



Scientific Gateways, Open OnDemand, Workflows

Karan Vahi

USC Information Sciences Institute



Cyberinfrastructure Training &
Education for
Synchrotron X-ray Science



X-CITE Workshop 2026



NSF OAC Award Numbers: 2320373, 2320375, 2320374



What Are Science Gateways?

- Community-developed platforms that give scientists, engineers, and students browser-based access to computing, data, and other research resources
- Provide a web-based interface that hides the complexity of the underlying cyberinfrastructure (CI) — researchers focus on the science, not on running HPC jobs or managing distributed data
- Gateways take many forms:
 - Web portals — submit and monitor jobs on HPC clusters
 - Domain-specific platforms (genomics, climate, materials science)
 - Workflow execution environments and data management platforms
- Examples: **Open OnDemand, Galaxy, ACCESS Pegasus, nanoHUB, CIPRES**





How Science Gateways Help Researchers

- **Ease of access**
 - Replace SSH clients, command-line use, scheduler syntax (SLURM, SGE, PBS), and manual scp transfers with intuitive browser interfaces — no software to install
- **Domain-specific interfaces**
 - Interfaces tailored to specific scientific workflows (e.g., protein structure analysis), hiding the complexity of running them on HPC clusters
- **Democratization of computing**
 - Opens advanced resources to researchers at smaller institutions, students and early-career researchers, domain scientists without a computational background, and occasional users of national resources
- **Reproducibility and collaboration**
 - Shared workflows, provenance tracking (what ran, when, and with what inputs), and community-contributed tools and pipelines





Open OnDemand: Browser-Based HPC Access

- Open-source, browser-based platform from the Ohio Supercomputer Center (OSC) — the standard web portal on most ACCESS resources
- Access HPC clusters with no client software: browse files, submit and monitor batch jobs, and launch interactive apps, all in your browser
- **Key features**
 - **File Manager** — browse the filesystem and upload, download, or edit files in the browser (no SFTP or scp)
 - **Job Composer** — form-based job creation with reusable templates; no scheduler syntax to memorize
 - **Interactive Apps** — launch Jupyter, RStudio, MATLAB, or a remote desktop on compute nodes
 - **Shell Access** — web terminal to the login node, no SSH client (handy with 2-factor login)





Open OnDemand on SDSC Expanse

Expansion Portal Apps Files Jobs Clusters Interactive Apps

SDS
Allocation and Usage Information
>_expansion Shell Access

The Expanse portal provides an integrated, and easy to use interface to access Expanse HPC resource.

With the portal, researchers can manage files (create, edit, move, upload, and download), view job templates for various applications, submit and monitor jobs, run interactive applications, and connect via SSH. The portal has no end-user installation requirements other than access to a modern up-to-date web browser

Pinned Apps A featured subset of all available apps

- Active Jobs System Installed App
- Home Directory System Installed App
- Job Composer System Installed App
- expansion Shell Access System Installed App
- MATLAB System Installed App
- RSTUDIO System Installed App
- Allocation and Usage Information System Installed App
- Jupyter System Installed App

Expansion Portal Apps Files Jobs Clusters Interactive Apps

Open in Terminal Refresh New File New Directory Upload Download Copy/Move Delete

Home Directory
Scratch

/ home / ux454545 / scratch / Change directory Copy path

Show Owner/Mode Show Dotfiles Filter: Showing 3 rows - 0 rows selected

Type	Name	Size	Modified at
Folder	karan	-	9/18/2025 1:06:43 PM
Folder	pegasuswfs	-	9/22/2025 1:50:53 PM
Folder	pilot.bUZ6bKCA	-	9/18/2025 7:13:59 AM

c2

c1





Launching a Jupyter Notebook via Open OnDemand – SDSC ExpansE

1. Go to the resource's portal URL and log in with your ACCESS credentials
2. Click Interactive Apps, then select Jupyter Notebook
3. Fill in the job form — account, partition (shared or compute), cores, memory, and walltime
4. Click Launch; wait for the job to move from Queued to Running
5. Click Connect to Jupyter to open the notebook in a new tab

Tip: interactive jobs consume allocation credits — request only what you need and close your session when done.

ACCESS Resource	Open OnDemand URL
Anvil (Purdue)	ondemand.anvil.rcac.purdue.edu
Bridges-2 (PSC)	ondemand.bridges2.psc.edu
Delta (NCSA)	login.delta.ncsa.illinois.edu
ExpansE (SDSC)	portal.expansE.sdsc.edu

portal.expansE.sdsc.edu/p/

Open OnDemand / Jupyter Session

Jupyter Session

Account:
TG-STA230005

Partition (Please choose the gpu, gpu-shared, or gpu-preempt as the partition if using gpus):
shared

Time limit (min):
30

Number of cores:
1

Memory required per node (GB):
2

GPUs (optional):
0

Singularity Image File Location: (Use your own or to include from existing container library at /cm/shared/apps/container e.g., /cm/shared/apps/containers/singularity/pytorch/pytorch-latest.sif)

Environment modules to be loaded (E.g., to use latest version of system Anaconda3 include cpu,gcc,anaconda3):

Conda Environment (Enter your own conda environment if any):

Conda Init (Provide path to conda initialization scripts)

Conda Yaml (Upload a yaml file to build the conda environment at runtime)
Choose File no file selected

Turn on use of mamba for speeding up conda-yml installs

Enable use of new caching mechanism that will store and reuse conda-yml created environments using conda-pack !!!!!

Reservation:

QoS:

Working directory:
home

Type:
JupyterLab

Submit

Open OnDemand: Job Composer

Expansive Portal / Job Composer Jobs Templates

Jobs

+ New Job ☆ Create Template

Edit Files Job Options Open Terminal Submit Stop Delete

Show 25 entries Search:

Created	Name	ID	Cluster	Status
No data available in table				

Showing 0 to 0 of 0 entries Previous Next

Build and submit jobs without the command line

The Job Composer builds jobs from reusable templates, then submits and monitors them — New Job, Edit Files, Submit / Stop, and a live status table, all in the browser.

Source: X-CITE — “Using Science Gateways” (xcitecourse.org)



Open OnDemand: Anvil Demo

Login: <https://ondemand.anvil.rcac.purdue.edu>





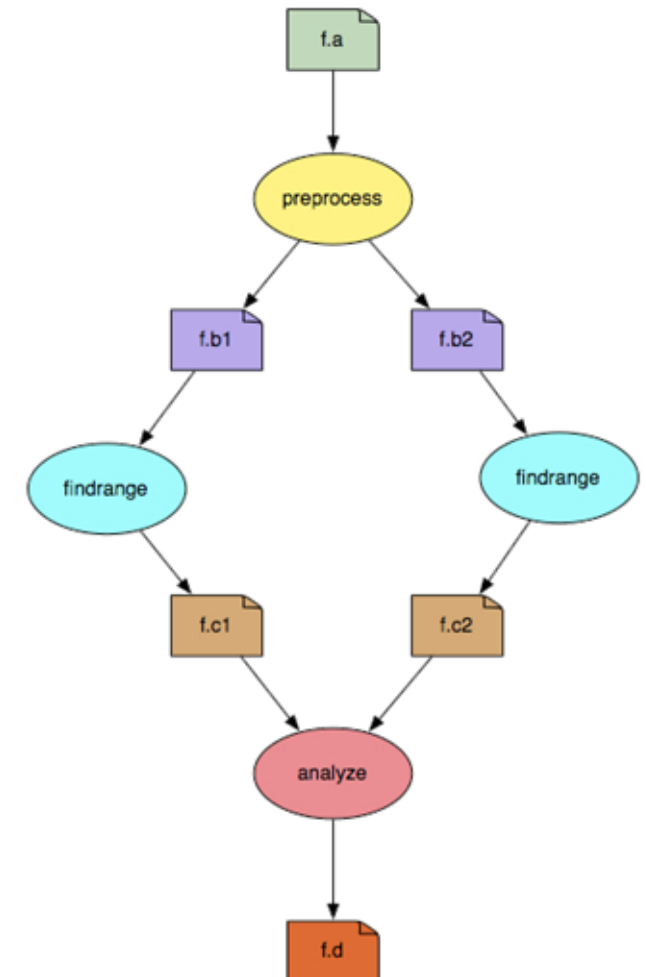
ACCESS Pegasus: A Gateway for Workflows

- A ready-to-use environment for running **Pegasus Workflows** on ACCESS resources — no need to install or configure Pegasus yourself
- Built on Open OnDemand and maintained by the Pegasus team at USC's Information Sciences Institute (USC/ISI)
- What it provides:
 - A pre-configured Pegasus submit host connected to ACCESS
 - Direct connectivity to Expanse, Anvil, Bridges-2, and other ACCESS resources, plus the OSG OSPool for high-throughput workloads
 - Notebooks run on pegasus.access-ci.org while jobs execute on HTCondor points — a clean separation of notebook and compute environments
- **Why use it:** zero configuration, multi-site execution, fault tolerance (retry and checkpointing), full provenance, and expert support from the Pegasus team



Scientific Workflows

- An abstraction to express ensemble of complex computational operations
 - Eg: retrieving data from remote storage services, executing applications, and transferring data products to designated storage sites
- A workflow is represented as a directed acyclic graph (DAG)
 - Nodes: tasks or jobs to be executed
 - Edges: depend between the tasks
- Have a monolithic application/experiment?
- The tasks in a scientific workflow can be everything from short serial tasks to very large parallel tasks (MPI for example) surrounded by a large number of small, serial tasks used for pre- and post-processing.
- Find the inherent DAG structure in your application to convert into a workflow



Workflow Challenges Across Domains

- Describe complex workflows in a simple way
- Access distributed, heterogeneous data and resources (heterogeneous interfaces)
- Deal with resources/software that change over time
- Ease of use. Ability to debug and monitor large workflows

Our Focus

- Separation between workflow description and workflow execution
- Workflow planning and scheduling (scalability, performance)
- Task execution (monitoring, fault tolerance, debugging, web dashboard)
- Provide additional assurances that a scientific workflow is not accidentally or maliciously tampered with during its execution.



Pegasus Workflow Management System



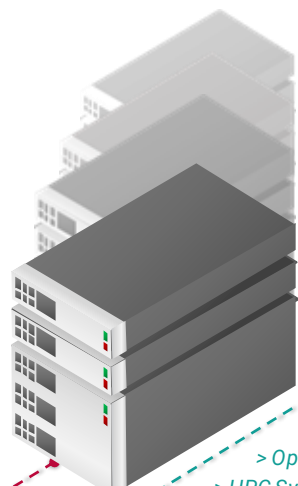
Pegasus WMS

Planner

Monitoring & Provenance

Engine

Scheduler



- > Cloud Resources
- > Open Science Grid
- > HPC Systems
- > HTCondor Pools

Submit Node

Compute Resources

End to End Workflow Management & Execution

- ▶ Develop portable scientific workflows in Python, Java, and R
- ▶ Compile workflows to be run on heterogeneous resources
- ▶ Monitor and debug workflow execution via CLI and web-based tools
- ▶ Recover from failures with built-in fault tolerance mechanisms
- ▶ Regular release schedule incorporating latest research and development

2001	2003	2005	2007	2009	2011	2013	2015	2017	2018	2020												
1.0	1.1	1.2	1.3	1.4	2.0	2.1	2.2	2.3	2.4	3.0	3.1	4.0	4.1	4.2	4.3	4.4	4.5	4.6	4.7	4.8	4.9	5.0
Development			support for GT4	task clustering	support for AWS	hierarchical workflows	pegasus-lite engine	monitoring dashboard	ensemble manager	support for containers	redesign of APIs											
Research LIGO, SCEC, and others			data cleanup algorithms	data footprint	cloud computing evaluation	MPI-based workflow engine design	Real time performance data capture	metadata capture	data integrity assurance													

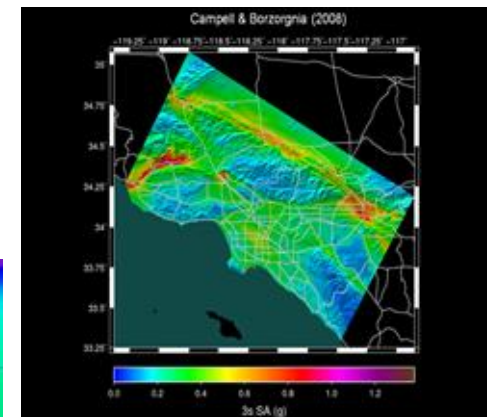
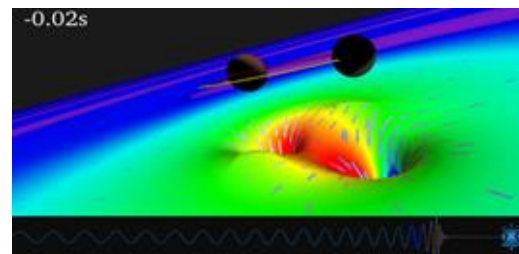
Pegasus in practice

- ▶ Laser Interferometer Gravitational Wave Observatory (LIGO) develops large scale analysis pipelines used for gravitational wave detection.
- ▶ Southern California Earthquake Center (SCEC) CyberShake project generates hazard maps using hierarchical workflows .
- ▶ The XENONnT project uses Pegasus for processing and monte carlo workflows, searching for dark matter

The XENONnT detector



LIGO observation of colliding black holes



Hazard map indicating maximum amount of shaking at a particular geographic location generated from SCEC's CyberShake Pegasus workflow

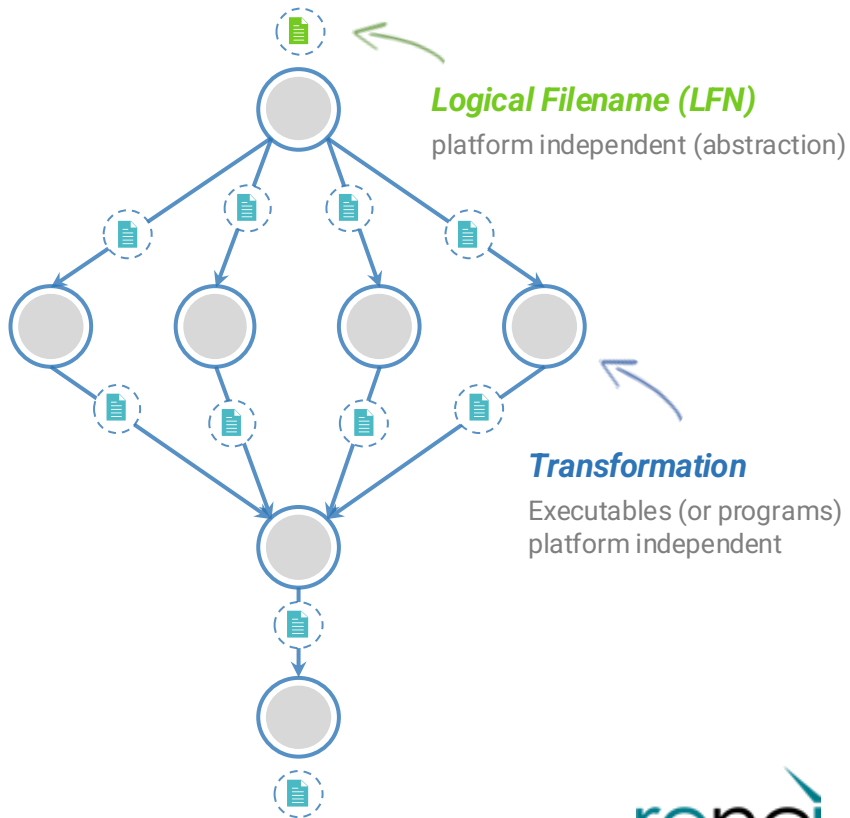


Input Workflow Specification **YAML formatted**

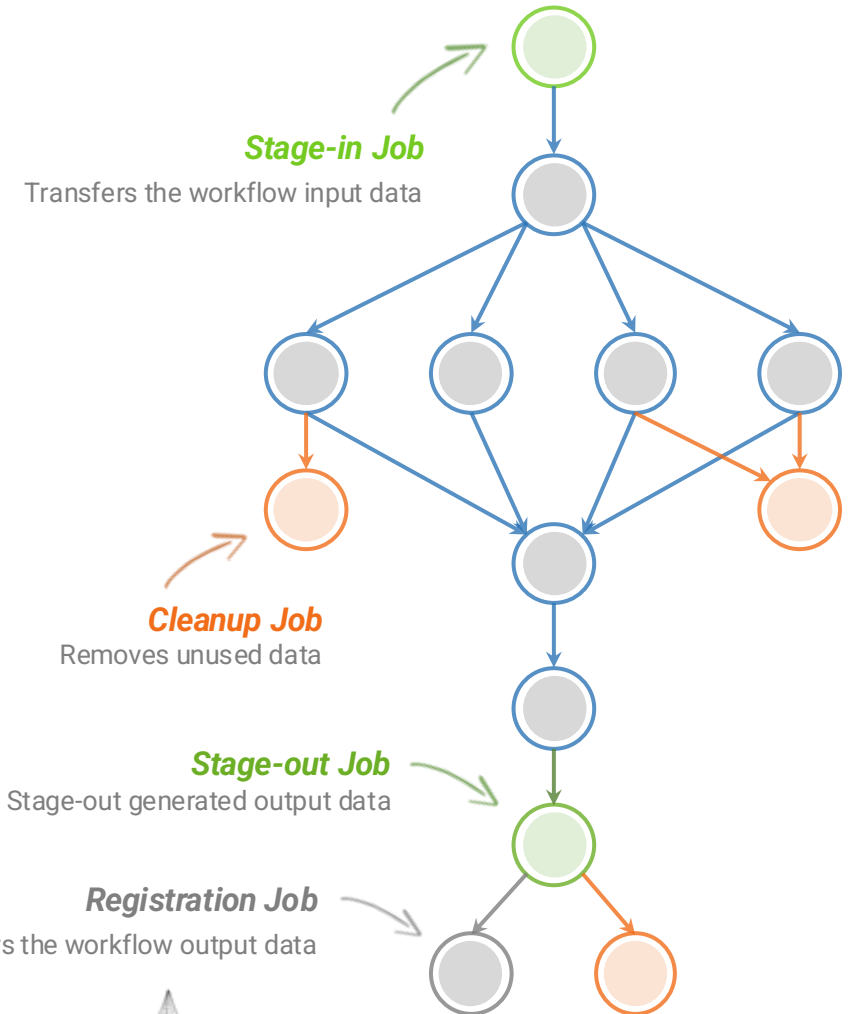
Portable Description

Users do not worry about low level execution details

ABSTRACT WORKFLOW



Output Workflow



EXECUTABLE WORKFLOW

Pegasus CHESS QM Workflow

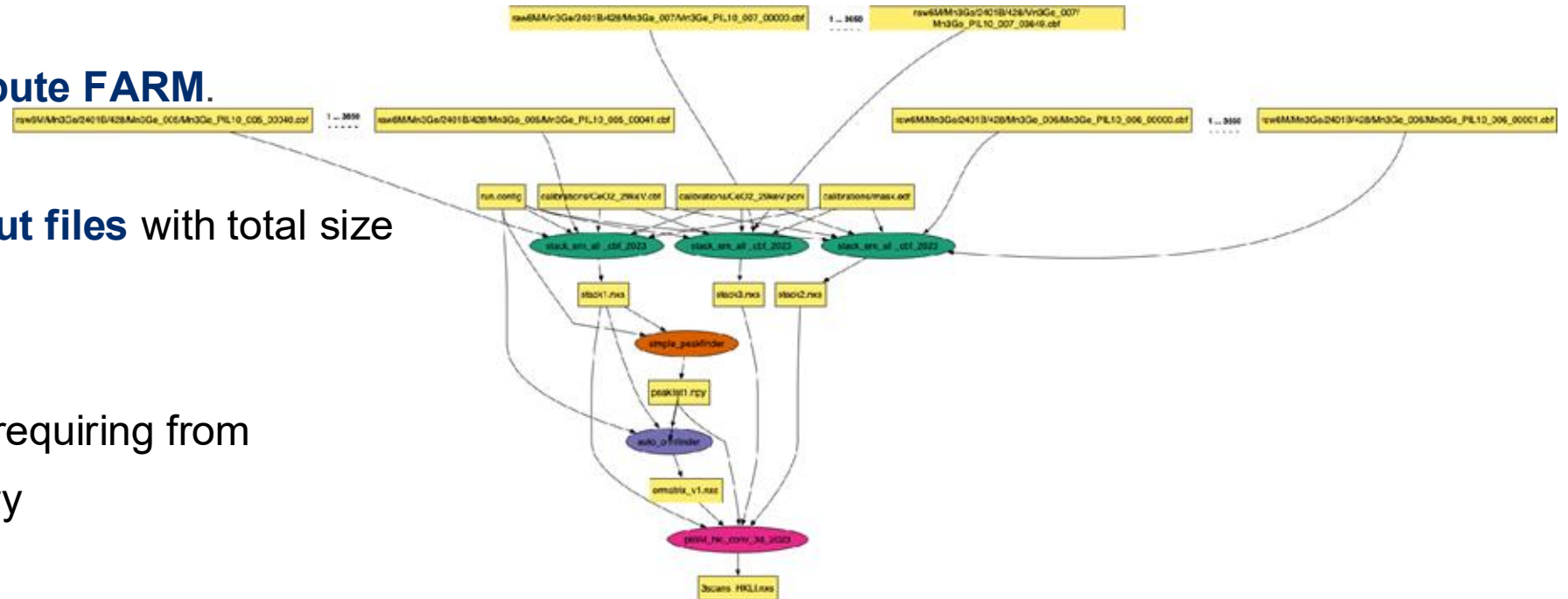
GitHub Repository: <https://github.com/pegasus-isi/chess-qmb-workflow>

Runs on the **CHESS Compute FARM**.

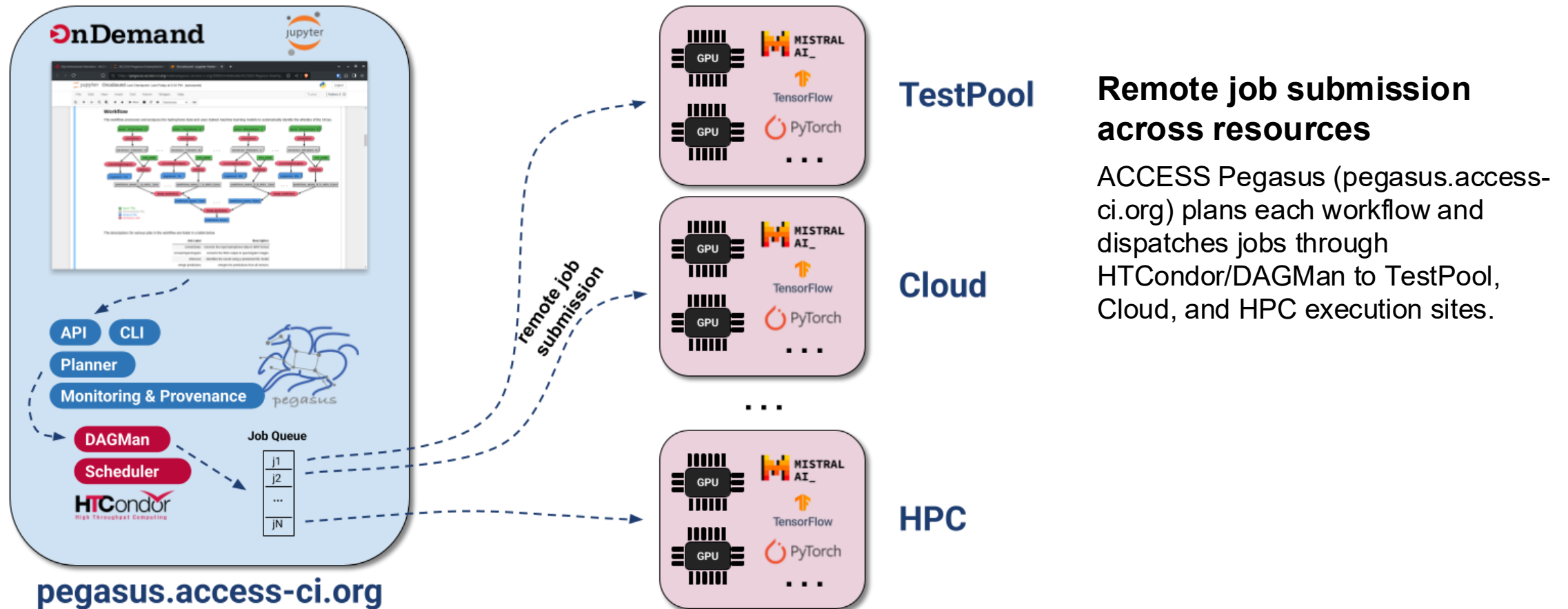
Requires about **11,000 input files** with total size of approximately **570GB**.

Mainly **high memory** jobs requiring from **10GB to 350 GB** of memory

The **pil6M_hkl_conv** job requires 56 cores.



ACCESS Pegasus: Remote Job Submission



Remote job submission across resources

ACCESS Pegasus (pegasus.access-ci.org) plans each workflow and dispatches jobs through HTCondor/DAGMan to TestPool, Cloud, and HPC execution sites.



Getting Started with ACCESS Pegasus

1. **Get an ACCESS account** — register at <https://operations.access-ci.org/identity/new-user>
2. **Get an allocation (optional)** — not required for the training notebooks; start with an EXPLORE allocation for production work (see the DC200 module)
3. **Try Pegasus** — visit <https://support.access-ci.org/tools/pegasus> and log in with your ACCESS ID
4. **Run the examples** — open a Jupyter notebook and work through the ACCESS-Pegasus-Examples directory (01-Quickstart through 09-Tutorial-SharedFS)

Get support:

- Email: pegasus-support@isi.edu
- Slack: pegasus-users.slack.com
- ACCESS Support Portal: <https://support.access-ci.org>





Galaxy: Web-Based Analysis Gateway

- Open-source, browser-based gateway — originally for bioinformatics, now used across many domains including X-ray science, genomics, and climate modeling
- Build, run, and share complex data-analysis workflows entirely in a browser — no programming or command-line experience required (use public usegalaxy.org or a local instance)
- **Key features:**
 - **Tool Integration** — a large, admin-managed library of analysis tools, with no software installs or version conflicts
 - **Workflow Builder** — a visual editor that connects tools on a canvas, encoding data dependencies automatically
 - **History** — records every dataset, tool, and parameter as a shareable, exportable audit trail
 - **Data Import & ToolShed** — import from uploads, URLs, shared libraries, and services like NCBI and Globus; reuse community tools via the Galaxy ToolShed
- **At CHESS:** a Galaxy instance configured for X-ray data analysis that submits jobs to the CHESS SGE cluster





Summary

- Science gateways — Open OnDemand, Galaxy, and ACCESS Pegasus — sharply lower the barriers to using advanced computing resources
- They span everything from interactive analysis in a Jupyter notebook on a national HPC cluster to orchestrating complex, multi-step workflows across distributed resources — with no deep HPC system-administration expertise required
- **For CHESSE researchers:**
 - **Open OnDemand** — the easiest way to get started on ACCESS HPC: open a browser, log in, and launch a Jupyter notebook on the cluster
 - **ACCESS Pegasus** — the best path for running Pegasus workflow pipelines on national CI without the setup overhead
 - Both work seamlessly with the ACCESS allocation described in the DC200 module





Thank you!

XCITE Gateways Module: <https://xcitecourse.org/theme3/DC102/using-science-gateways.html>

XCITE Workflows Module: <https://xcitecourse.org/theme3/DC101/scientific-workflow-management.html>

Pegasus Website: <https://pegasus.isi.edu>



National
Science
Foundation

NSF Cybertraining Award # 2320373

