

# Scalable Earth Observation ML Workflow in Climate Applications

**Takuya Kurihana**, Tyler J. Skluzacek, Rafael Ferreira da Silva, Valentine Anantharaj  
Oak Ridge National Laboratory



U.S. DEPARTMENT OF  
**ENERGY**

# High quality & large-volume of Earth Observation (EO) datasets are distributed across multiple platforms

<https://www.earthdata.nasa.gov/eosdis/cloud-evolution>

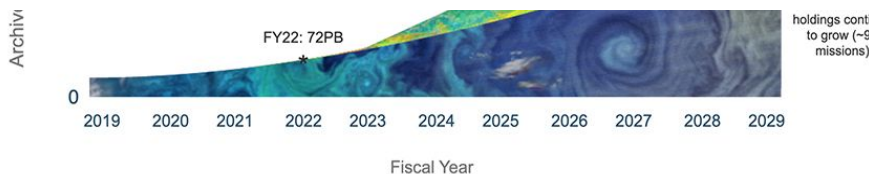
## ESDS Program Earthdata Cloud Evolution

More About EOSDIS

11,000+ Data products



With the launch of new, high-data-volume Earth observation missions, NASA's ability to effectively ingest, and archive large amounts of data requires the most cost-effective, flexible, and scalable data-management architectures and technologies. To meet these demands, NASA's Earth Science Data Systems (ESDS) Program implemented a strategic vision to develop and operate multiple components of NASA's Earth Observing System and Information System (EOSDIS) in a commercial cloud environment.



<https://www.earthdata.nasa.gov/esds/esds-highlights/2022-esds-highlights>

72PB

## Overview and NCICS/CISESS Contributions

NOAA's Big Data Program (BDP) is designed to facilitate public use of key environmental datasets by providing copies of NOAA's information in the Cloud, allowing users to do analyses of data and extract information without having to transfer and store these massive datasets themselves.

<https://ncics.org/data/noaa-big-data-project/>

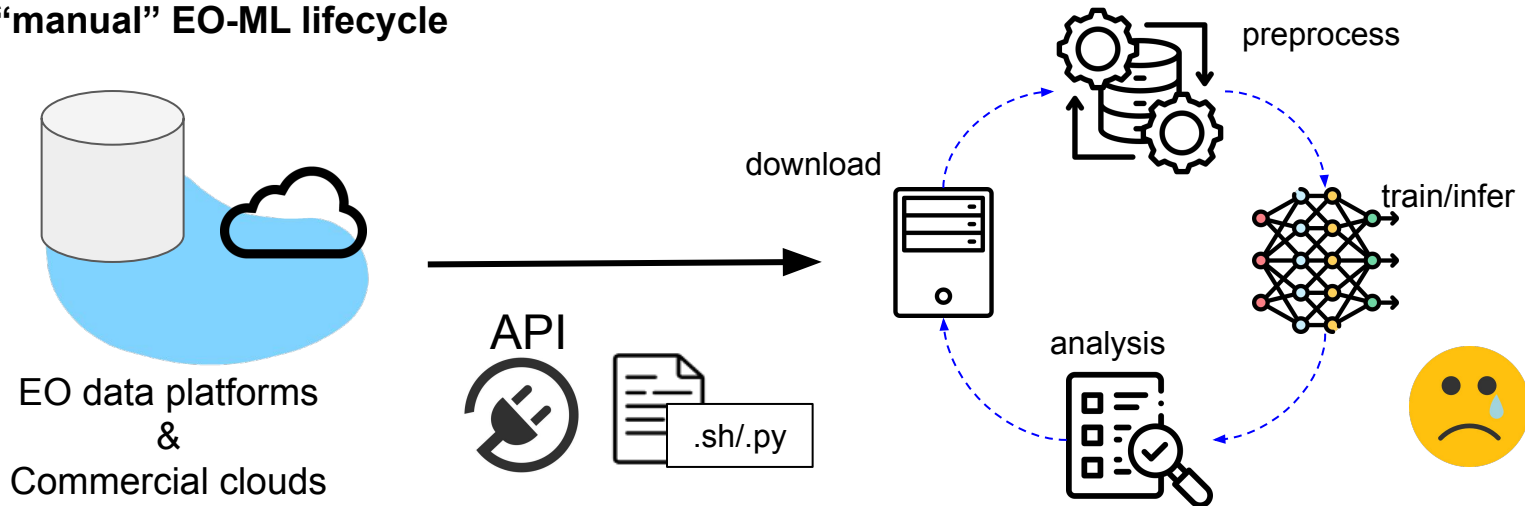
<https://www.arm.gov/news/features/post/88126>



# Manual & disconnected processes in ML workflow are time-consuming & computationally intensive

Complex data curation & post-analysis process are often required in EO-ML pipeline

## Example of “manual” EO-ML lifecycle

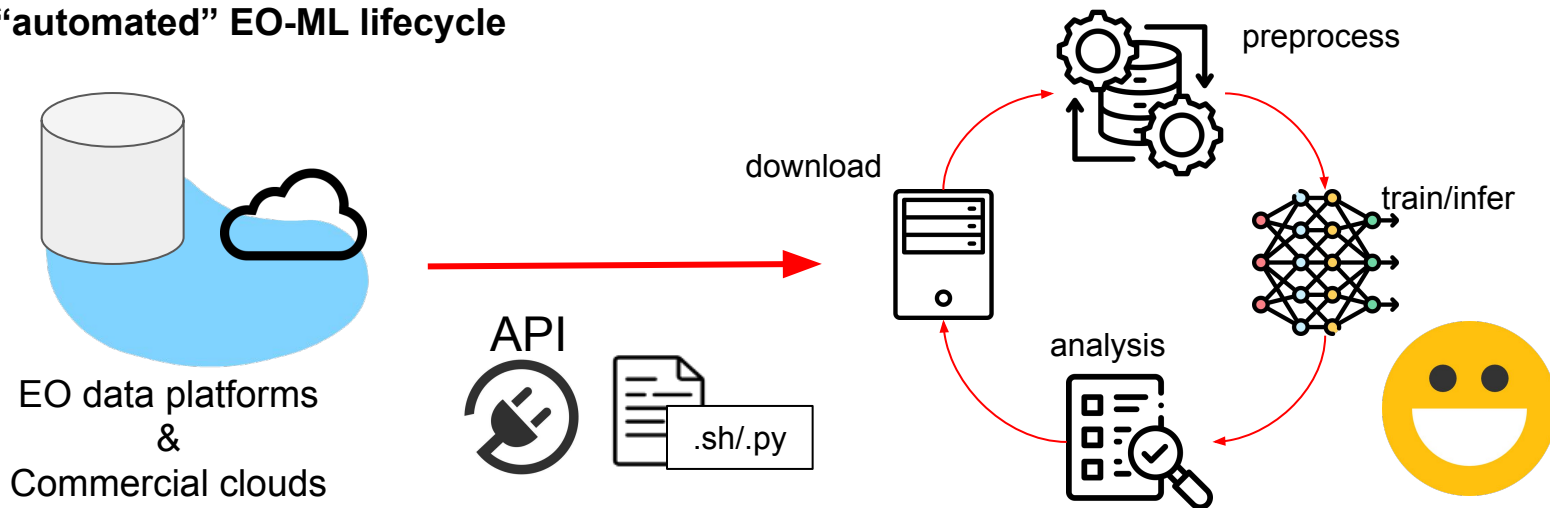


**Disconnected** among workflow components

# Globus software tools to harness the potential of EO-ML workflow for accelerating scientific discovery

Orchestrate data collection, movement, and processing across multiple computational platforms.

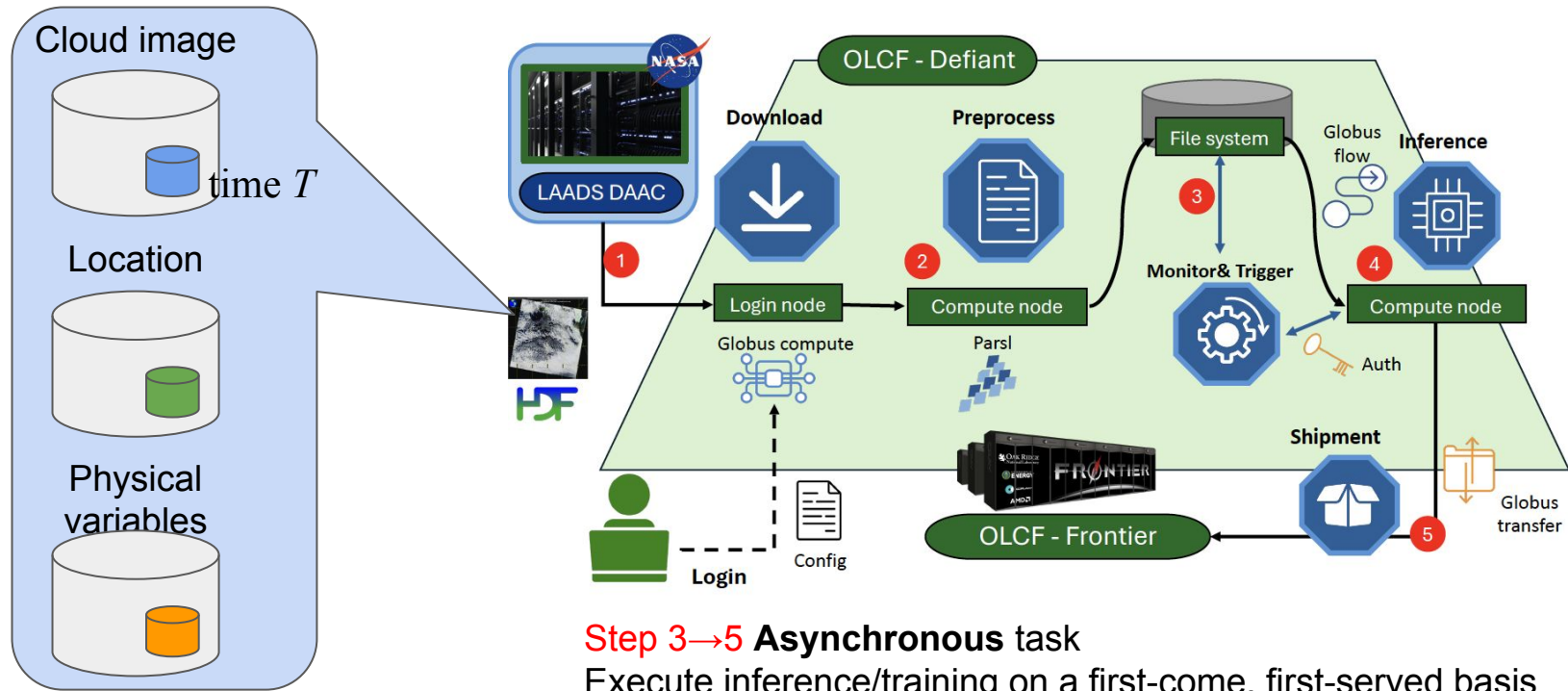
## Example of “automated” EO-ML lifecycle



# Target: Self-supervised deep learning to classify satellite clouds images

Total volume ~ 1PB

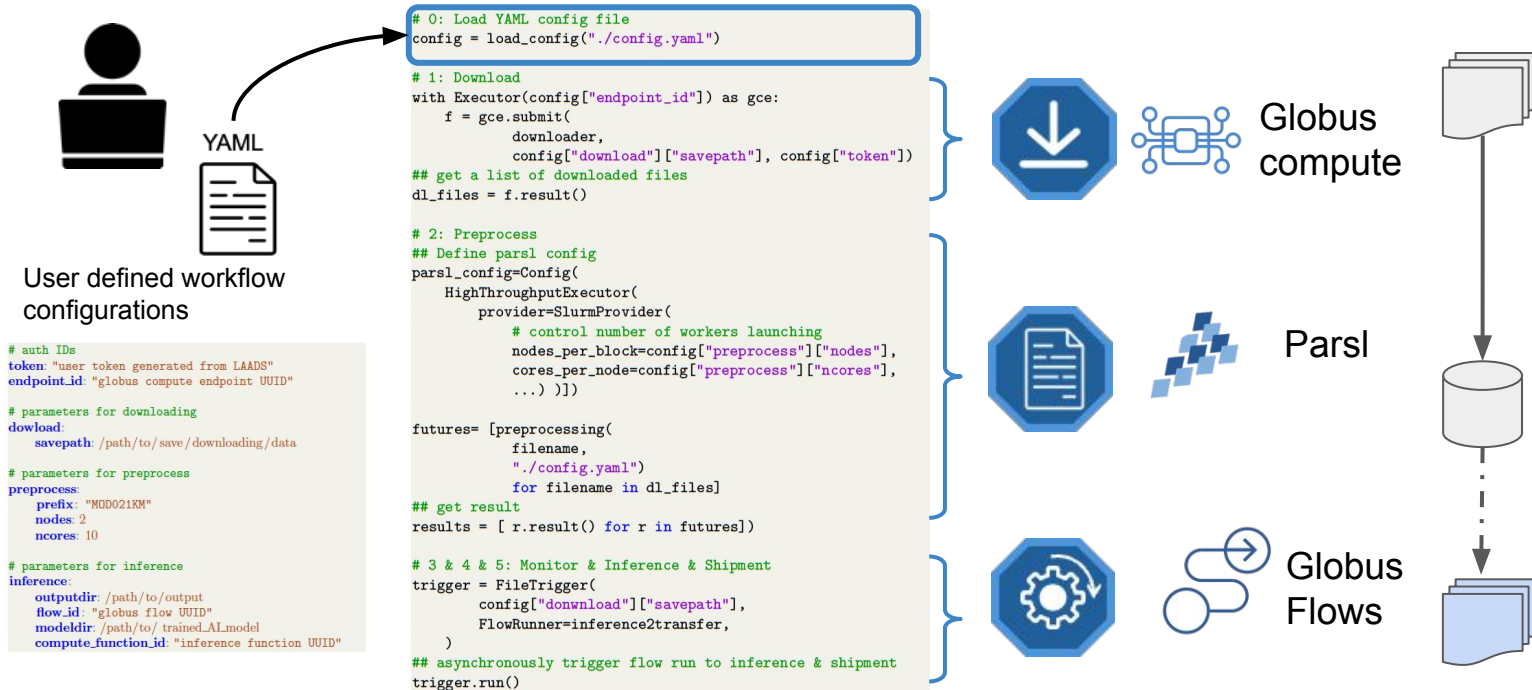
Step 1→2 **Synchronous** task  
Wait for files with same timestamp



Step 3→5 **Asynchronous** task  
Execute inference/training on a first-come, first-served basis



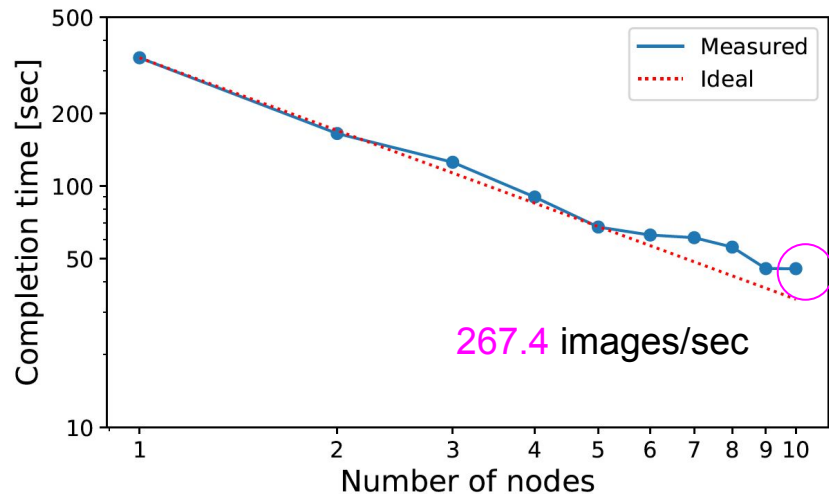
# Globus Computes/Flow enable to easily create a seamless & flexibly multi-facility ML workflow for domain scientists



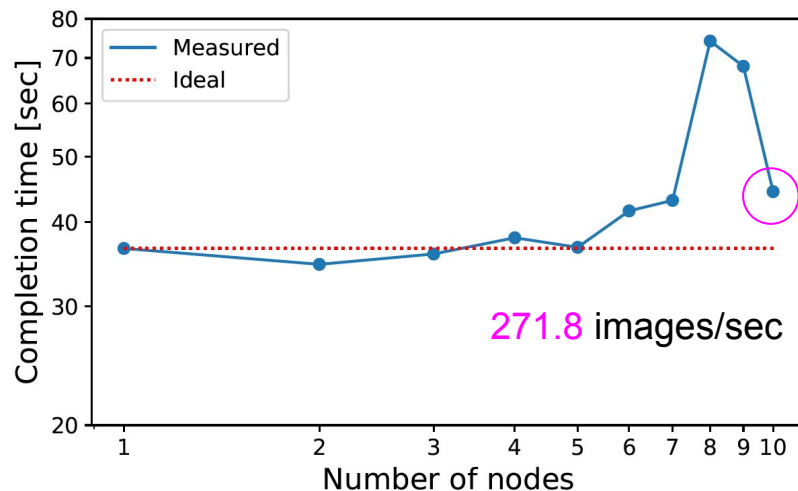
# Parsl enables scaling data intensive pipeline

The preprocessing step of our EO-ML workflow achieves optimal scaling by increasing the number of nodes

### Strong scaling

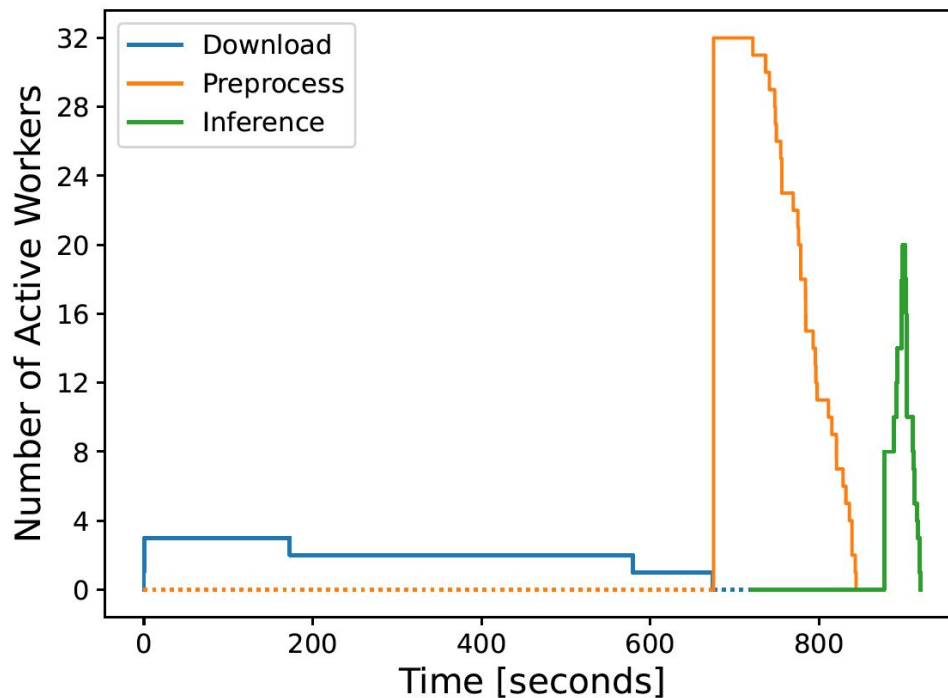


### Weak scaling



# Flexibly manage resources for workflow components

Specify different computing resources based on nature of tasks & machine availability.





# Summary, Future works & Question

kurihanat@ornl.gov

Xloop @SC'24



## Summary:

- Demonstrated a deployment of EO-ML workflow on OLCF for peta-scale EO dataset
- Leveraged Globus Compute, Parsl, and Globus Flows to automate ML lifecycle
- Globus software tools provide excellent scaling performance and fast overhead times.

## Future works:

- Complete more seamless workflow with Globus Flows
- Shareable and publicly accessible repository of complete workflows or individual workflow steps
- Multi-DOE facility EO-ML training workflow



# Globus software ecosystem automates **multi-facility** EO-ML workflow execution on OLCF Advanced Computing Ecosystem

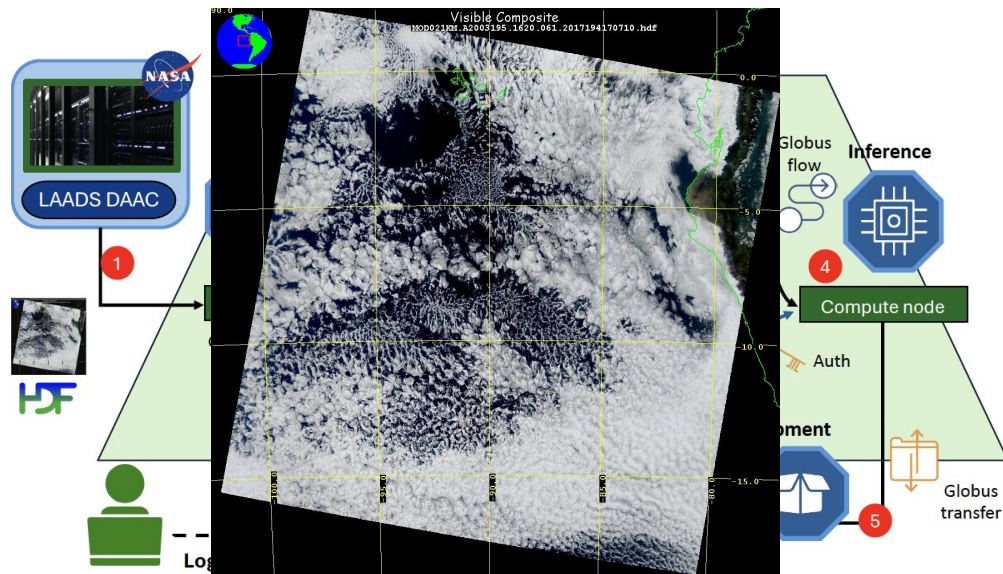
## Target application:

*Self-supervised cloud classification → classify satellite clouds images*

Problem state *before* automation:

- 3 different satellite image products (~850TB) from NASA
- Downloading was remotely executed by **funcX**
- Preprocess was accelerated by **Parsl**

**Question:** How to connect data pipeline over multiple facilities?



# Globus software ecosystem automates **multi-facility** EO-ML workflow execution on OLCF Advanced Computing Ecosystem

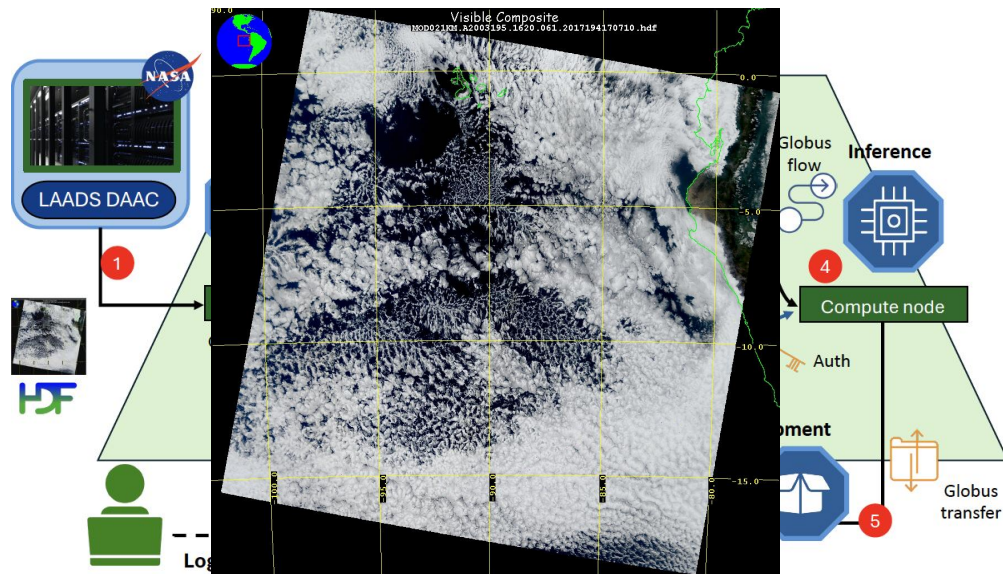
## Target application:

*Self-supervised cloud classification → classify satellite clouds images*

Problem state *before* automation:

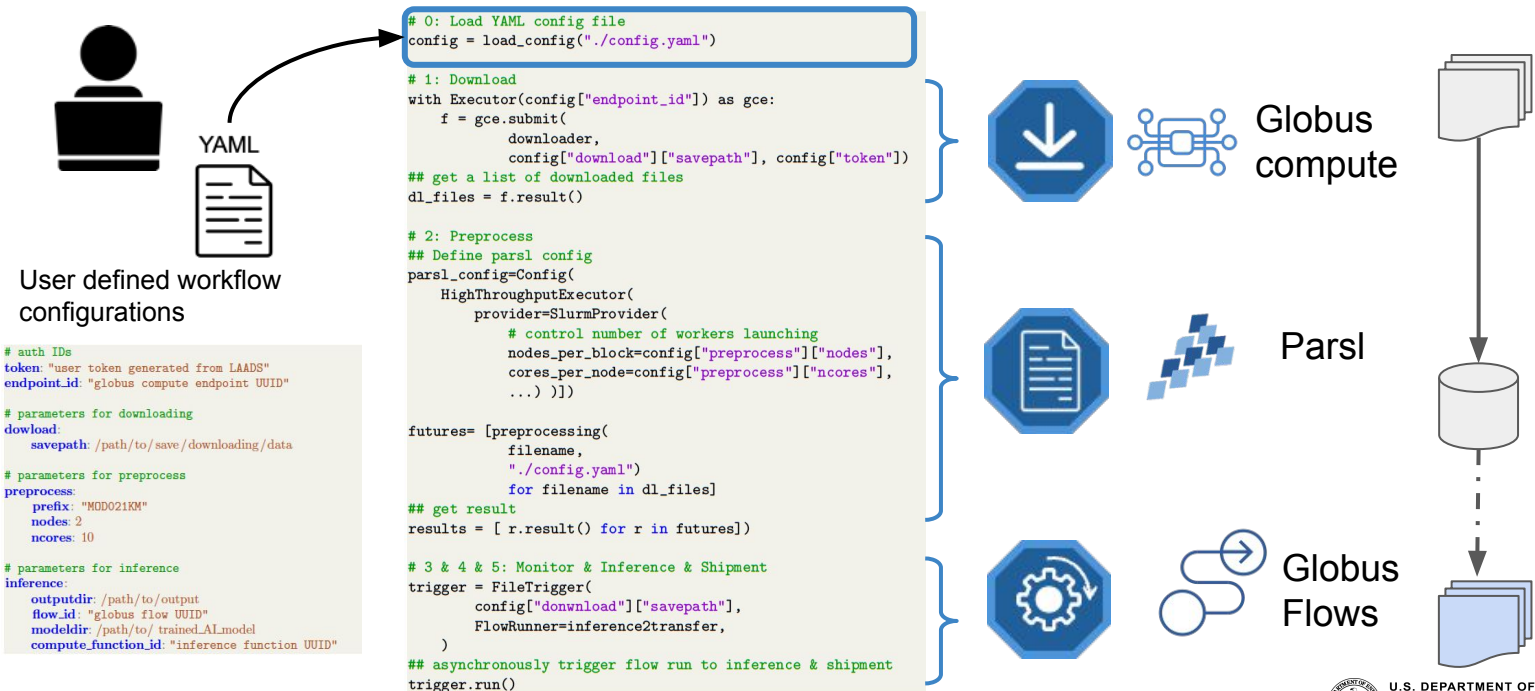
- 3 different satellite image products (~850TB) from NASA
- Downloading was remotely executed by **funcX**
- Preprocess was accelerated by **Parsl**

**Question:** How to connect data pipeline over multiple facilities?



# Use Globus Compute/Flow to easily create a seamless & flexibly multi-facility ML workflow for domain scientists

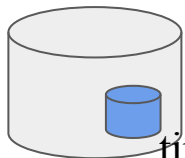
Previously disconnected ML pipelines are now integrated with Globus software tools ecosystem.



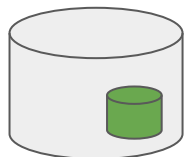
# Target application:

*Self-supervised deep learning → classify satellite clouds images*

Cloud image



Location



Physical variables

