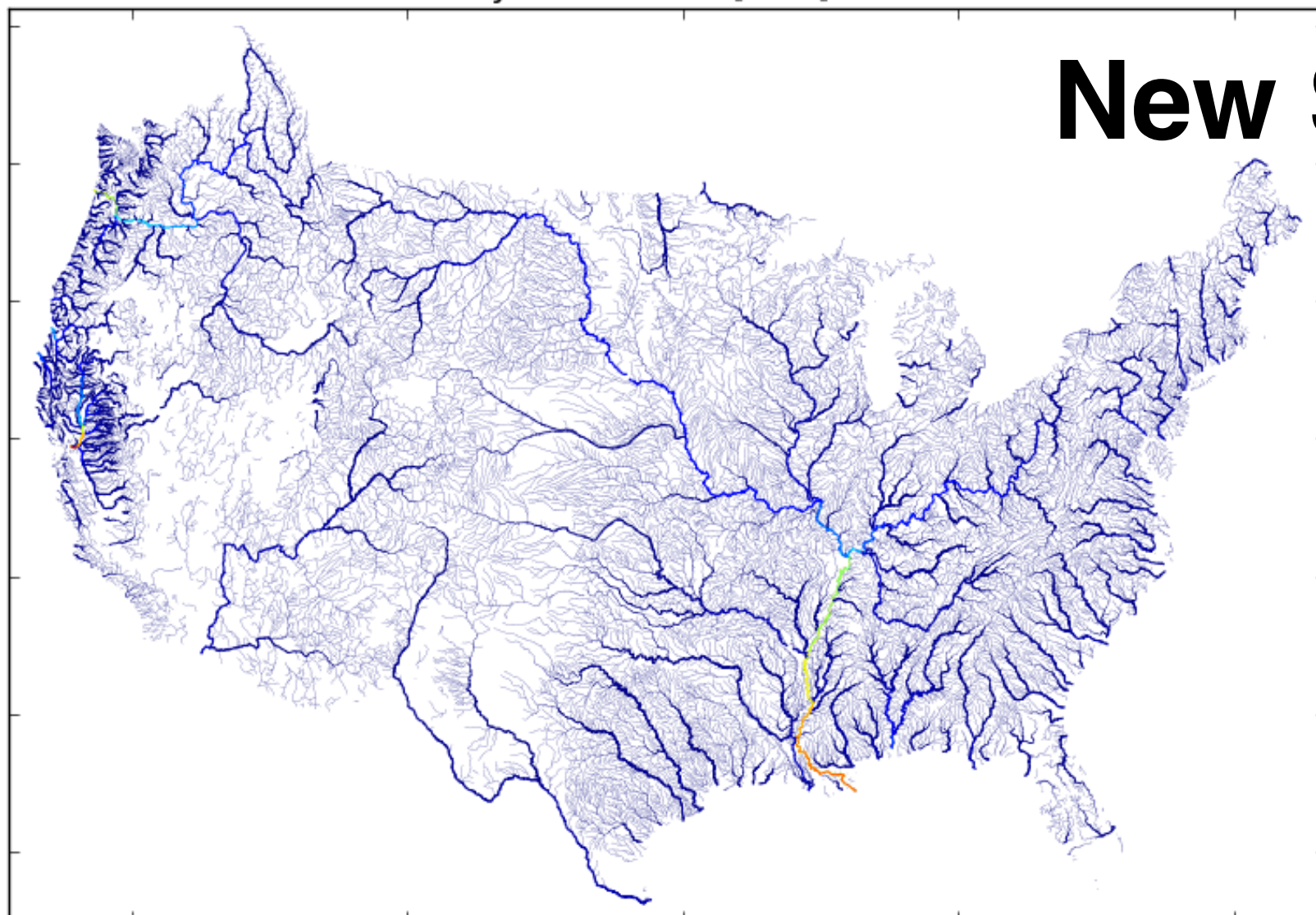
**New Observations**

**New Simulations**

Monthly Streamflow [cms] 2000-1

Left: The Soil Moisture and Ocean Salinity (SMOS)
www.smos-mode.eu

Right: The Soil Moisture Active/Passive (SMAP)
mission www.jpl.nasa.gov

# FRAGMENTATION PROBLEMS

1. **Software**

   - Few tangible incentives to share source code (funding agencies, journals)

   - Lack of extensible development patterns; often it is easier to "home grow" your own solution, rather than using someone else's.

   - Result is that most geoscientific research is effectively unreproducible and prone to failure.

2. **Data sprawl**

   - Inefficiencies of many copies of the same datasets ("dark replicas")

   - Lessons learned from the CMIP archives (CMIP3 was duplicated > 30x)

3. **Local vs. High-performance vs. Cloud Computing**

   - Traditional scientific computing workflows are difficult to port from a laptop, to HPC, to the cloud

# PANGEO PROJECT GOALS

- Foster collaboration around the open source scientific Python ecosystem for ocean / atmosphere / land / climate science.

- Support the development with domain-specific geoscience packages.

- Improve scalability of these tools to to handle petabyte-scale datasets on HPC and cloud platforms.

# PANGEO COLLABORATORS



And many more…

# PANGEO ARCHITECTURE



Cloud / HPC

Distributed storage

"Analysis Ready Data" stored on globally-available distributed storage.

Jupyter for interactive access remote systems

end user

web browser

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.

# BUILD YOUR OWN PANGEO

| | | | |
|---|---|---|---|
| **Storage Formats** | HDF | OPeNDAP | Cloud Optimized COG/Zarr/Parquet/etc. |
| **ND-Arrays** | NumPy | DASK | More coming… |
| **Data Models** | xarray | Iris | pandas $y_it = \beta' x_{it} + \mu_i + \epsilon_{it}$ |
| **Processing Mode** | jupyter Interactive | Batch | Serverless |
| **Compute Platform** | HPC | Cloud | Local |

# PANGEO DEPLOYMENTS

HTTP://PANGEO.IO/DEPLOYMENTS.HTML

PANGEO.PYDATA.ORG
BINDER.PANGEO.IO

**NASA Pleiades**

**Over 1000 unique users since March!**

Google Cloud Platform

**NCAR Cheyenne**

Microsoft Azure

aws

(SCALE USING JOB QUEUE SYSTEM)

(SCALE USING KUBERNETES)

# FOO.PANGEO.IO
## DEPLOY YOUR OWN PANGEO

- What's in a typical Pangeo?

  - JupyterHub interface

  - Tools to deploy dask clusters

  - Customizable software/hardware environment

- Current effort to federate pangeo deployments for problem specific uses (e.g. cds.pangeo.io?)

- Custom deployments:

  - polar.pangeo.io

  - solar.pangeo.io

  - ocean.pangeo.io

  - hydroshare.pangeo.io

  - And more coming…

# BINDER.PANGEO.IO

- BinderHub

  ▸ Highly customizable Jupyter environment

  ▸ Automates Git repo -> docker image -> Jupyter notebook

  ▸ Automates deployment of Dask clusters

- Easiest way to share Pangeo workflows

- Try it: https://bit.ly/2O9qJr3

PANGEO IN A NUTSHELL

- **Scientific Python ecosystem**

  ‣ flexible, open-source, community driven

- **Interoperable**

  ‣ integrates with existing/developing tools used by science community

- **Analysis ready data formats**

  ‣ cloud optimized data (e.g. zarr)

- **Intuitive self-describing data models**

  ‣ e.g. xarray, Iris

- **Scalable**

  ‣ e.g. Dask, Kubernetes

- **Interactive**

  ‣ e.g. Jupyter, JupyterHub, BinderHub

- **Cross platform**

  ‣ HPC, Cloud, local computing

# WHAT'S COMING FOR PANGEO

- Governance (https://github.com/pangeo-data/governance)

- Funding (new projects from NASA and NSF)

- AWS Open Datasets Program and Pangeo compute resources

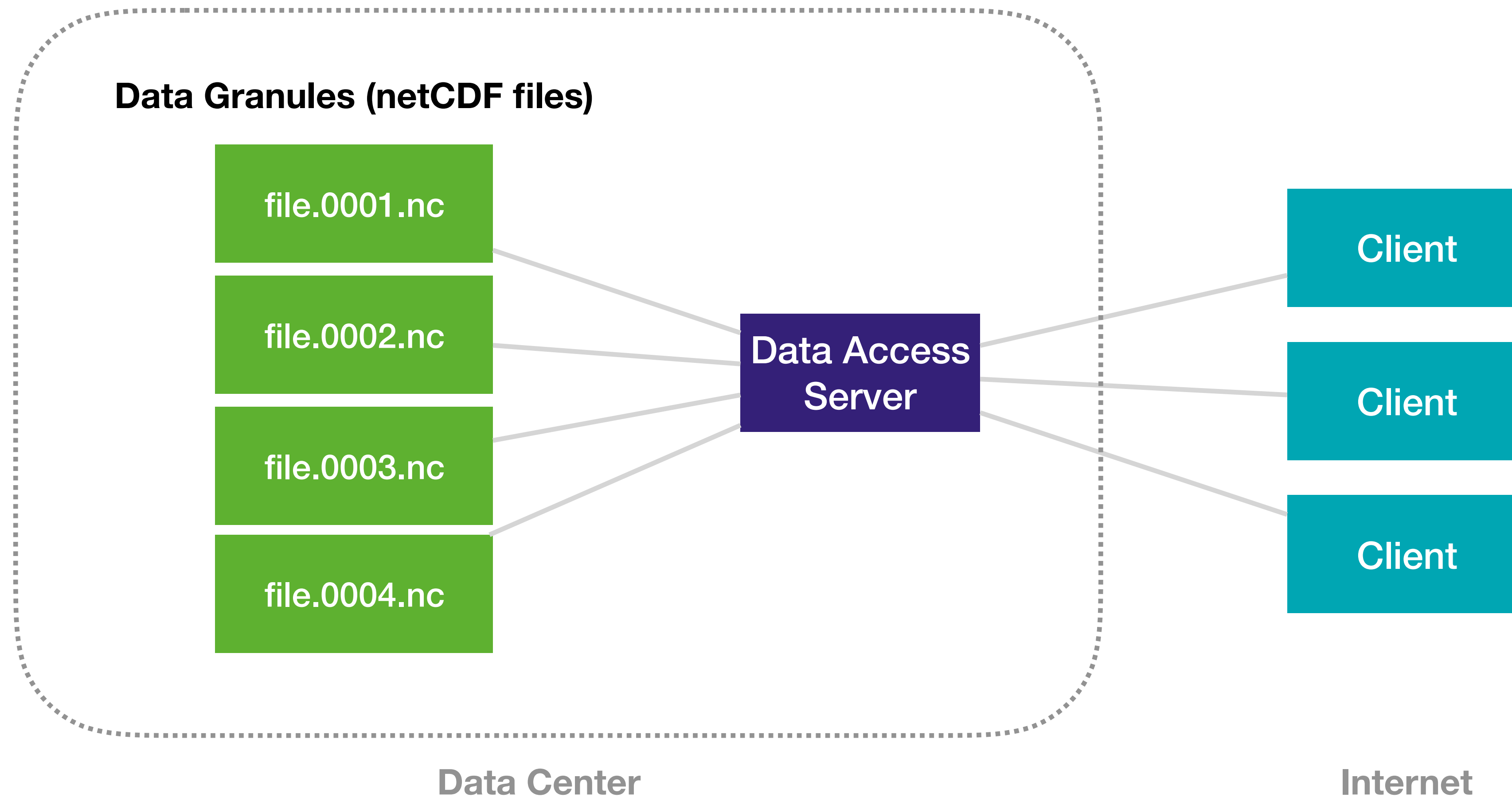- Science focus on remote sensing datasets

- Looking for new community partners

# HOW TO GET INVOLVED

## HTTP://PANGEO.IO

- Access and existing Pangeo deployment on an HPC cluster, or cloud resources (eg. binder.pangeo.io)

- Adapt Pangeo elements to meet your projects needs (data portals, etc.) and give feedback via GitHub: github.com/pangeo-data/pangeo

- Participate in open-source software development!

# SHARING DATA IN THE CLOUD

## Traditional Approach: A Data Access Portal

PANGEO

# ON-DEMAND ANALYSIS-READY DATA

- **Too big to move**: assume data is to be used but not copied

- **Self-describing**: data and metadata packaged together

- **On-demand**: data can be read/used in its current form from anywhere

- **Analysis-ready**: no pre-processing required

# SHARING DATA IN THE CLOUD

## Direct Access to Cloud Object Storage

**Zarr**

**?**

**Data Granules
(netCDF files or something new)
Cloud Object Storage**

chunk.0.0.0

chunk.0.0.1

chunk.0.0.2

chunk.0.0.3

**Cloud Compute
Instances**

Client

Client

Client

Catalog

**Cloud Data Center**