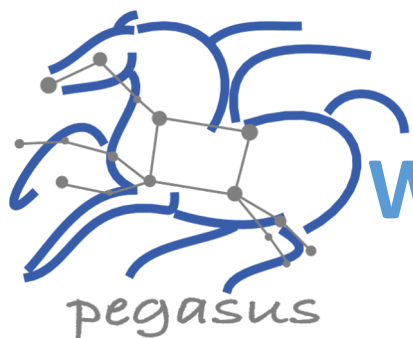


eScience '22 | DEMOCRATIZING SCIENCE



Pegasus 5.0 Workflows with Containers

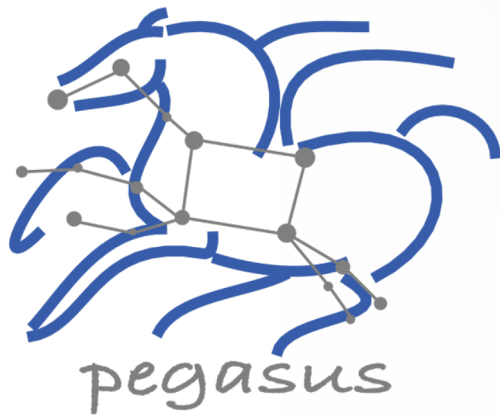
Karan Vahi, Mats Rynge

Information Sciences Institute
University of Southern California, School of Engineering
vahi@isi.edu , rynge@isi.edu



U.S. DEPARTMENT OF
ENERGY





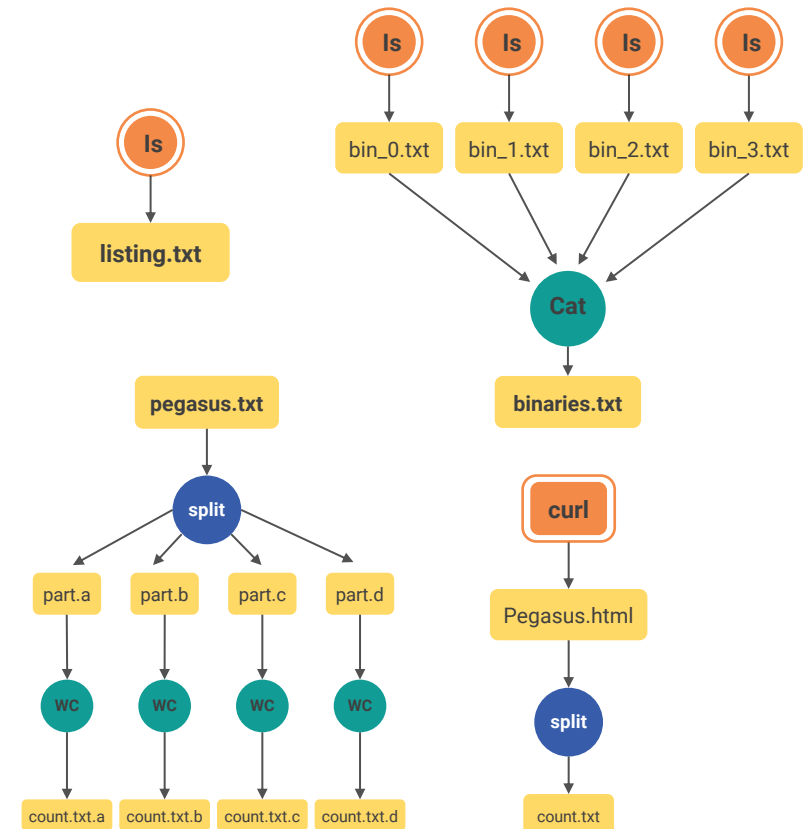
1. Introduction



What are Scientific Workflows

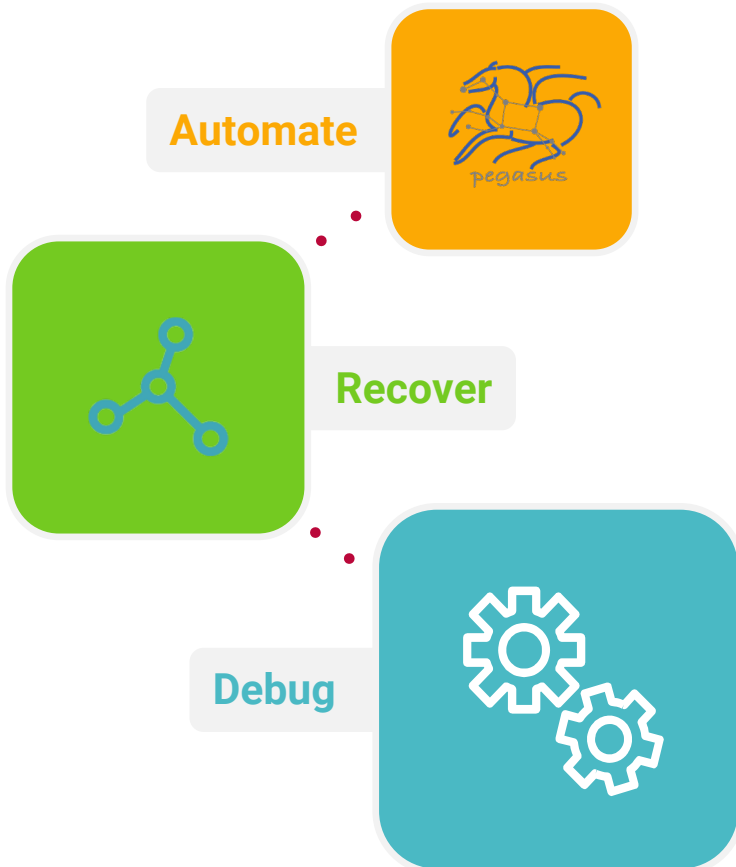
- ▲ **Conducts a series of computational tasks.**
 - Resources distributed across Internet.
- ▲ **Chaining (outputs become inputs) replaces manual hand-offs.**
 - Accelerated creation of products.
- ▲ **Ease of use - gives non-developers access to sophisticated codes.**
 - Resources distributed across Internet.
- ▲ **Provides framework to host or assemble community set of applications.**
 - Honors original codes. Allows for heterogeneous coding styles.
- ▲ **Framework to define common formats or standards when useful.**
 - Promotes exchange of data, products, codes. Community metadata.
- ▲ **Multi-disciplinary workflows can promote even broader collaborations.**
 - E.g., ground motions fed into simulation of building shaking.
- ▲ **Certain rules or guidelines make it easier to add a code into a workflow.**

Workflow Building Blocks



Slide Content Courtesy of David Okaya, SCEG, USC

Why Pegasus?



- ▶ **Automates Complex**, Multi-stage Processing Pipelines
- ▶ Enables Parallel, **Distributed Computations**
- ▶ **Automatically Executes** Data Transfers
- ▶ Reusable, Aids **Reproducibility**
- ▶ Records How Data was Produced (**Provenance**)
- ▶ Handles **Failures** with to Provide Reliability
- ▶ Keeps Track of Data and **Files**
- ▶ Ensures **Data Integrity** during workflow execution



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

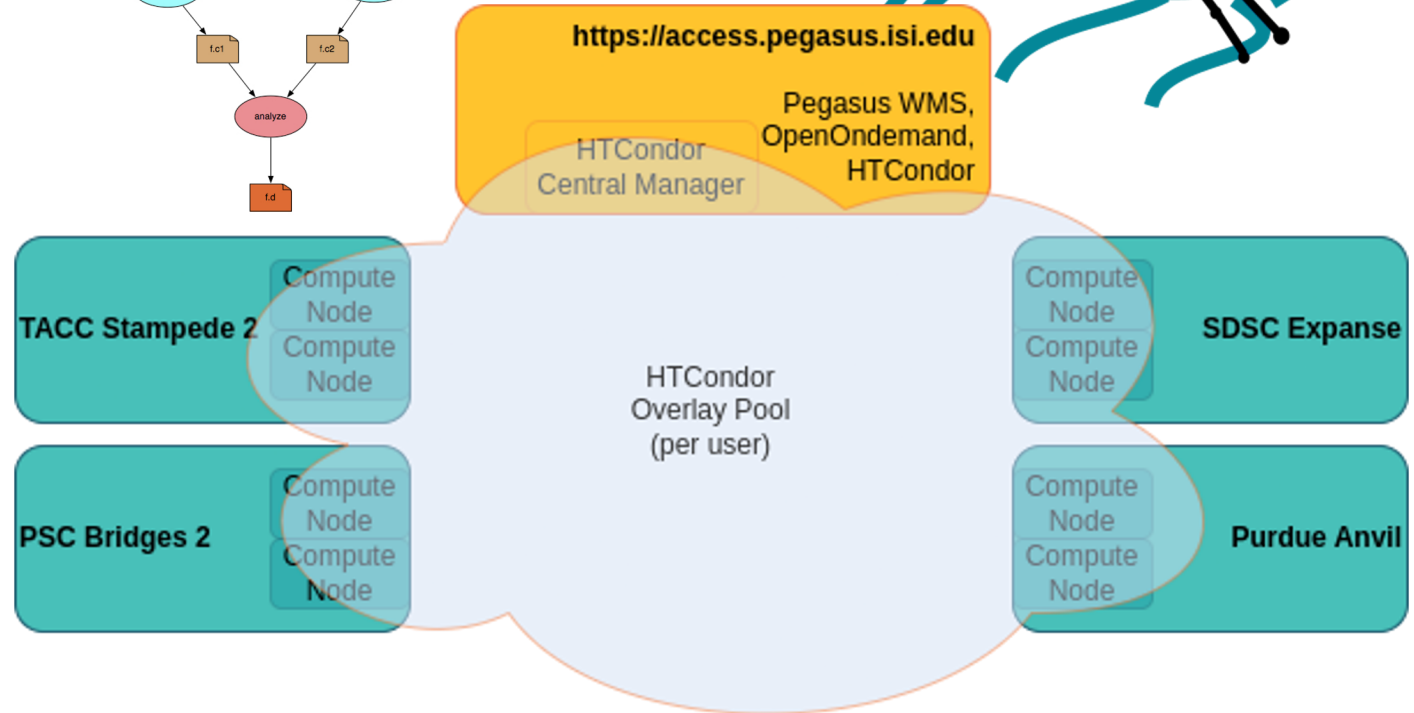
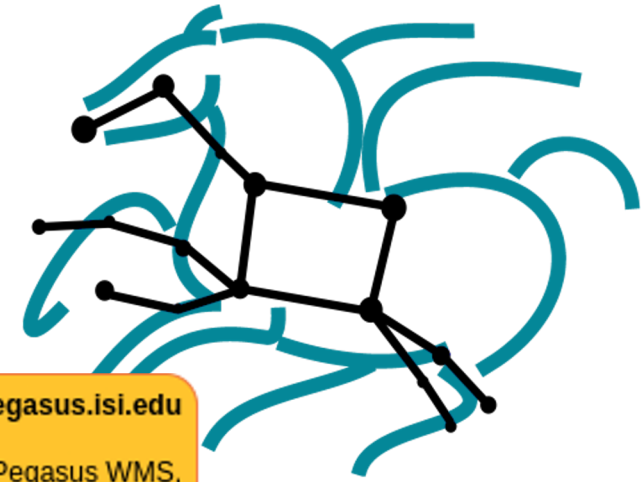
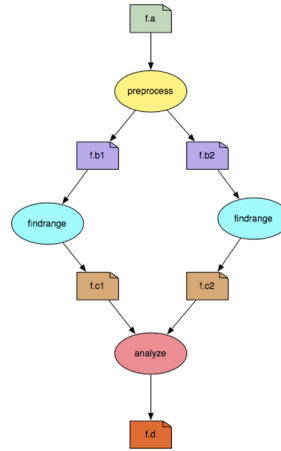
Our Focus

- ▶ 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.

ACCESS Pegasus

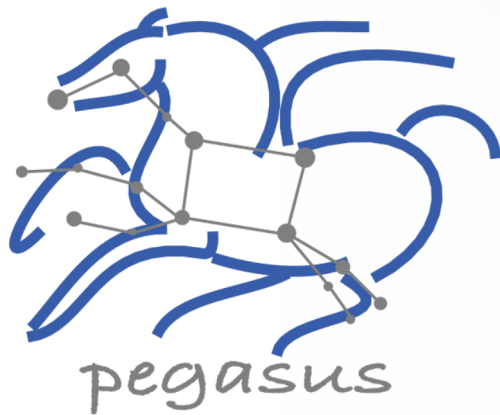
Bring your workflows to ACCESS!

- Execute scientific workflows across ACCESS resources
- OpenOnDemand Portal: **has all you need**: Jupyter Notebooks, ACCESS authentication, Pegasus workflow management, and HTCondor job management
- **Bring your own ACCESS capacity**: HTCondor Annex - pilot jobs automatically create a virtual HTCondor pool



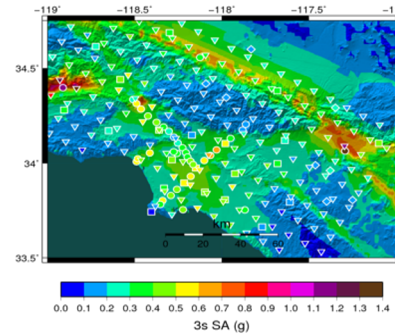
<http://access.pegasus.isi.edu>

More at: support.access-ci.org/pegasus



Some of The Success Stories...

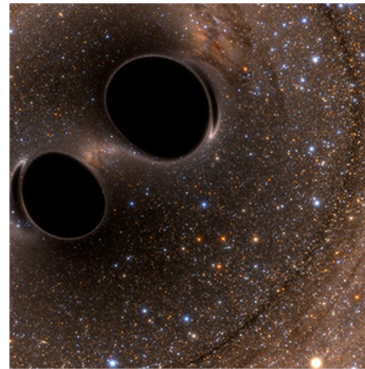
Southern California Earthquake Center's CyberShake



First Physics-Based "Shake map" of Southern California

Mix of MPI and single-core jobs, mix of CPU, GPU codes.
Large data sets (10s of TBs), ~300 workflows with 420,000 tasks each
Supported since 2005: changing CI, x-platform execution

Laser Interferometer Gravitational-Wave Observatory (LIGO)



First direct detection of a gravitational wave (colliding black holes)

High-throughput computing workload, access to HPC resources, ~ 21K Pegasus workflows, ~ 107M tasks

Supported since 2001, distributed data, opportunistic computing resources

XENONnT - Dark Matter Search



Custom data management
Rucio for data management
MongoDB instance to track science runs and data products.

Monte Carlo simulations and the main processing pipeline.

Southern California Earthquake Center's CyberShake



Builders ask seismologists:

What will the peak ground motion be at my new building in the next 50 years?



Seismologists answer this question

using Probabilistic Seismic Hazard Analysis (PSHA)

CPU jobs
(Mesh generation, seismogram synthesis)
1,094,000 node-hours



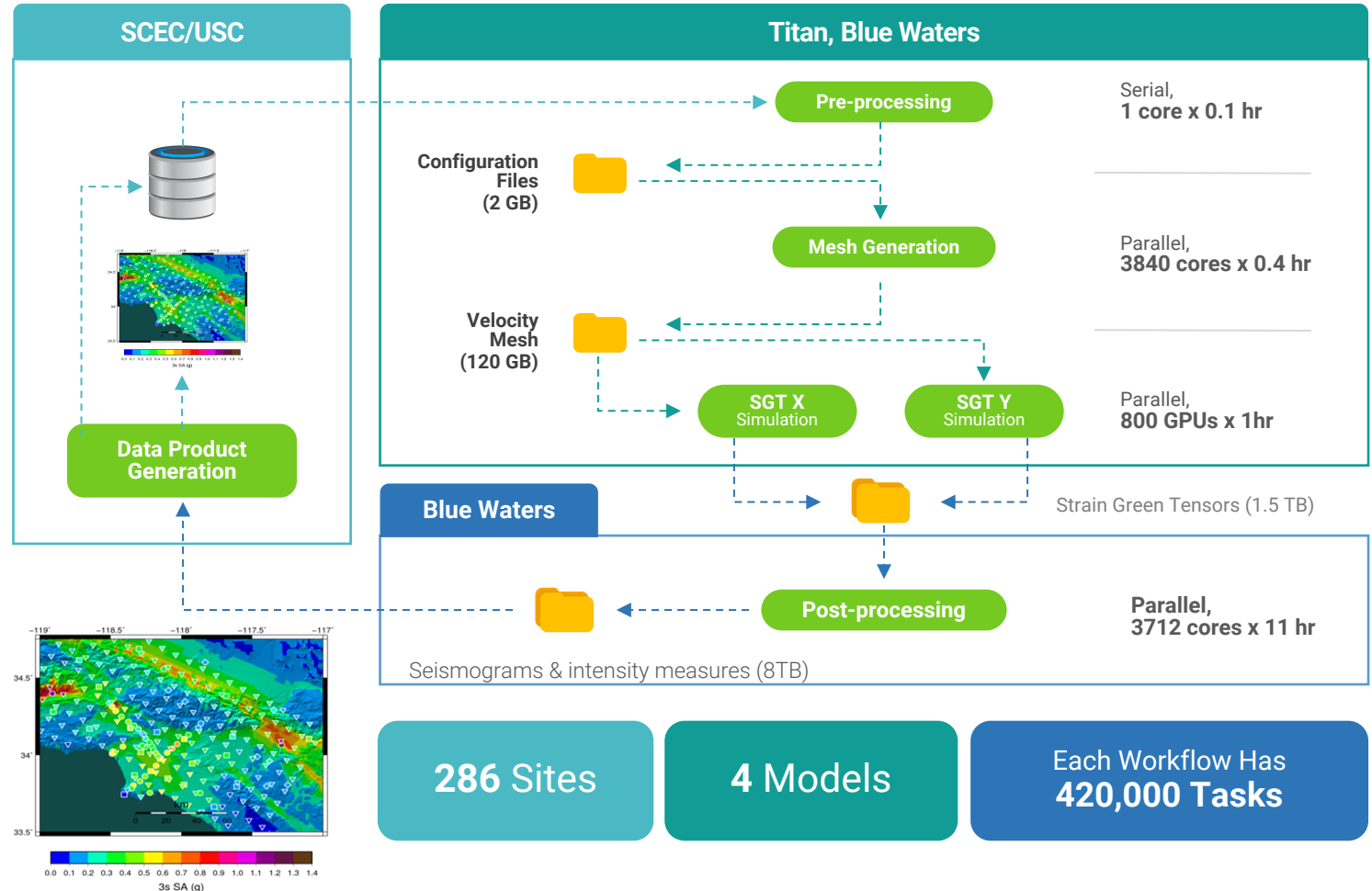
GPU jobs:
439,000 node-hours
AWP-ODC finite-difference code
5 billion points per volume, 23,000 timesteps
200 GPUs for 1 hour



Titan:
421,000 CPU node-hours, 110,000 GPU node-hours



Blue Waters:
673,000 CPU node-hours, 329,000 GPU node-hours



286 Sites

4 Models

Each Workflow Has 420,000 Tasks



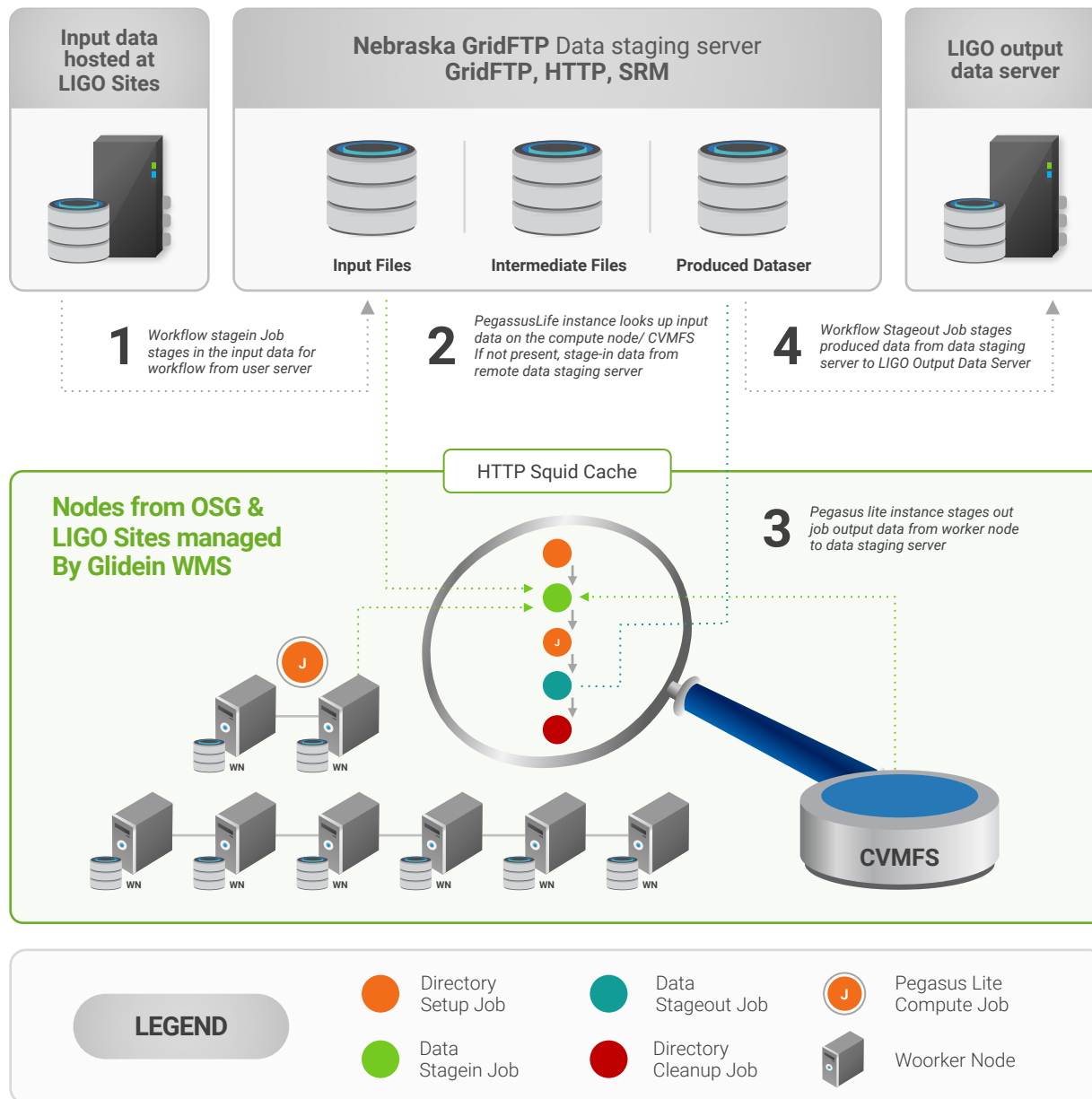
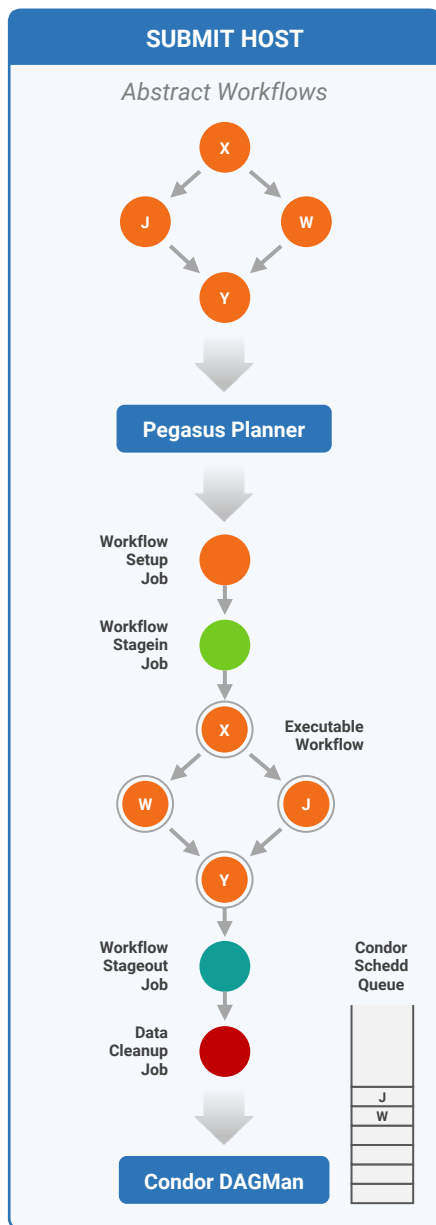
Data Flow for LIGO Pegasus Workflows in OSG

Advanced LIGO Laser Interferometer Gravitational Wave Observatory



60,000 Compute Tasks
 Input Data: 5000 files (10GB total)
 Output Data: 60,000 files (60GB total)
 Processed Data: 725 GB

Executed on LIGO Data Grid, EGI, Open Science Grid and XSEDE



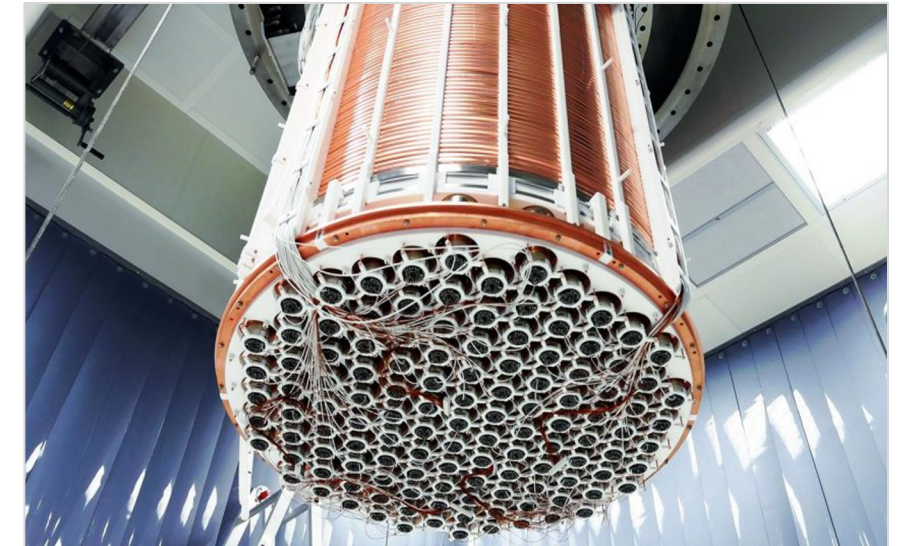
XENONnT - Dark Matter Search



Two Workflows

Monte Carlo simulations and the main processing pipeline.

- Workflows execute across Open Science Grid (OSG) & European Grid Infrastructure (EGI)
- Rucio for data management
- MongoDB instance to track science runs and data products.

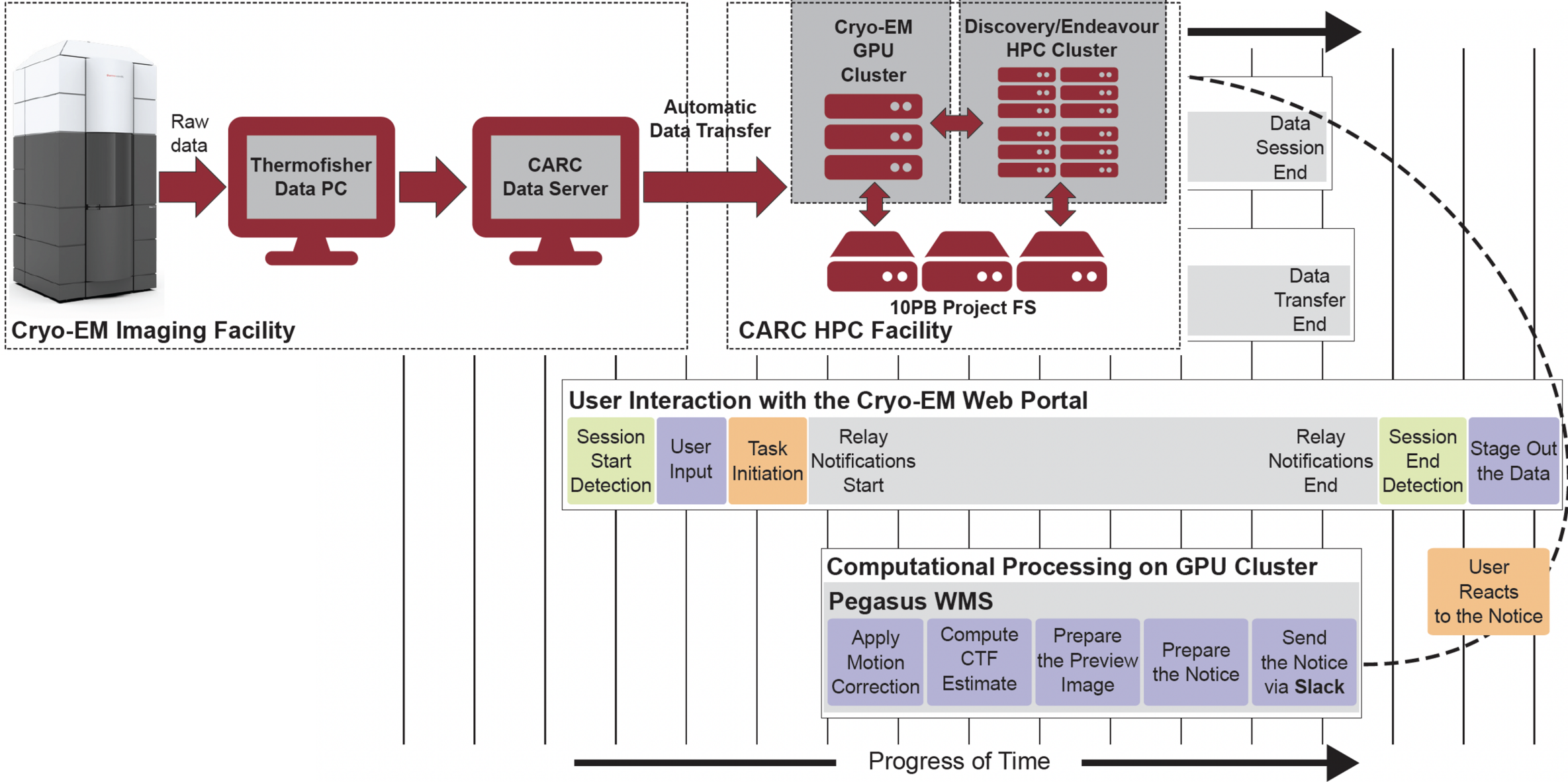


Type	Succeeded	Failed	Incomplete	Total	Retries	Total+Retries
Tasks	4000	0	0	4000	267	4267
Jobs	4484	0	0	4484	267	4751
Sub-Workflows	0	0	0	0	0	0

Workflow wall time	: 5 hrs, 2 mins
Cumulative job wall time	: 136 days, 9 hrs
Cumulative job wall time as seen from submit side	: 141 days, 16 hrs
Cumulative job badput wall time	: 1 day, 2 hrs
Cumulative job badput wall time as seen from submit side	: 4 days, 20 hrs

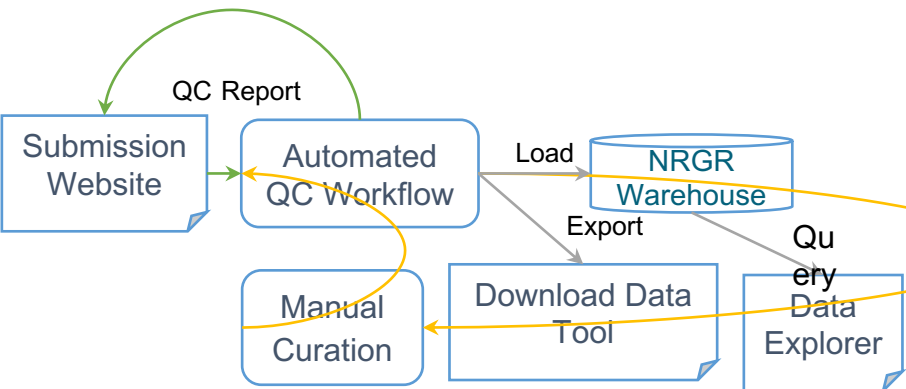
Main processing pipeline is being developed for XENONnT - data taking will start at the end of 2019. Workflow in development:

Processing instrument data in real time





The NIMH Center for Collaborative Genomic Studies on Mental Disorders, now known as the NIMH Repository and Genomics Resource (NRGR), maintains biomaterials, demographic, and phenotypic data from over 200,000 well-characterized individuals with a range of psychiatric illnesses, their family members, and unaffected controls.

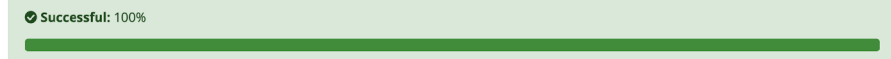


- **Easy to Use Web-Based Interface**
 - Simple Submission
 - Real-time Monitoring and Error Reports
 - After automated QC, submit corrected files for expert curation
- **Scalable**
 - Workflow based architecture using Pegasus WMS
- **Extensible Design**
 - Easily add new QC steps, and checks
- **Enables Complex checks**
 - Pedigree Checks
 - QC Checks validating data with external sources
 - QC Checks can correlate data across multiple files and across multiple fields within files
- Ensures high-quality uniform data deposited at NRGR
- Better resource utilization: solve most QC problems automatically, use expert curation for hard cases

Auto QC Status

New Validation [Help](#)

[Back to Previous Validations](#)



Summary

UID	5e6a6ddd95f6e
Disorder	Depression
Study Id	149
File	shaptest7.zip
User	JaclynVitanza
Email	jv607@dls.rutgers.edu
Started On	Mar 12, 2020 10:14 AM
Workflow Directory	/web/data/qc/runs/5e6a6ddd95f6e

Sanity Check Status

[Download All Files](#)

File	Submission Validation	Pedigree Validation
study_149_sub.csv	Standardized File Log	Log Log
File	ID Validation	
study_149_id.csv	Standardized File Log	Log Log
File	Phenotypic Validation	
shaps01_phen.csv	Standardized File Log	Log Log
File	Advanced QC	
study_149_sub.canon.csv	Corrected Submission File Log	
study_149_id.canon.csv	Corrected ID File Log	
Corrections Log	Corrections Log Log	
Advanced QC Report	Advanced QC Report Log	

Validate with AutoQC

[Previous Validations](#) [Help](#)

OVERVIEW HOW TO VALIDATE AND SUBMIT DATA SUBMISSION REQUIREMENTS **VALIDATE WITH AUTOQC**

Validate your data for sanity checks and quality control.

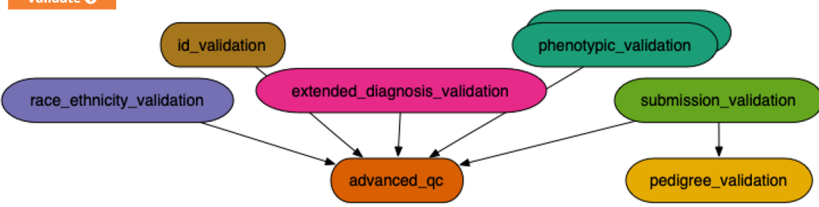
Choose File [Browse](#)

What data are you submitting?
-- Choose a Disorder --

Study Id
256

Email Notification
email@address.com

[Validate](#)



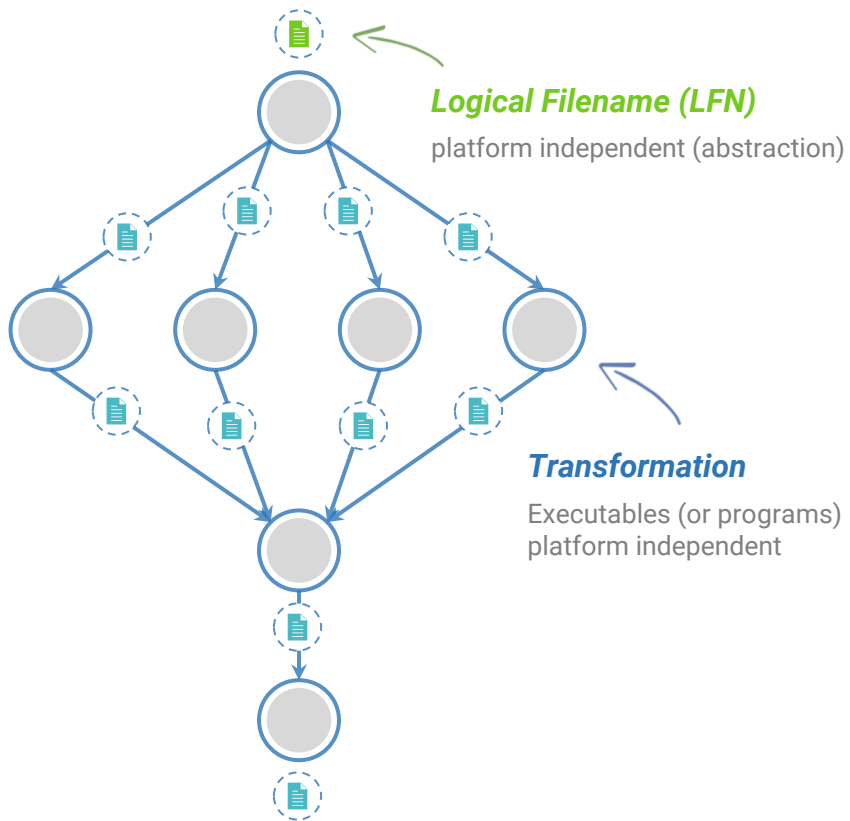


Input Workflow Specification **YAML formatted**

Portable Description

Users do not worry about low level execution details

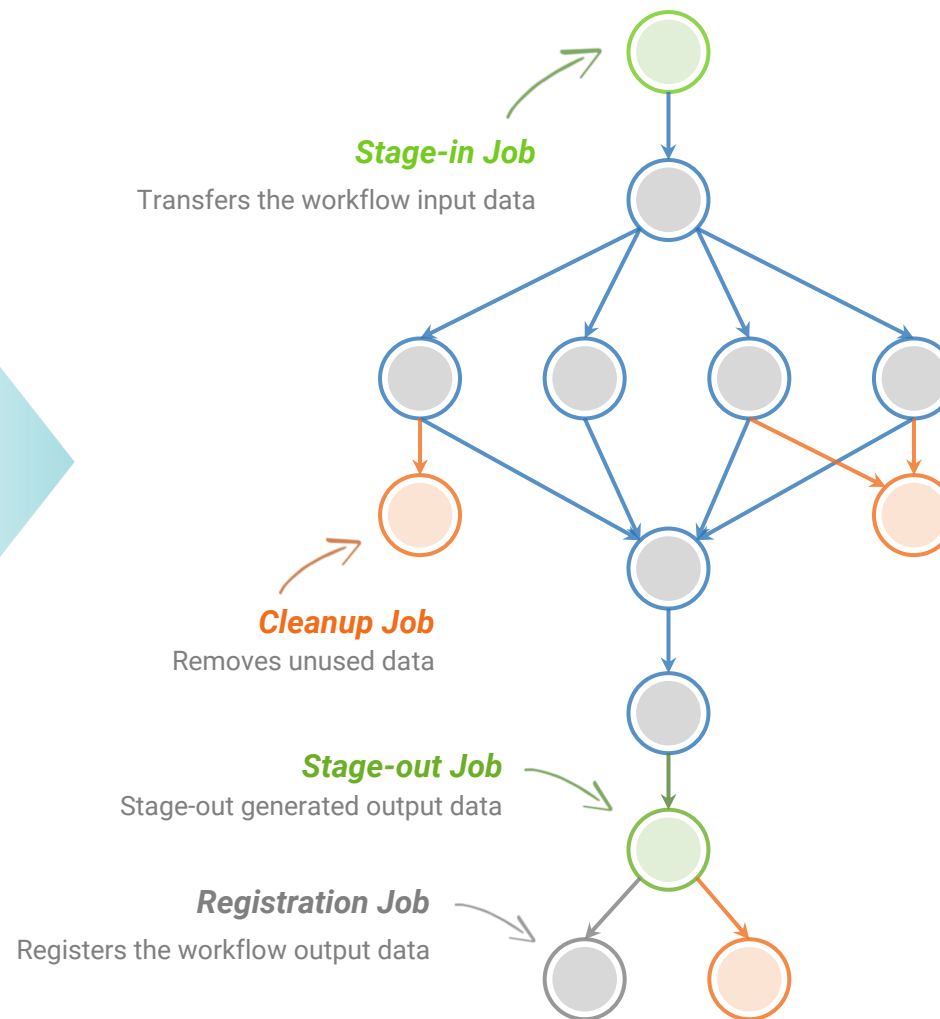
ABSTRACT WORKFLOW



directed-acyclic graphs

Output Workflow

EXECUTABLE WORKFLOW



Pegasus Deployment



Workflow Submit Node

- Pegasus WMS
- HTCondor

One or more Compute Sites

- Compute Clusters
- Cloud
- OSG

Input Sites

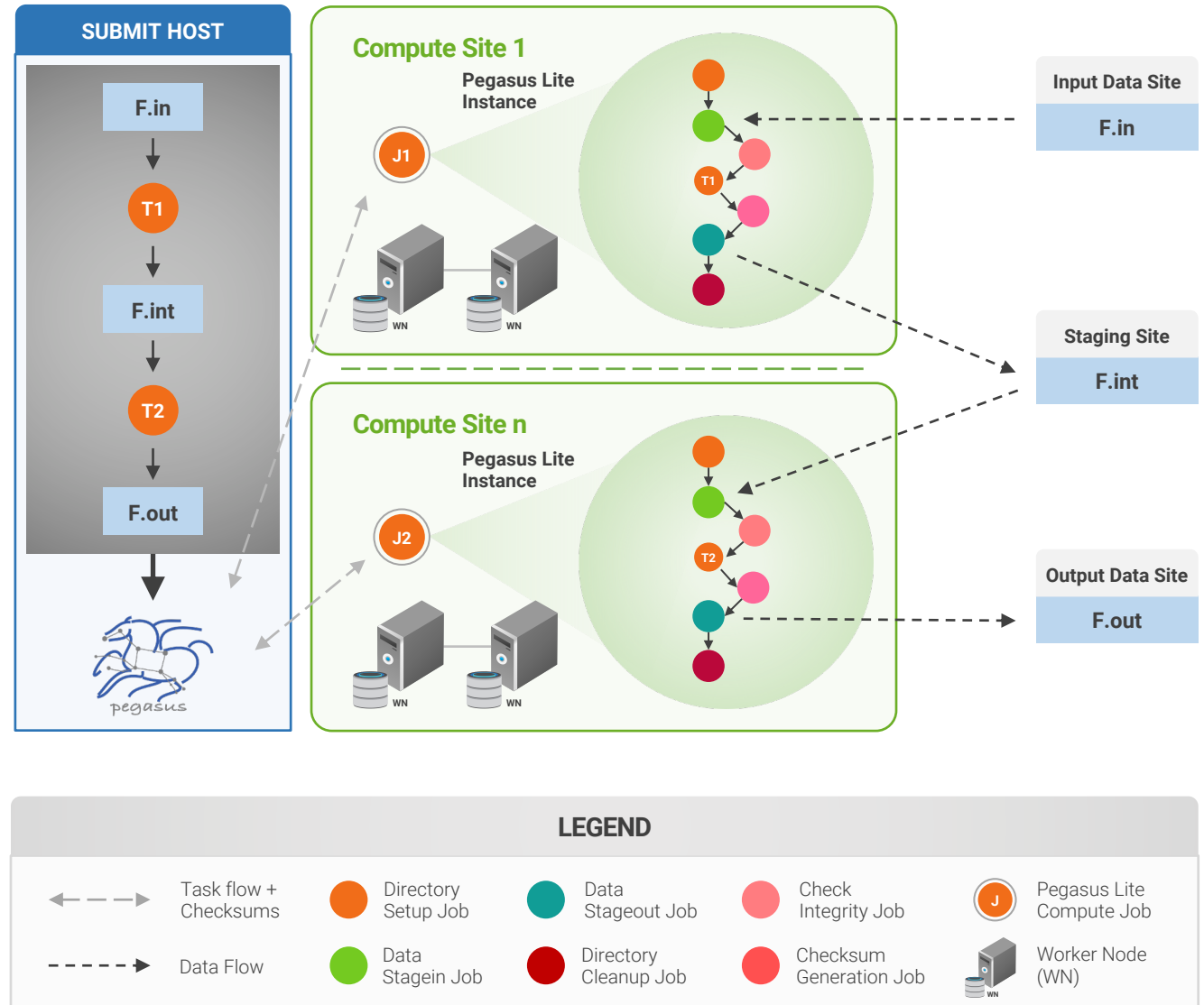
- Host Input Data

Data Staging Site

- Coordinate data movement for workflow

Output Site

- Where output data is placed



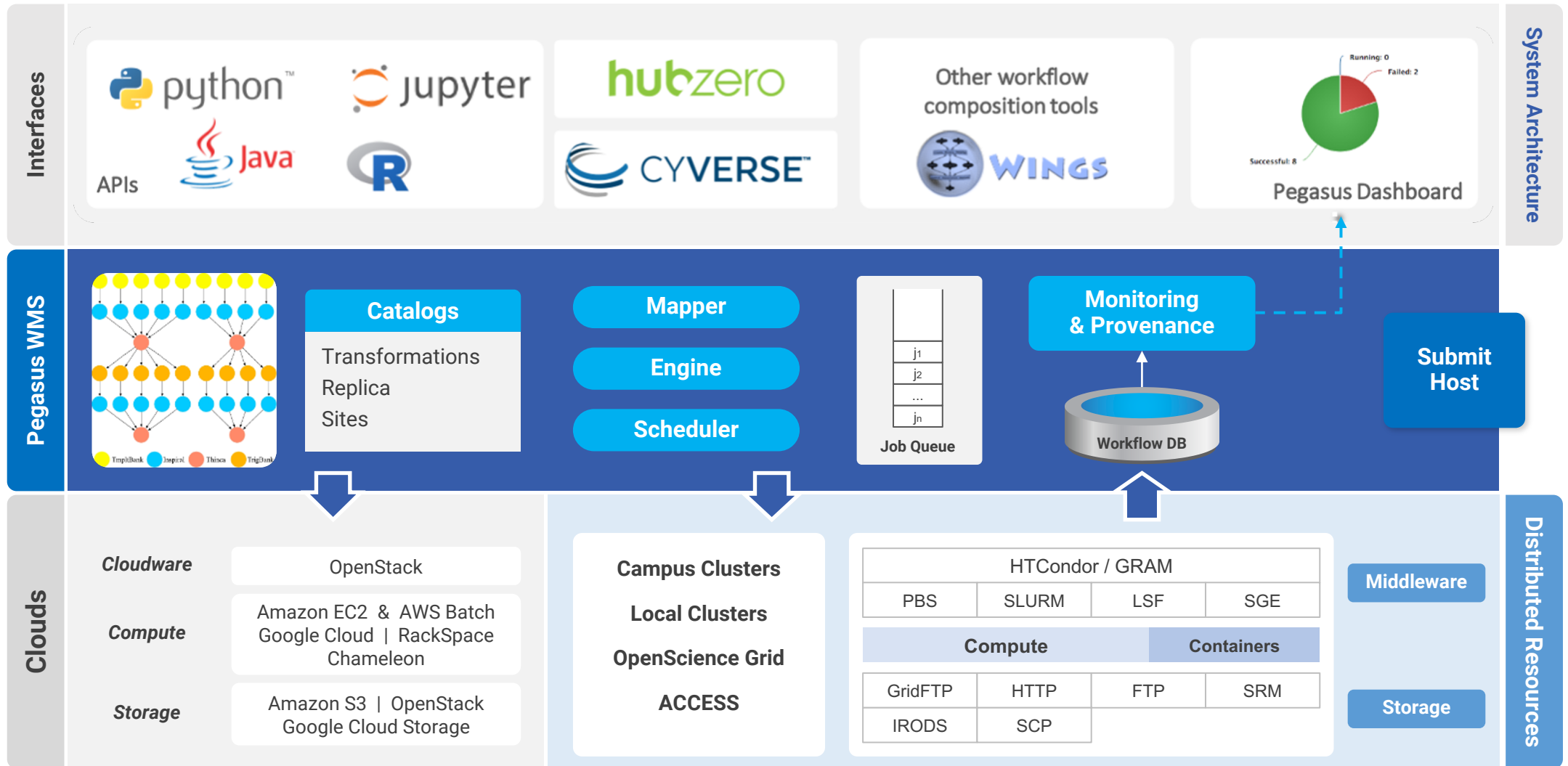


Pegasus-transfer

Pegasus' internal data transfer tool with support for a number of different protocols

- **Directory creation, file removal**
 - If protocol can support it, also used for cleanup
- **Two stage transfers**
 - e.g., GridFTP to S3 = GridFTP to local file, local file to S3
- **Parallel transfers**
- **Automatic retries**
- **Credential management**
 - Uses the appropriate credential for each site and each protocol (even 3rd party transfers)

```
HTTP
SCP
GridFTP
Globus
Online
iRods
Amazon S3
Google
Storage
SRM
FDT
Stashcp
Rucio
cp
ln -s
```



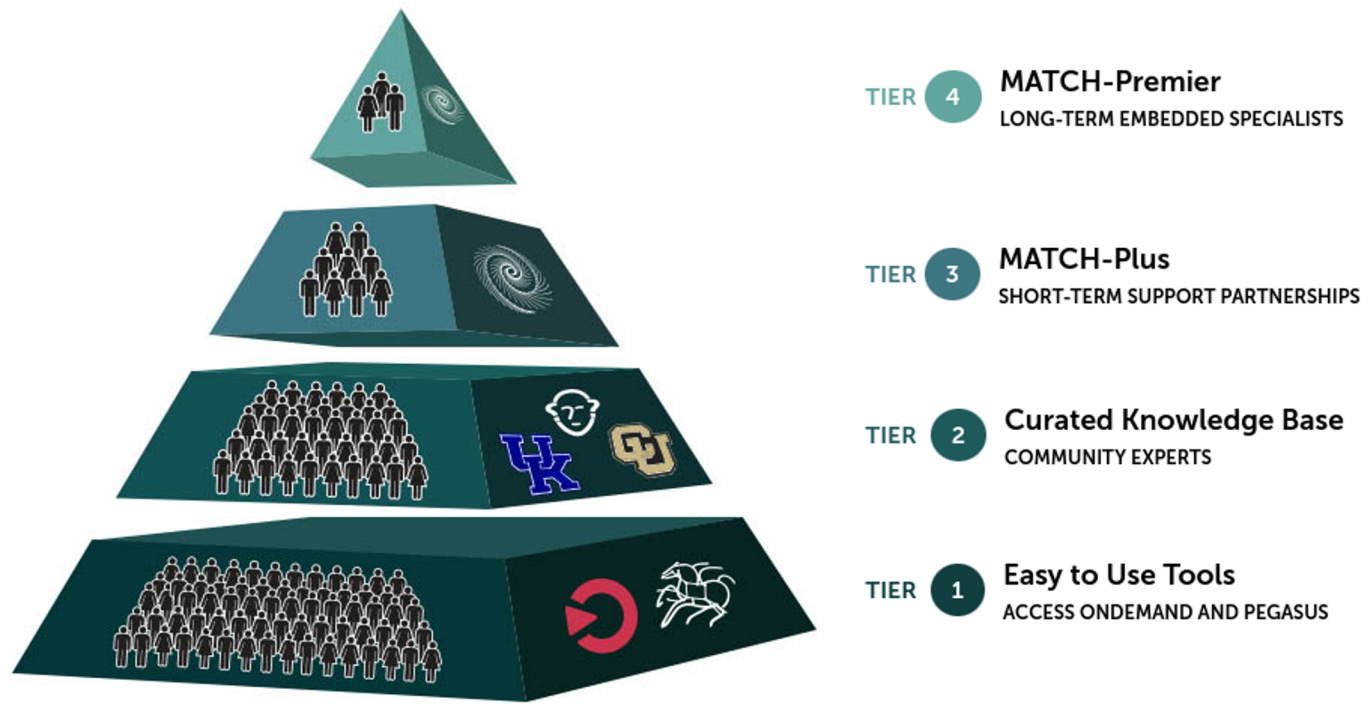
Pegasus is part of the ACCESS support strategy

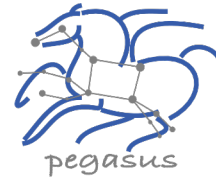
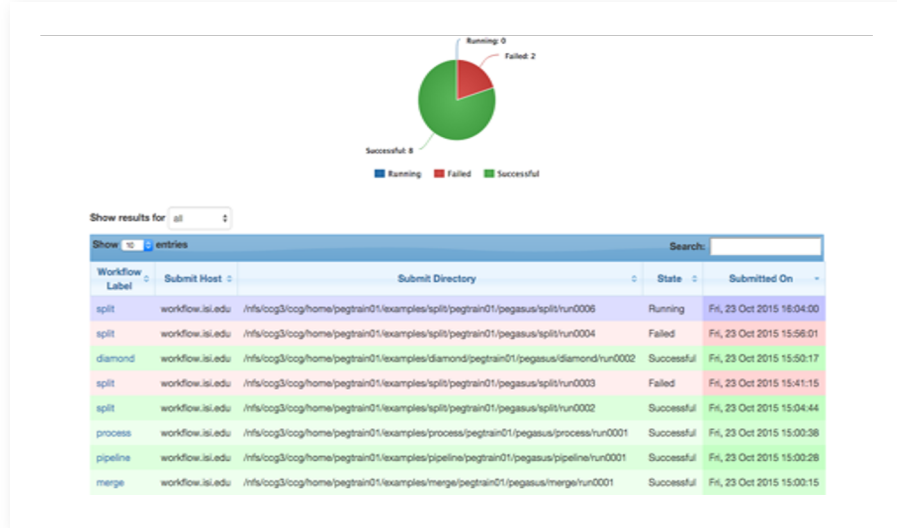
Pegasus is be used as a tier 1 tool

Central Open OnDemand instance with Pegasus, HTCondor and Jupyter

It is be easy to run HTC workflows across ACCESS sites

Tiered Support Strategy





PEGASUS DASHBOARD

web interface for monitoring and debugging workflows

Statistics

Workflow Wall Time	12 mins 23 secs
Workflow Cumulative Job Wall Time	9 mins 34 secs
Cumulative Job Walltime as seen from Submit Side	9 mins 35 secs
Workflow Cumulative Badput Time	9 mins 23 secs
Cumulative Job Badput Walltime as seen from Submit Side	9 mins 20 secs
Workflow Retries	1

Workflow Statistics

Type	Succeeded	Failed	Incomplete	Total	Retries	Total + Retries
Tasks	5	0	0	5	0	5
Jobs	16	0	0	16	2	18
Sub Workflows	0	0	0	0	0	0

Type	Succeeded	Failed	Incomplete	Total	Retries	Total + Retries
Tasks	5	0	0	5	0	5
Jobs	16	0	0	16	2	18
Sub Workflows	0	0	0	0	0	0

Real-time **monitoring** of workflow executions. It shows the **status** of the workflows and jobs, job **characteristics, statistics** and **performance** metrics.

Provenance data is stored into a relational database.

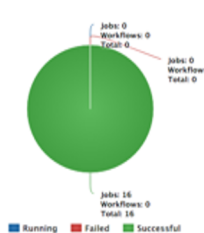
Workflow Details

Label	split
Type	root-wf
Progress	Successful
Submit Host	workflow.isi.edu
User	pegtrain01
Submit Directory	/ifs/ccg3/ccg/home/pegtrain01/examples/split/split/run0002
DAGMan Out File	split-0.dag.dagman.out
Wall Time	12 mins 23 secs
Cumulative Wall Time	9 mins 34 secs

Job Status (Entire Workflow)



Job Status (Per Workflow)



Real-time Monitoring

Reporting

Debugging

Troubleshooting

RESTful API



command-line...

```
$ pegasus-status pegasus/examples/split/run0001
STAT IN_STATE JOB
Run 00:39 split-0 (/home/pegasus/examples/split/run0001)
Idle 00:03 └─split_ID0000001
Summary: 2 Condor jobs total (I:1 R:1)

UNRDY READY PRE IN_Q POST DONE FAIL %DONE STATE DAGNAME
14      0      0      1      0      2      0      11.8 Running *split-0.dag
```

```
$ pegasus-analyzer pegasus/examples/split/run0001
pegasus-analyzer: initializing...

*****Summary*****

Total jobs : 7 (100.00%)
# jobs succeeded : 7 (100.00%)
# jobs failed : 0 (0.00%)
# jobs unsubmitted : 0 (0.00%)
```

```
$ pegasus-statistics -s all pegasus/examples/split/run0001
-----
Type           Succeeded Failed Incomplete Total Retries Total+Retries
Tasks           5         0         0         5         0         5
Jobs            17        0         0        17         0        17
Sub-Workflows   0         0         0         0         0         0
-----

Workflow wall time : 2 mins, 6 secs
Workflow cumulative job wall time : 38 secs
Cumulative job wall time as seen from submit side : 42 secs
Workflow cumulative job badput wall time :
Cumulative job badput wall time as seen from submit side :
```

**Provenance Data
can be Summarized
pegasus-statistics
or
Used for Debugging
pegasus-analyzer**



And if a job fails?



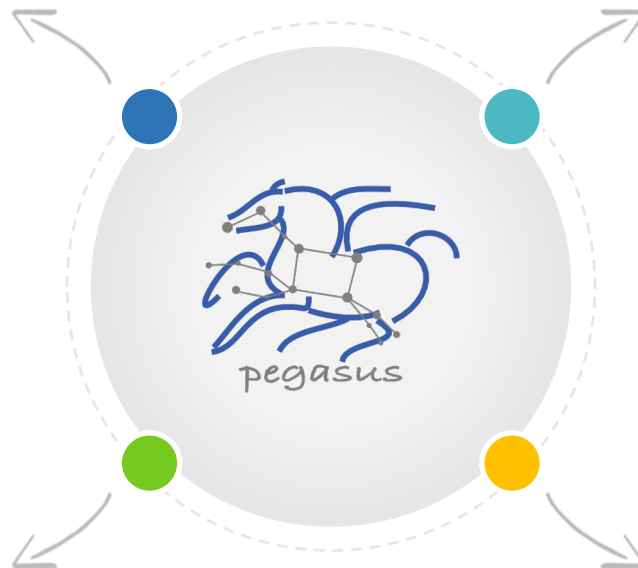
Postscript

detects non-zero exit code output parsing for success or failure message exceeded timeout do not produced expected output files



Checkpoint Files

job generates checkpoint files staging of checkpoint files is automatic on restarts



Job Retry



helps with transient failures set number of retries per job and run

Rescue DAGs



workflow can be restarted from checkpoint file recover from failures with minimal loss

Pegasus 5.0



- **New and fresh Python3 API to compose, submit and monitor workflows, and configure catalogs**
- **New Catalog Formats**
- **Python 3 Support**
 - ▶ All Pegasus tools are Python 3 compliant
 - ▶ Python PIP packages for workflow composition and monitoring
- **Zero configuration required to submit to local HTCondor pool.**
- **Data Management Improvements**
 - ▶ New output replica catalog that registers outputs including file metadata such as size and checksums
 - ▶ Improved support for hierarchical workflows
- **Reworked Documentation and Tutorial**
 - ▶ <https://pegasus.isi.edu/documentation/>

```
#!/usr/bin/env python3
import logging
import sys

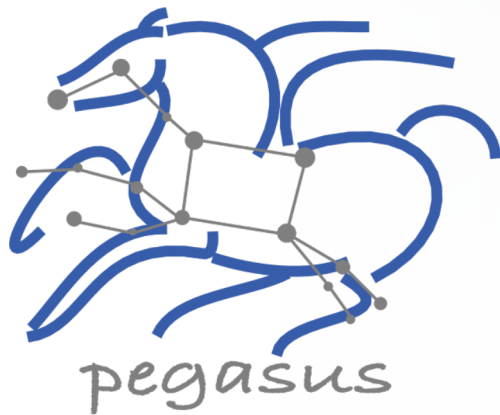
from Pegasus.api import *

# logs to be sent to stdout
logging.basicConfig(level=logging.DEBUG, stream=sys.stdout)

# --- Transformations -----
echo = Transformation(
    "echo",
    pfn="/bin/echo",
    site="condorpool"
)

tc = TransformationCatalog()\
    .add_transformations(echo)

# --- Workflow -----
Workflow("hello-world", infer_dependencies=True)\
    .add_jobs(
        Job(echo)
        .add_args("Hello World")
        .set_stdout("hello.out")
    ).add_transformation_catalog(tc)\
    .plan(submit=True)\
    .wait()
```



2. Hands on Exercises

Hands on Tutorial Exercises: Setup

Please claim an instance by putting your name next to an unused instance in:

<https://shorturl.at/CMV15>

Follow the link next to your name.

(This is the same (but hosted) as the self-guided tutorial available in the Pegasus documentation: <https://pegasus.isi.edu/documentation/user-guide/tutorial.html>)

Docker Container / Jupyter Notebook

Container is for tutorial purposes - most production workflows have dedicated submit hosts.

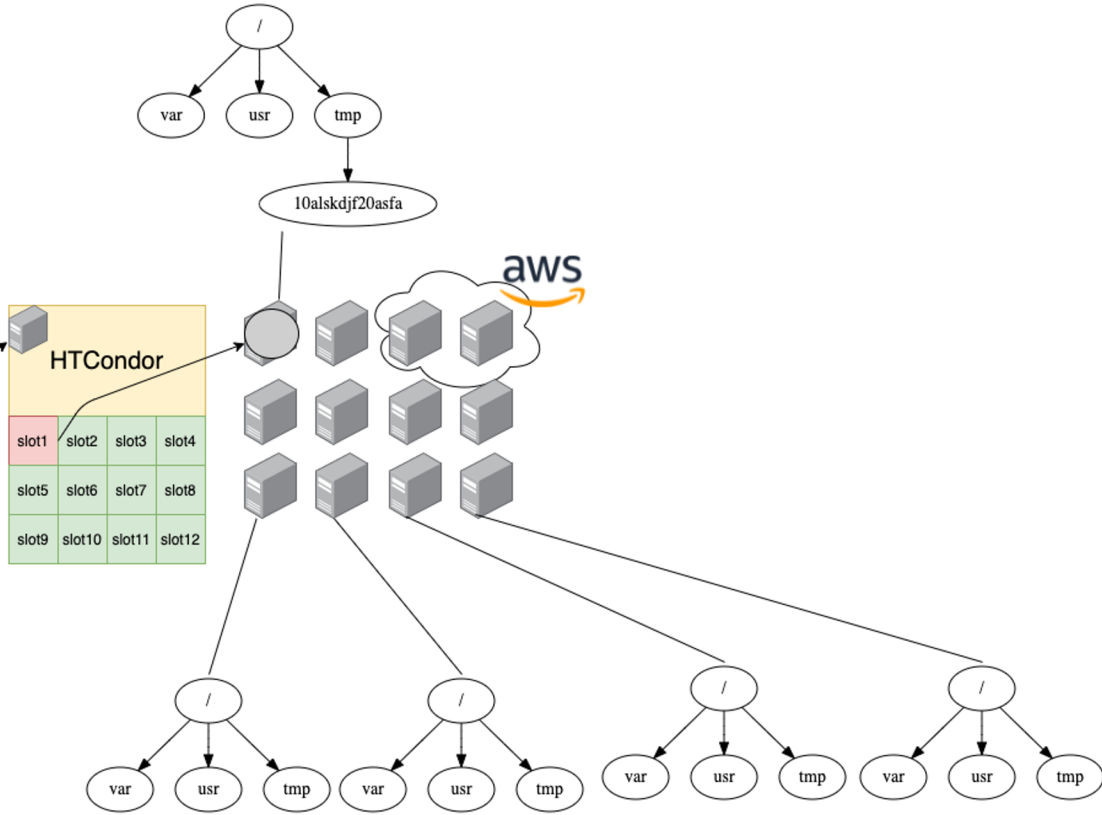
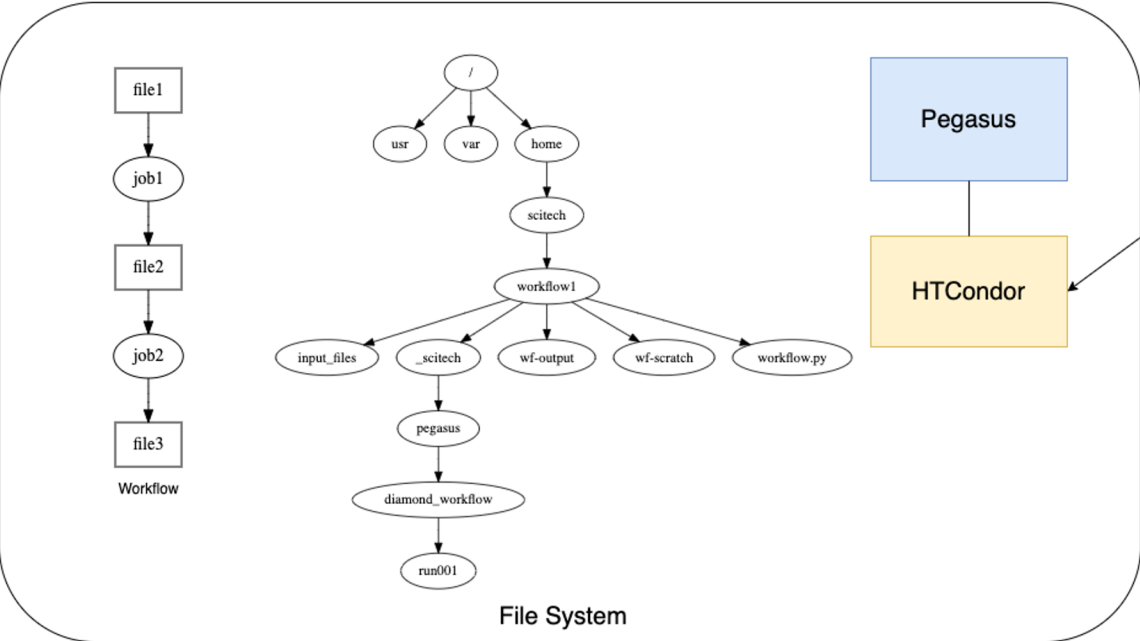
- One such host is available via **ACCESS** <http://access.pegasus.isi.edu>
- Submit workflows using Open OnDemand Portal
- Similar tutorial content available for that host

Jupyter is optional. You can choose to use just the workflow abstraction API, the full workflow management API, inside or outside Jupyter.

Docker Container / Jupyter Notebook

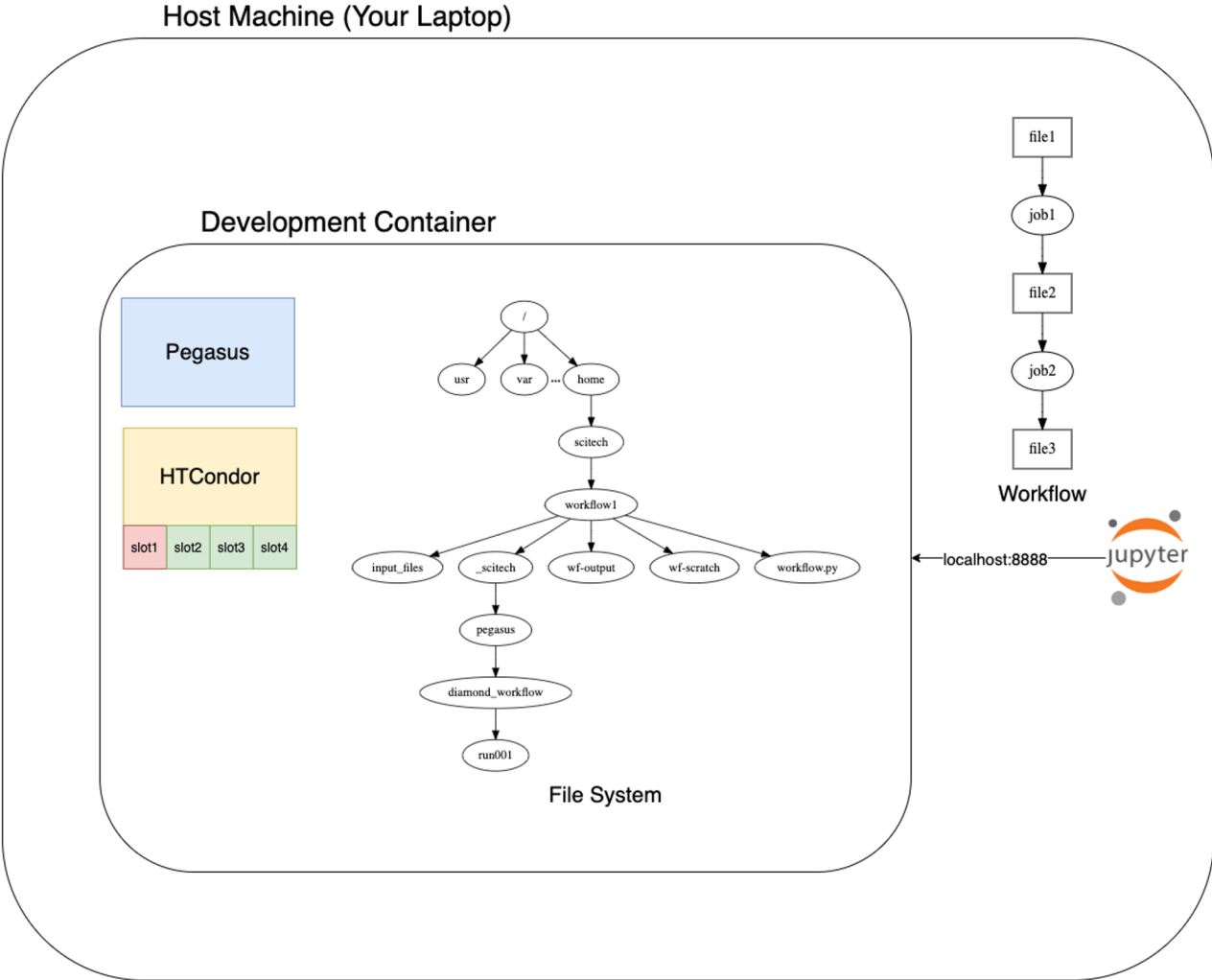
Production Setup

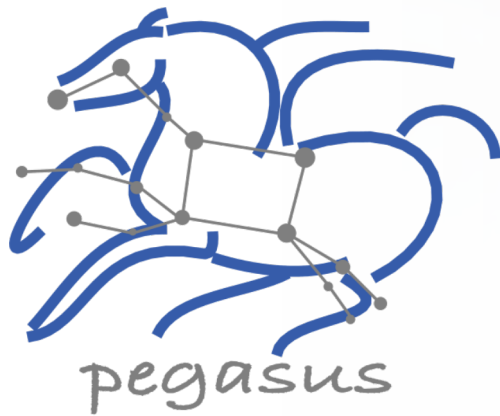
Host Machine (Your Laptop)



Docker Container / Jupyter Notebook

Development Setup





2.1 API



Key Pegasus Concepts

▲ Pegasus WMS == Pegasus planner (mapper) + DAGMan workflow engine + HTCondor scheduler/broker

- Pegasus maps workflows to infrastructure
- DAGMan manages dependencies and reliability
- HTCondor is used as a broker to interface with different schedulers

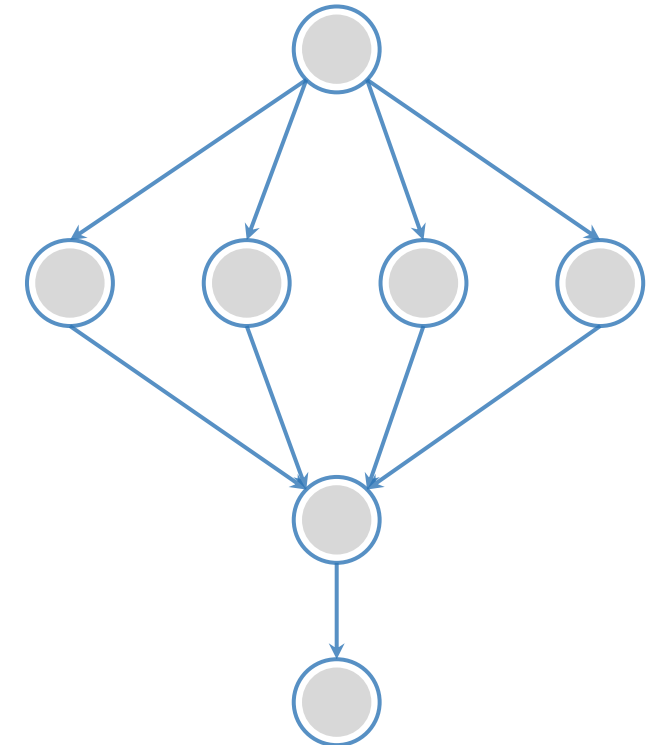
▲ Workflows are DAGs

- Nodes: jobs, edges: dependencies
- No while loops, no conditional branches
- Jobs are standalone executables

▲ Planning occurs ahead of execution

▲ Planning converts an abstract workflow into a concrete, executable workflow

- Planner is like a compiler



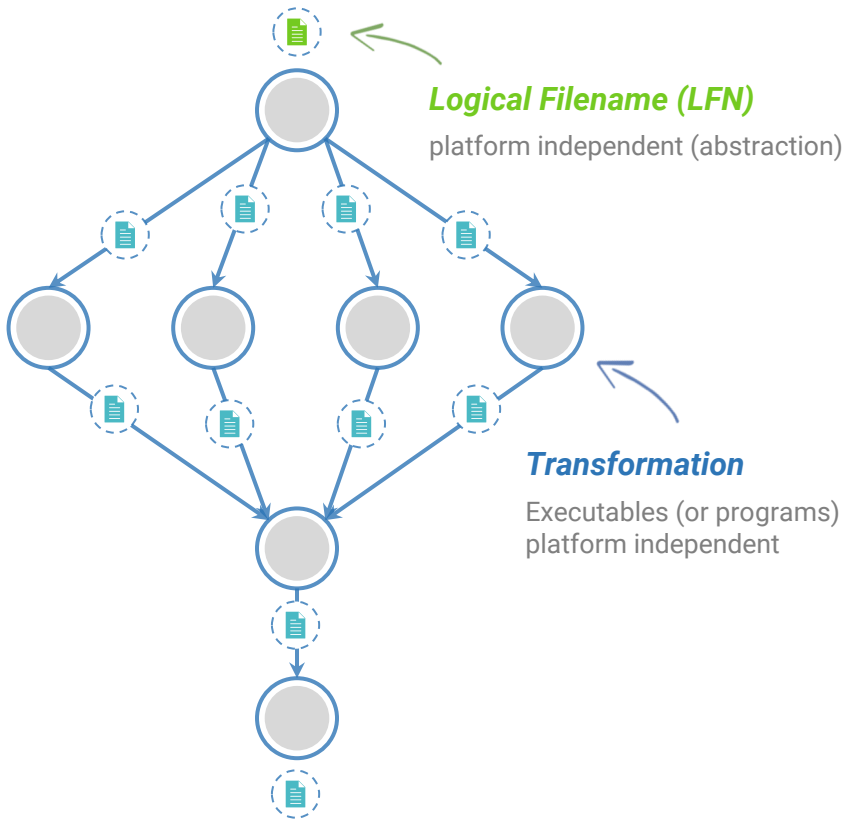


Input Workflow Specification **YAML formatted**

Portable Description

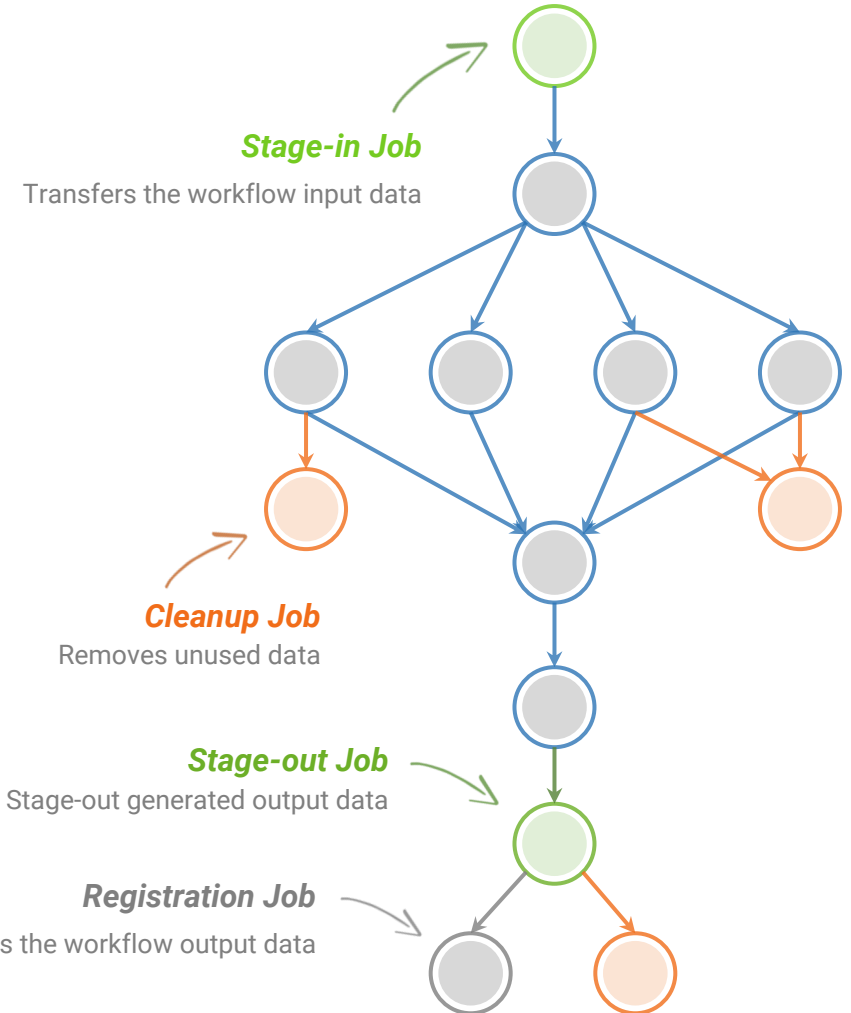
Users do not worry about low level execution details

ABSTRACT WORKFLOW

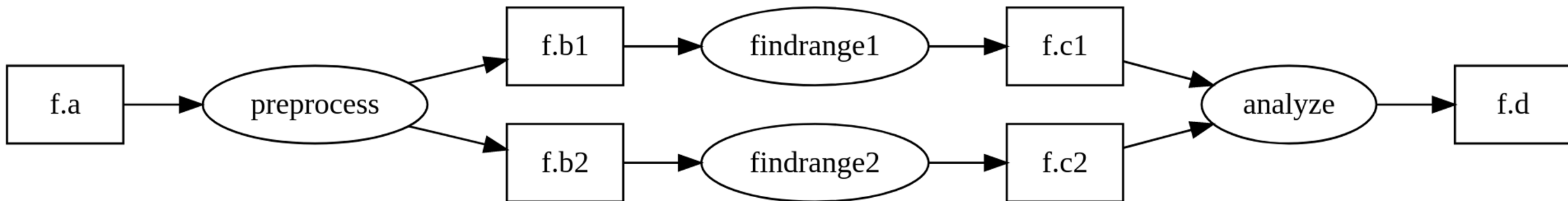


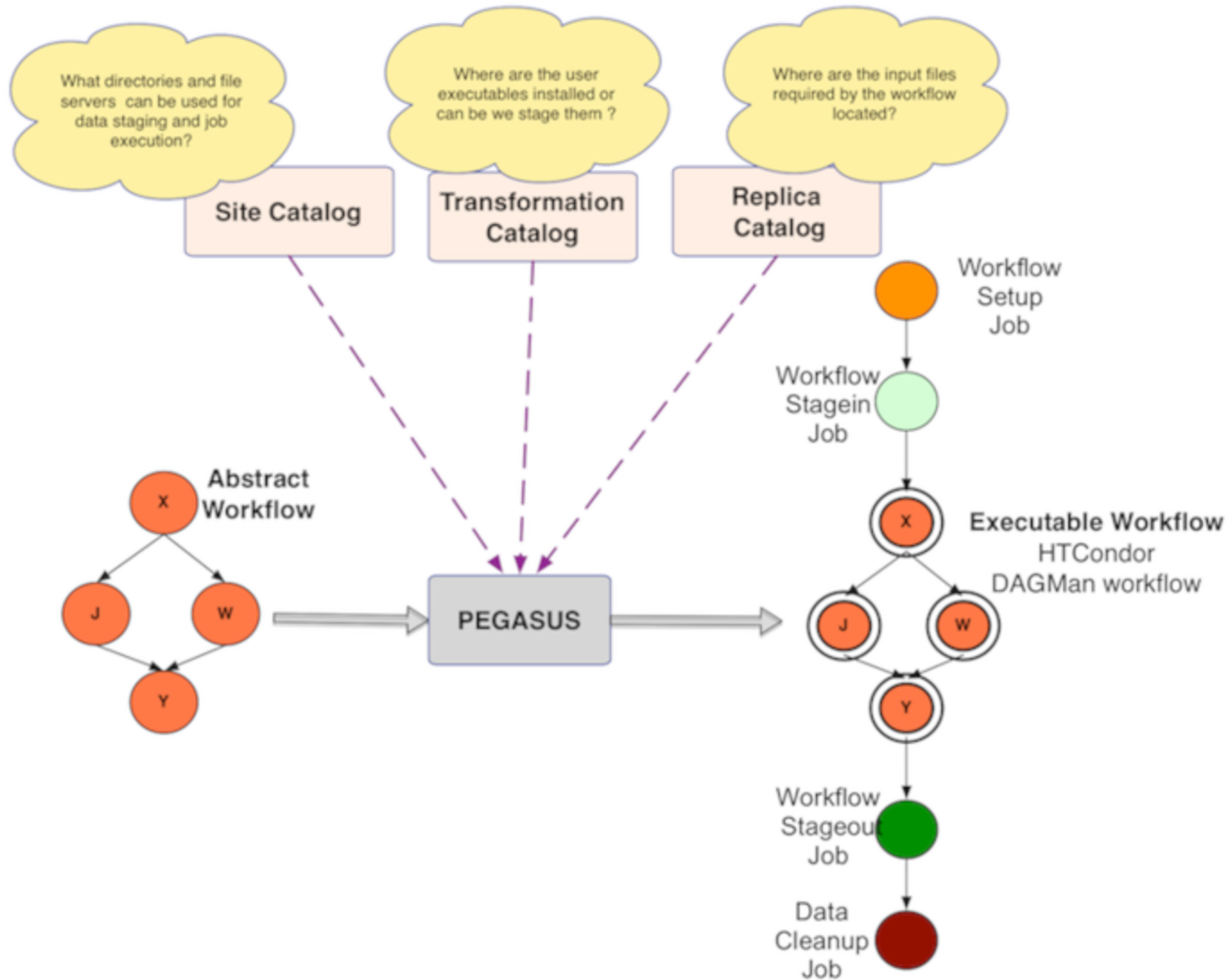
directed-acyclic graphs

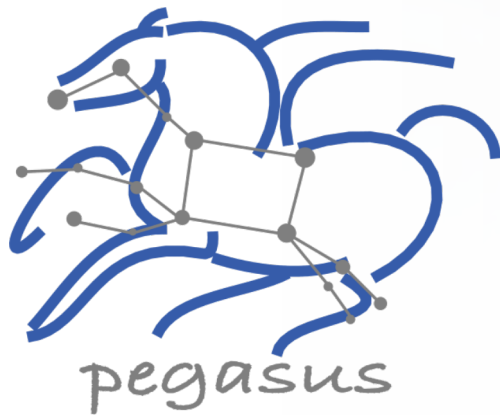
Output Workflow



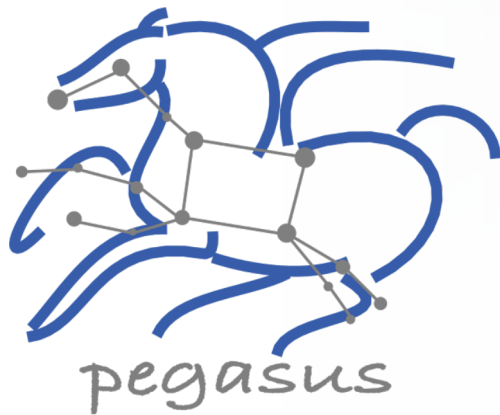
EXECUTABLE WORKFLOW



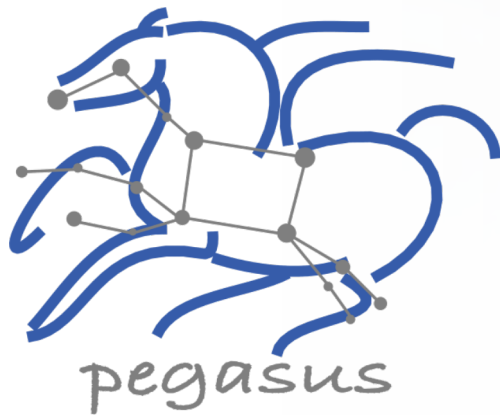




2.2 Debugging



2.3 Command Line Tools



30 Minute Break

Pegasus Container Support



Users can refer to **containers** in the **Transformation Catalog** with their executable preinstalled



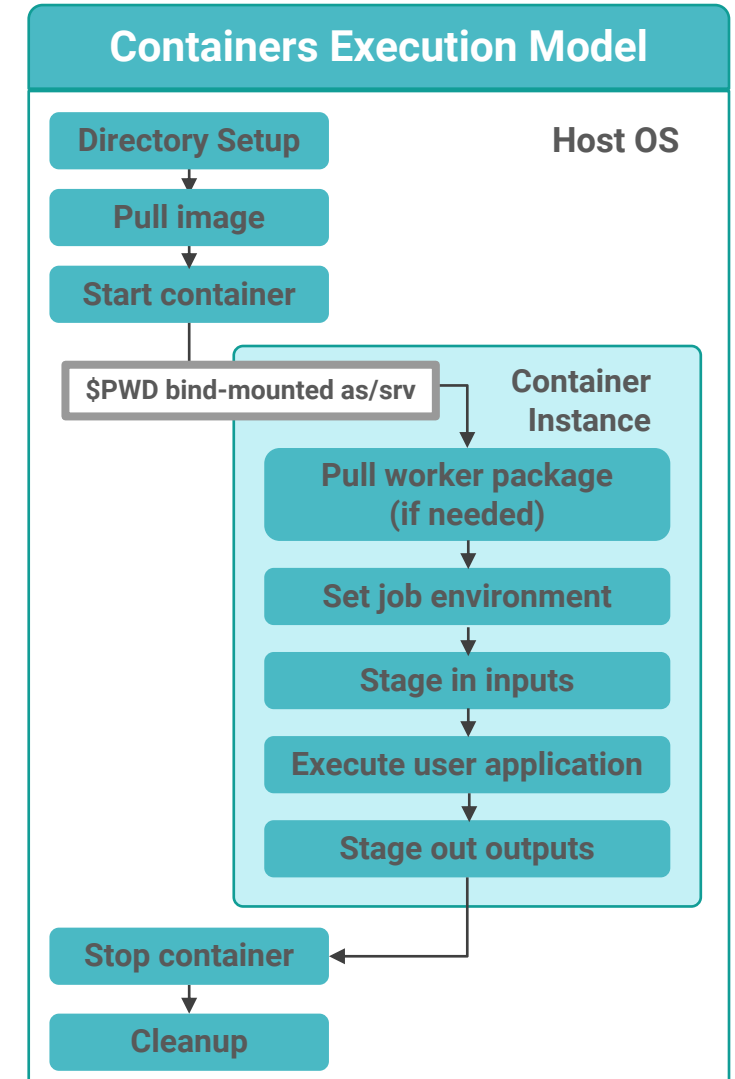
Users can **refer** to a **container** they want to **use** – **Pegasus stages** their executables and containers to the node

- Useful if you want to use a site recommended/standard container image.
- Users are using generic image with executable staging.



Future Plans

- Users can **specify an image buildfile** for their jobs.
- *Pegasus will build the Docker image as separate jobs in the executable workflow, export them as a tar file and ship them around*



Data Management for Containers



Containers are data too!

Pegasus treats containers as input data dependency

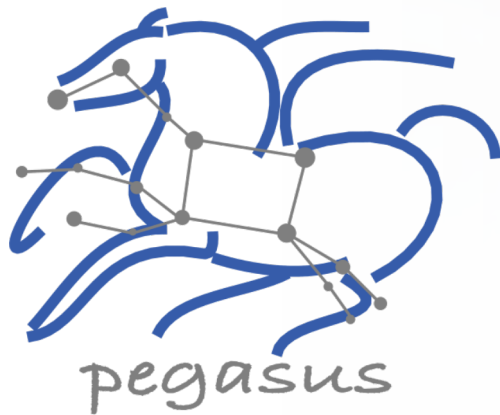
- Staged to compute node if not present
- Docker or Singularity Hub URL's
- Docker Image exported as a TAR file and available at a server, just like any other input dataset

Scaling up for larger workflows

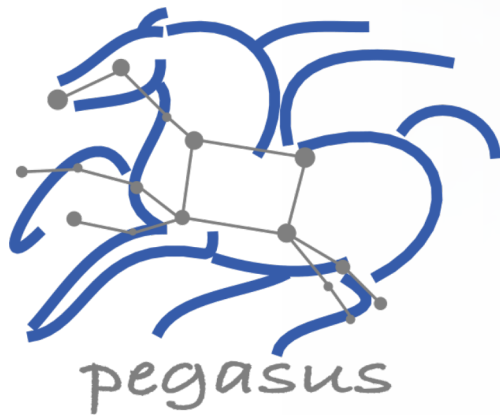
- The image is pulled down as a tar file as part of data stage-in jobs in the workflow
- The exported tar file is then shipped with the workflow and made available to the jobs
- Pricing considerations. You are now charged if you exceed a certain rate of pulls from Hubs

Other Optimizations

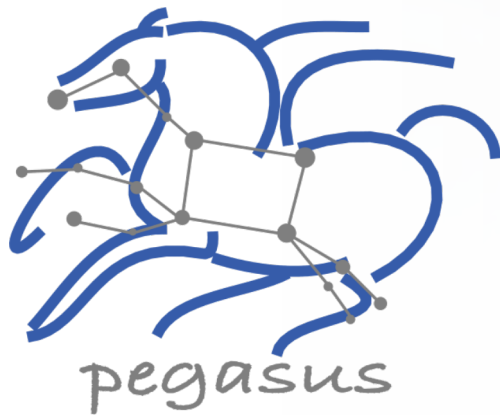
- **Symlink** against **existing images** on shared file system such as **CVMFS**
- The exported tar file is then shipped with the workflow and made available to the jobs



2.4 Containers



2.5 Summary



3. Advanced Topics

Data Staging Configurations

HTCondor I/O (HTCondor pools, OSG, ...)

- Worker nodes do not share a file system
- Data is pulled from / pushed to the submit host via HTCondor file transfers
- Staging site is the submit host

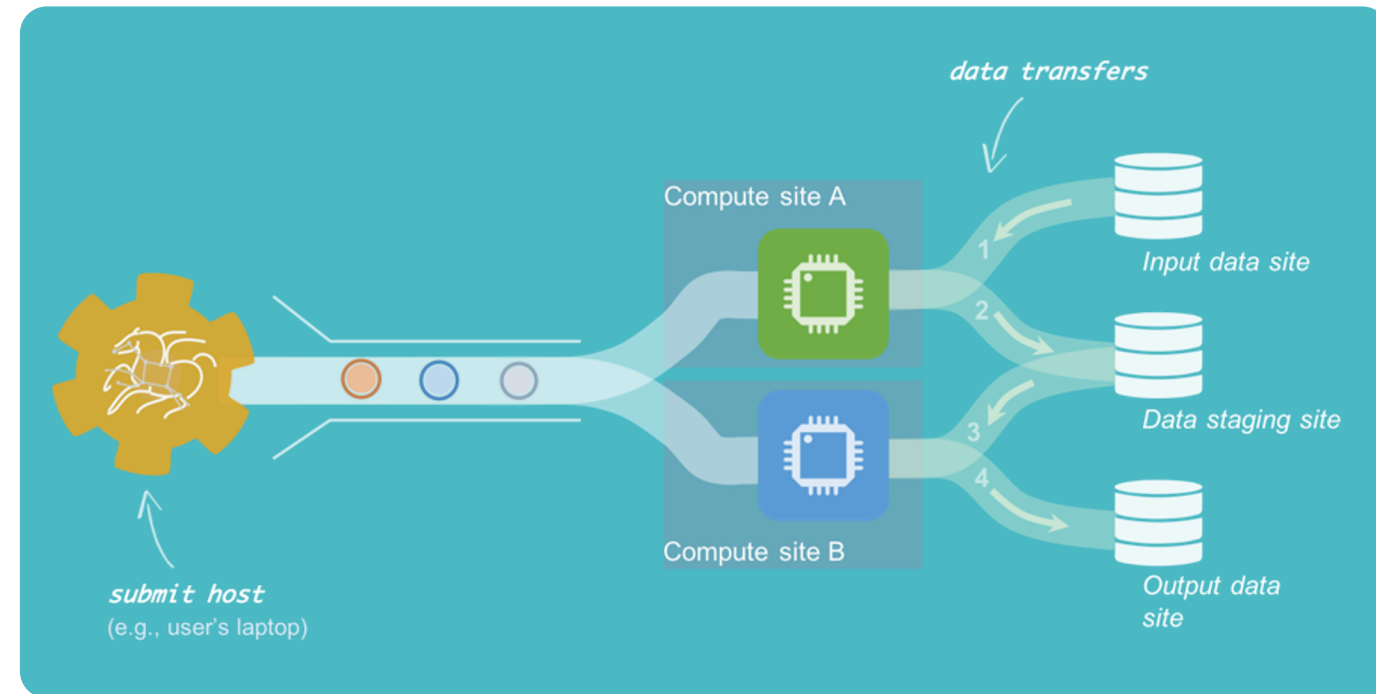
Non-shared File System (clouds, OSG, ...)

- Worker nodes do not share a file system
- Data is pulled / pushed from a staging site, possibly not co-located with the computation

Shared File System

(HPC sites, XSEDE, Campus clusters, ...)

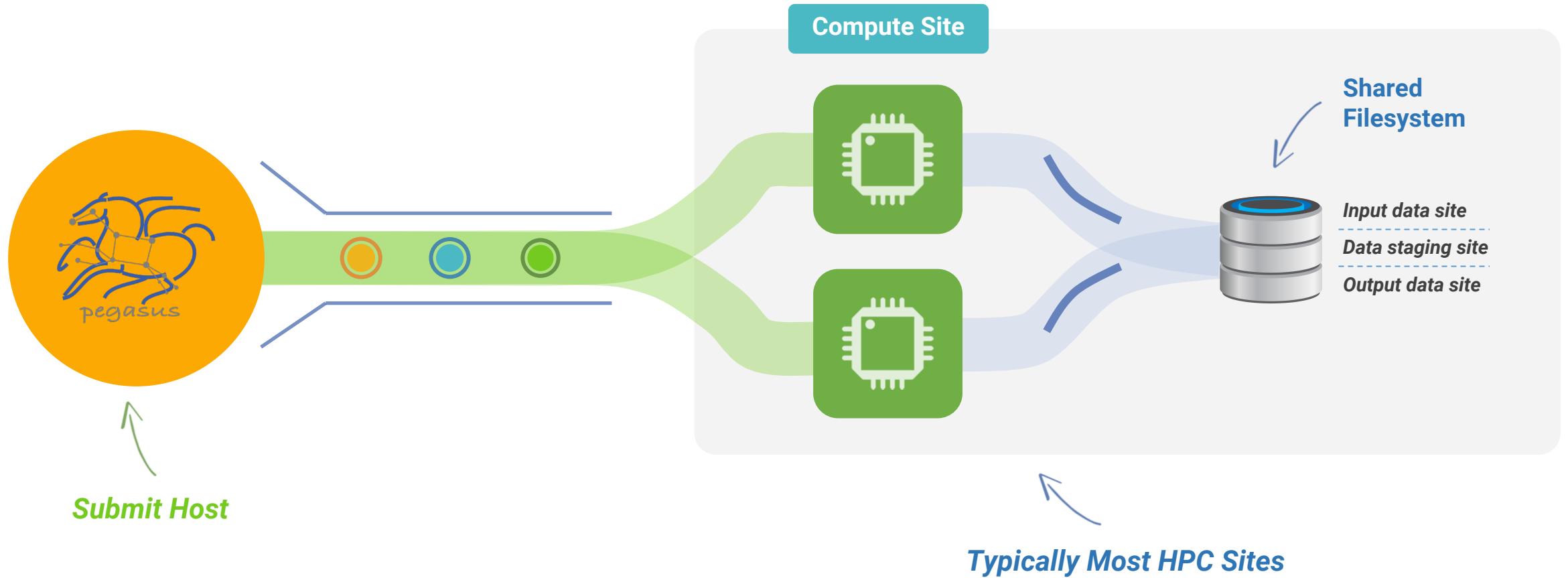
- I/O is directly against the shared file system





High Performance Computing

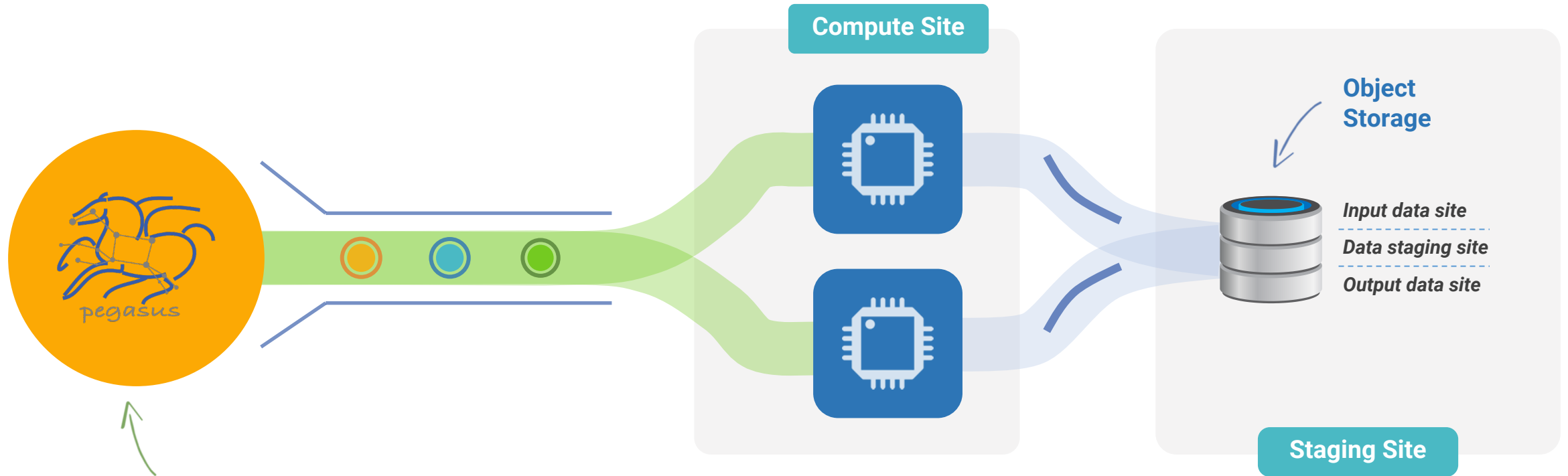
There are several possible configurations...





Cloud Computing

High-scalable object storages

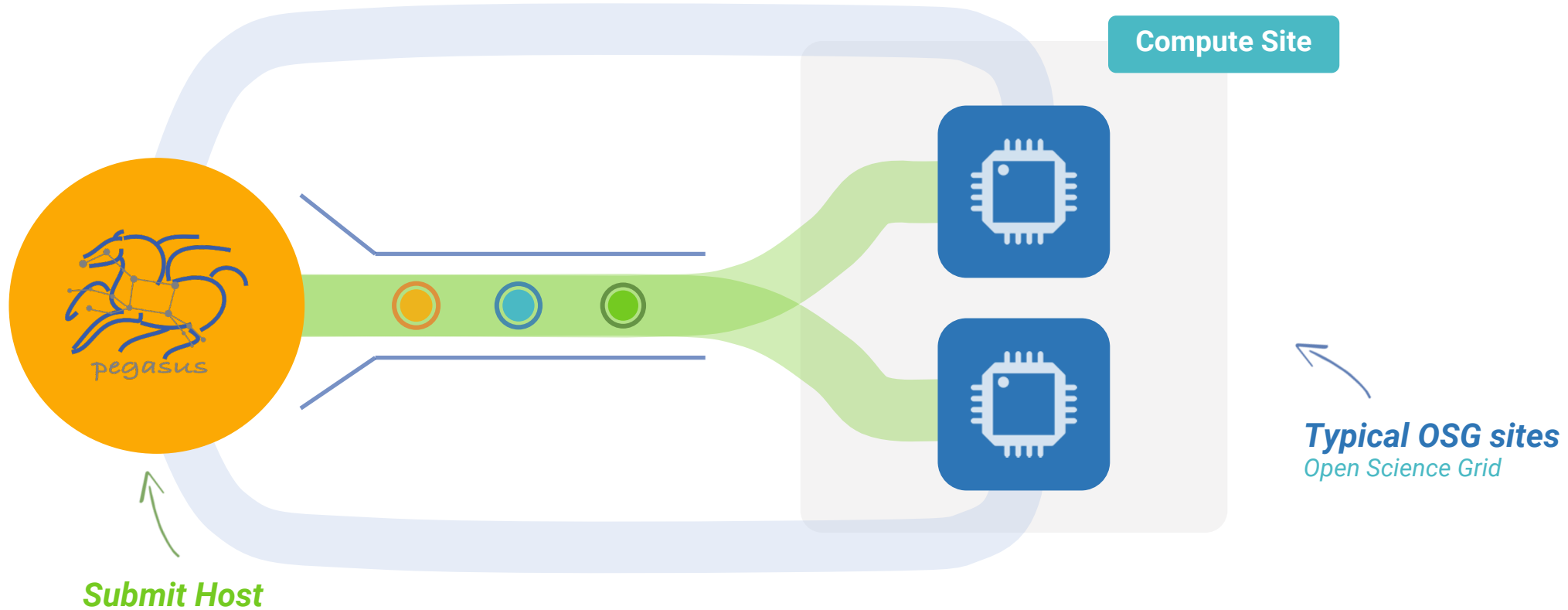


Typical cloud computing deployment
(Amazon S3, Google Storage)

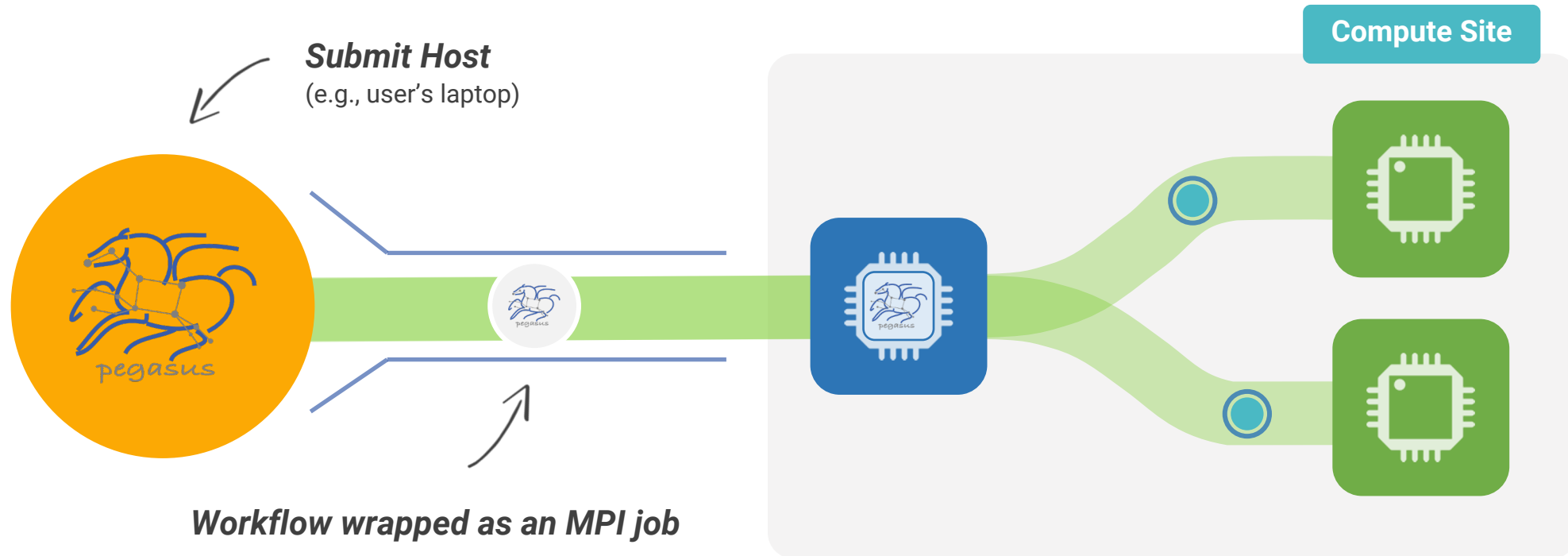


Grid Computing

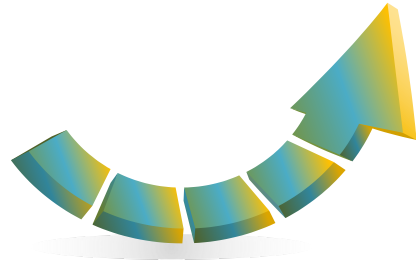
Local data management



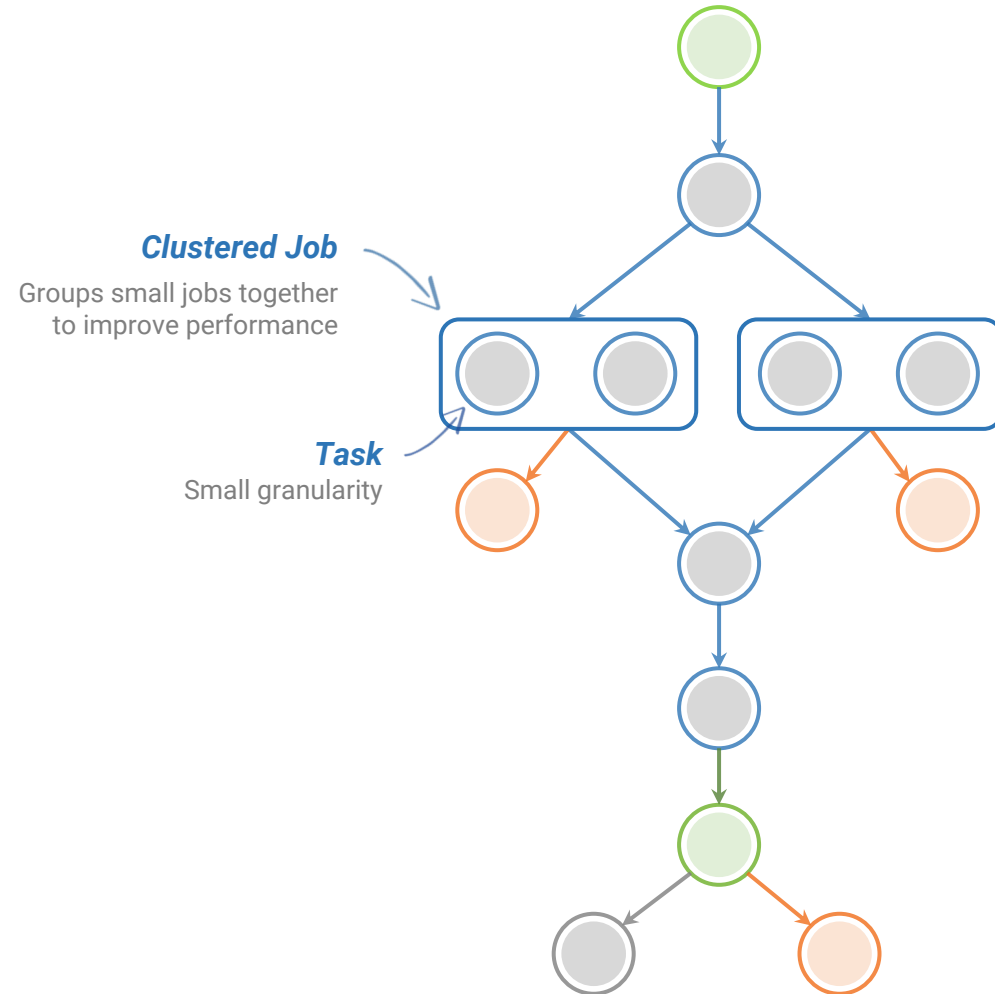
Running fine-grained workflows on HPC systems...



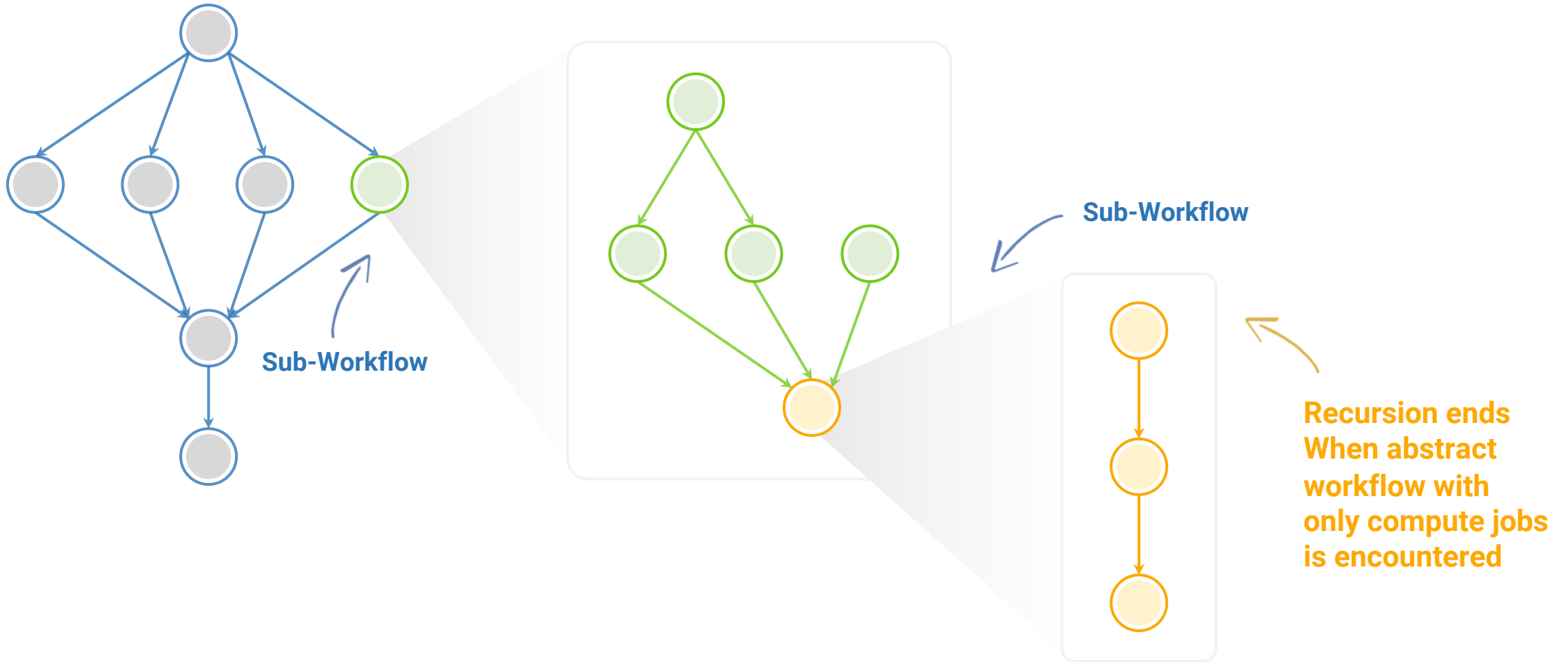
Allows sub-graphs of a Pegasus workflow to be submitted as monolithic jobs to remote resources



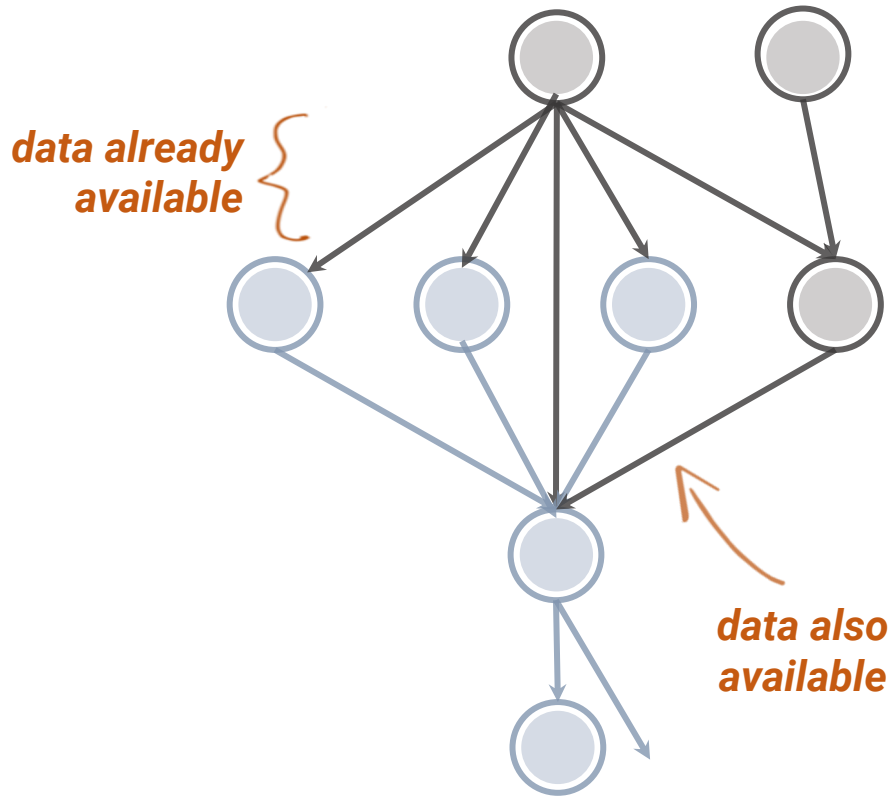
Performance. Why not improve it?



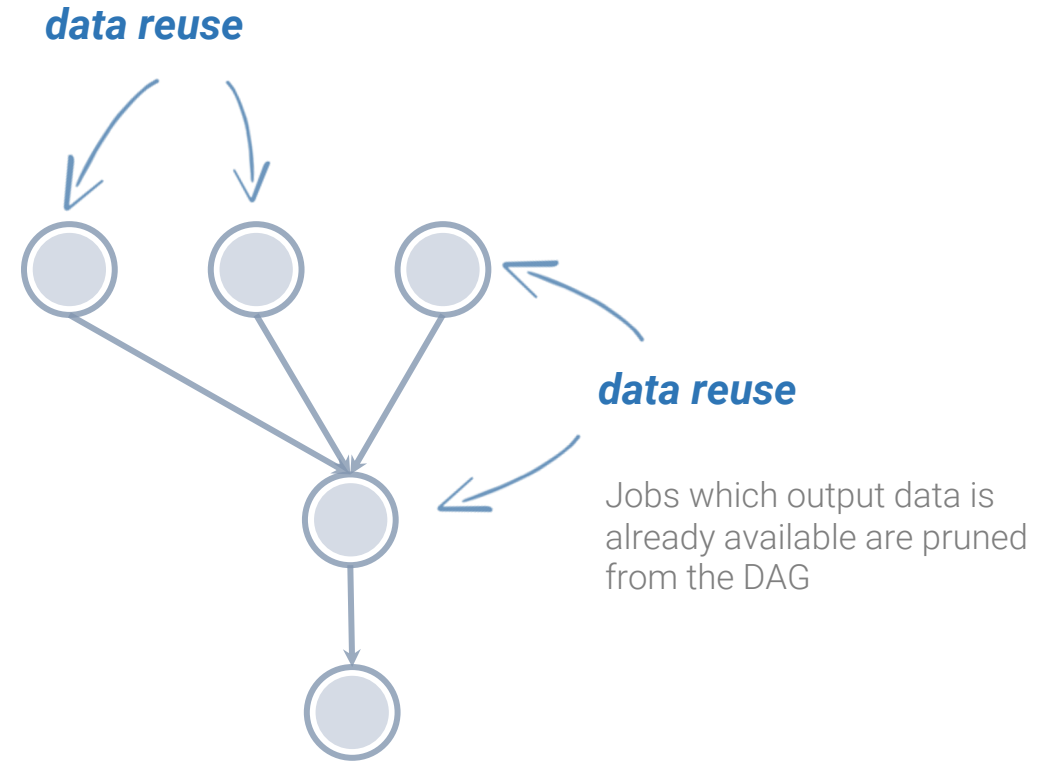
Pegasus also handles **large-scale workflows**



Data Reuse **prune jobs if output data already exists**



workflow
reduction





And if a job fails?



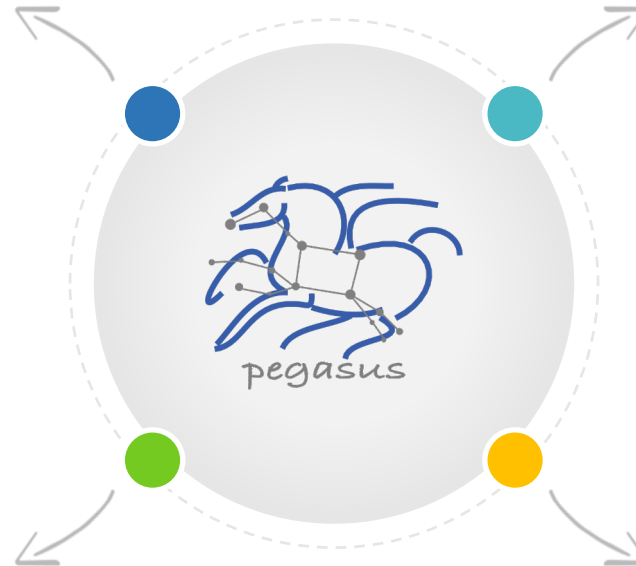
Postscript

detects non-zero exit code output parsing for success or failure message exceeded timeout do not produced expected output files



Checkpoint Files

job generates checkpoint files staging of checkpoint files is automatic on restarts



Job Retry



helps with transient failures set number of retries per job and run

Rescue DAGs



workflow can be restarted from checkpoint file recover from failures with minimal loss



Metadata

Can associate arbitrary key-value pairs with workflows, jobs, and files

Data Registration

Output files get tagged with metadata on registration in the workflow database

Workflow,
Job, File

```

x-pegasus:
apiLang: python
createdBy: vahi
createdOn: 12-08-20T10:08:48Z
pegasus: "5.0"
name: diamond
metadata:
  experiment: "par_all127_prot_lipid"
jobs:
- type: "job"
  name: "namd"
  id: "ID0000001"
  arguments: ["equilibrate.conf"]
  uses:
  - lfn: "Q42.psf"
    metadata:
      type: "psf"
      charge: "42"
    type: "input"
  - lfn: "eq.restart.coord"
    type: "output"
    metadata:
      type: "coordinates"
      stageOut: true
      registerReplica: true
  metadata:
    timesteps: 500000
    temperature: 200
    pressure: 1.01353

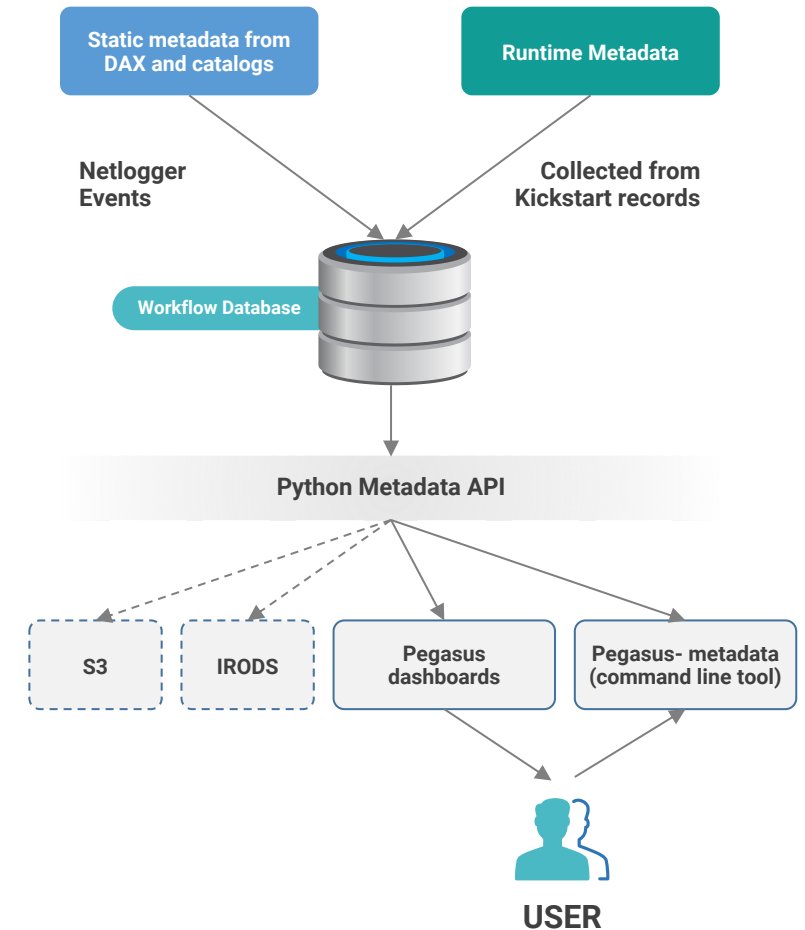
```

Static and Runtime Metadata

Static: application parameters
Runtime: performance metrics

Select Data
Based on Metadata

Register Data
With Metadata





Challenges to Scientific Data Integrity

Modern IT systems
are not perfect
- errors creep in.

At modern “**Big Data**” sizes we
are starting to see checksums
breaking down.

Plus there is the threat
of intentional changes:
*malicious attackers,
insider threats, etc.*

User Perception: “Am I not already protected? I have heard about TCP checksums, encrypted transfers, checksum validation, RAID and erasure coding – is that not enough?”

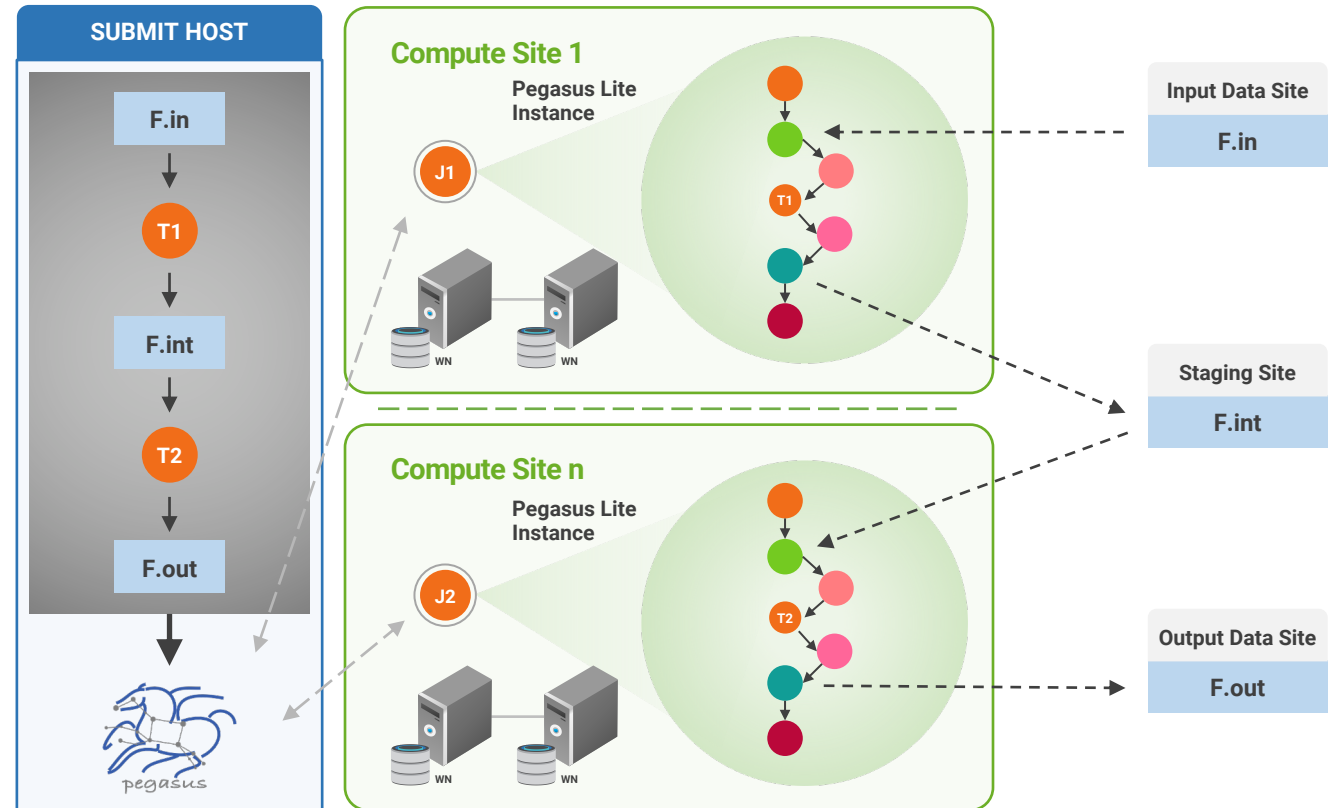


Automatic Integrity Checking in Pegasus

Pegasus performs integrity checksums on input files right before a job starts on the remote node.

- ▶ For raw inputs, **checksums specified in the input replica catalog** along with file locations
- ▶ All **intermediate** and **output** files checksums are generated and tracked within the system.
- ▶ Support for **sha256** checksums

Job failure is triggered if checksums fail



LEGEND									
	Task flow + Checksums		Directory Setup Job		Data Stageout Job		Check Integrity Job		Pegasus Lite Compute Job
	Data Flow		Data Stagein Job		Directory Cleanup Job		Checksum Generation Job		Worker Node (WN)



Job Submissions

LOCAL

Submit Machine

Personal HTCondor

Local Campus Cluster accessible via Submit Machine **

HTCondor via BLAHP

**** Both Glite and BOSCO build on HTCondor BLAHP**

**Currently supported schedulers:
SLURM SGE PBS MOAB**

REMOTE

BOSCO + SSH**

Each node in executable workflow submitted via SSH connection to remote cluster

BOSCO based Glideins**

SSH based submission of glideins

PyGlidein

IceCube glidein service

OSG using glideinWMS

Infrastructure provisioned glideins

CREAMCE

Uses CondorG

Globus GRAM

Uses CondorG

Credentials Management

▲ Credentials required for two purposes



- Job Submission
- Data transfers to **stage-in** input and **stage-out** generated outputs when a job executes

▲ Specifying Credentials

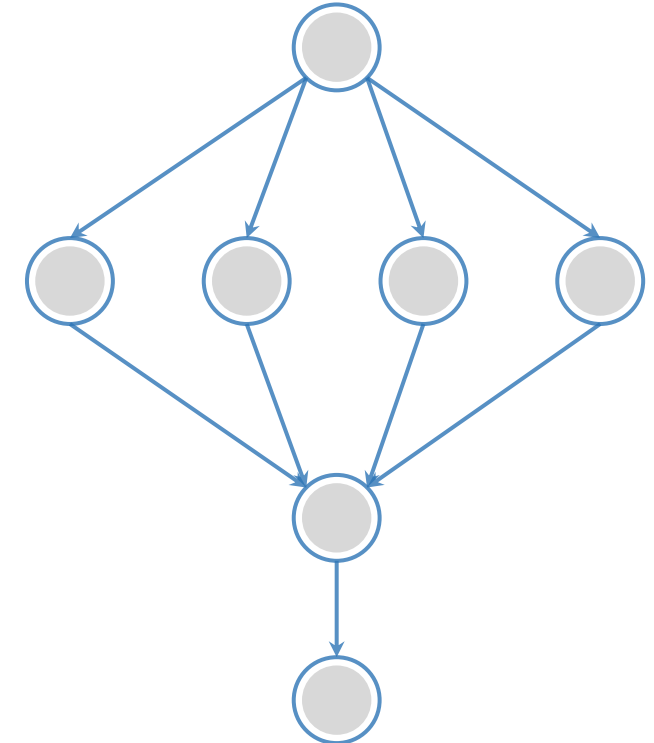
- Users can specify credentials in a **generic credentials file** on submit host
- Associate credentials with sites in site catalog

▲ Approach

- Planner will **automatically** associate the **required credentials** with each job
- The credentials are **transferred** along with the job
- Usually available **only for the duration** of the job **execution**

▲ Supported Credentials

- X.509 grid proxies
- Amazon AWS S3 keys,
- Google Cloud Platform OAuth token (.boto file),
- iRods password
- SSH keys
- Web Dav



Amazon AWS Batch

AWS Batch

- ▶ Container based, dynamically scaled and efficient batch computing service
- ▶ Automatically launches compute nodes in Amazon based on demand in the associated job queue
- ▶ Users can specify compute environment that dictates what type of VM's are launched

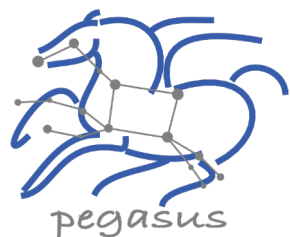
Pegasus will **allow clusters of jobs** to be run on **Amazon EC2** using **AWS Batch Service**

New command
line tool:

pegasus-aws-batch

Automates most of the batch setup programmatically

- **Sets up and Deprovisions**
 - Compute Environment
 - Job Queues
- **Follows AWS Batch HTTP specification**



Pegasus

est. 2001

Automate, recover, and debug scientific computations.

▶ Get Started

▶ Pegasus Website

<https://pegasus.isi.edu>

▶ Users Mailing List

pegasus-users@isi.edu

▶ Support

pegasus-support@isi.edu

▶ Slack

Ask for an invite by trying to join pegasus-users.slack.com in the Slack app

▶ Pegasus Online Office Hours

<https://pegasus.isi.edu/blog/online-pegasus-office-hours/>



YouTube Channel

<https://www.youtube.com/channel/UCwJQln1CqBvTJqiNr9X9F1Q/featured>



[Pegasus in 5 Minutes](#)