# Integrated Data Analytics Needs in ESGF2-US

Forrest M. Hoffman (ORNL) and the ESGF2-US Team

**ESnet Confab25 – San Francisco, California**

*April 9, 2025*
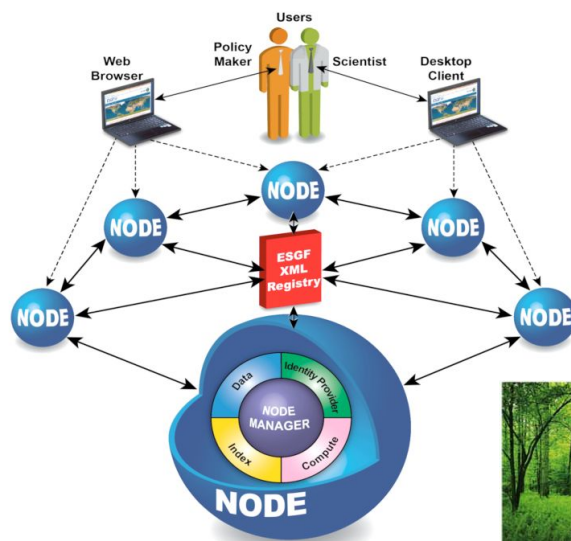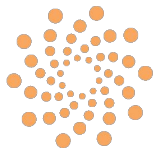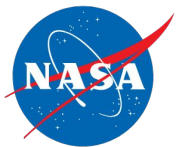
# What is the Earth System Grid Federation?

- **Earth System Grid Federation (ESGF)** is an _international consortium_ and a _globally distributed peer-to-peer network of data servers_ using a common set of protocols & interfaces to archive and distribute Earth system model output and related input, observational, and reanalysis data

- **Open Science data** are used by scientists all over the world to investigate Earth system variability and feedbacks and to inform research and assessments
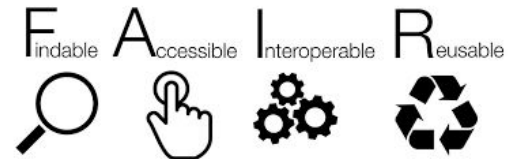
## ESGF Conceptual Diagram



_Model data from ESGF are used to understand key Earth system processes and interactions_

Logos represent primary international contributors: US Department of Energy, NASA, NOAA, NSF, European IS-ENES Project, and Australian NCI

# ESGF2 US ESGF Holdings are Open and Large

F indable   A ccessible   I nteroperable   R eusable

- CMIP5 totals >1.5 PB (>5 PB including replicas)

- CMIP6 totals >16.1 PB (>27 PB including replicas)

- CMIP7 is expected to have more experiments, high resolution output, and ensembles, totaling ~100 PB

- ESGF is concerned with the _full stack security_ and the _integrity of the data_, but we are **not** concerned about controlling _access to the data_

| CMIP6 | 14,893,892 total datasets 27,983.73 TB | CMIP6 | 7,670,309 distinct datasets 16,120.41 TB | CMIP6 | 7,223,583 replica datasets 11,863.32 TB |
|---|---|---|---|---|---|
| CORDEX | 187,785 total datasets 1,473.33 TB | CORDEX | 187,513 distinct datasets 1,472.77 TB | CORDEX | 272 replica datasets 0.56 TB |
| CMIP5 | 201,130 total datasets 5,293.61 TB | CMIP5 | 52,163 distinct datasets 1,525.07 TB | CMIP5 | 148,967 replica datasets 3,768.55 TB |
| INPUT4MIPS | 5,871 total datasets 10.84 TB | INPUT4MIPS | 21 distinct datasets 0.9 TB | INPUT4MIPS | 5,850 replica datasets 9.95 TB |
| OBS4MIPS | 126 total datasets 0.2 TB | OBS4MIPS | 108 distinct datasets 0.2 TB | OBS4MIPS | 18 replica datasets 0.01 TB |

As of April 7, 2025

ESGF2
US

**DOE's Next Generation ESGF**

- As many as 3 nodes located at DOE's major computing facilities

- Replicating data from the worldwide Federation

- Providing scalable cloud indexing and tape archiving

# The Solution? Add a Globus-Compute layer!

# Full Flows: Automating with Globus Flows

- Use Case:
  - User would like all yearly averages climate simulations from 2050 to 2070, over the United States
    - They would call a request to globus-compute
    - The output would be saved on that remote machine
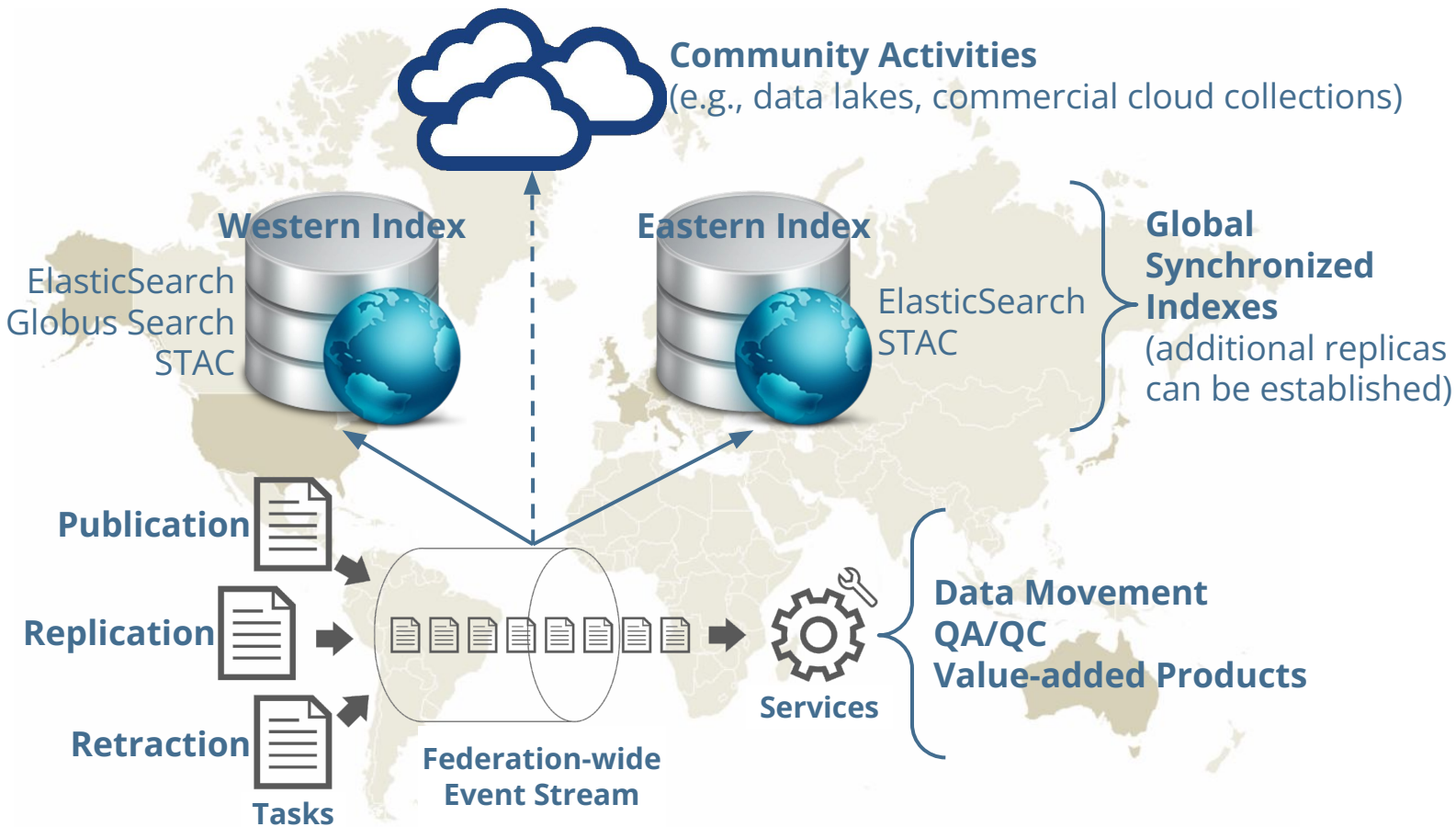    - A **guest collection** would be returned to the user, which they could either
      - Automatically transfer to their local machine (if a local endpoint is specified)
      - Extract the guest collection URL, which they can share with collaborators!
- This allows
  - A more secure method for running the WPS and gathering metrics of users
  - A more streamlined method of saving output, without filling up temporary space
  - Users can share this with collaborators easily, develop workflows around it, etc.

# What's Our Goal?

**Objective**: Remove the barriers and accelerate science with ESGF-hosted data

**Data access**: Develop improved APIs and services to access and analyze data;

**Server-side functions**: make it easy to run core operators (averaging, selecting, regridding) next to the data;
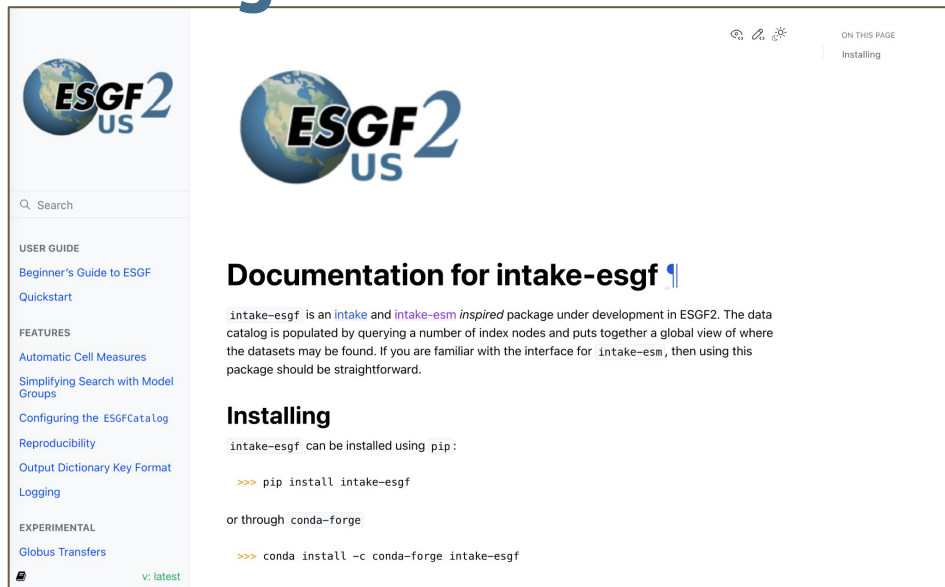
**User computing services:** Data proximate computing resources; reproducible/relocatable workflows;

**Community development**: Don't reinvent the wheel - use and improve existing solutions, entrain the community;

# Integrating with intake-esgf

- Improve the APIs to access data; simplify searching for data programmatically across the federation

- Provide STAC-based index query in addition to the existing Solr and Globus indices

- Extend the interface to provide capability for data streaming (OPeN-DAP, Kerchunk, Virtual Zarr) as available

- Integrate the errata service provided by es-doc into intake-esgf catalogs



- Intelligently determines the quickest way to access data (download, Globus Transfer, stream, load locally)

- **Provides method to package compute + flows**

Repository: https://github.com/esgf2-us/intake-esgf
Documentation: https://intake-esgf.readthedocs.io/
Installation: PyPI and Conda-forge

# Summary of Integration Activities

- All **ESGF development is being performed collaboratively** with Federation partners
- **New data projects** for downscaled projections (LOCA2, STAR-ESDM) were added; we will add large-scale AI/ML data, large ensembles simulations and intercomparisons
- **User computing** approaches initiated in the commercial cloud and deployed through on-premise cloud infrastructure will enable computing near the data
- Specific **integration activities**:
  - **Sharing data indexes** across DOE-BER platforms (ARM Data Center, ESS-DIVE, etc.)
  - Unifying on **Federated authentication** (*Globus Auth*) to simplify data access and enable cross-platform/cross-facility data access and analysis
  - **Integrating software stacks** for data access, analysis, and visualization for Jupyter
  - A few global **scalable data index** and search instances (*Globus Search*)
  - **Managed automation** of data publishing workflows (*Globus Flows*)
  - **Server-side computing** spawned by web or Jupyter/Python (*Web Processing Service* and *Globus Compute*) for generating value-added products and subsetting & summarizing data across platforms and facilities