

# Pegasus 5.0 Workflows

Karan Vahi

Information Sciences Institute

University of Southern California, School of Engineering

[vahi@isi.edu](mailto:vahi@isi.edu)

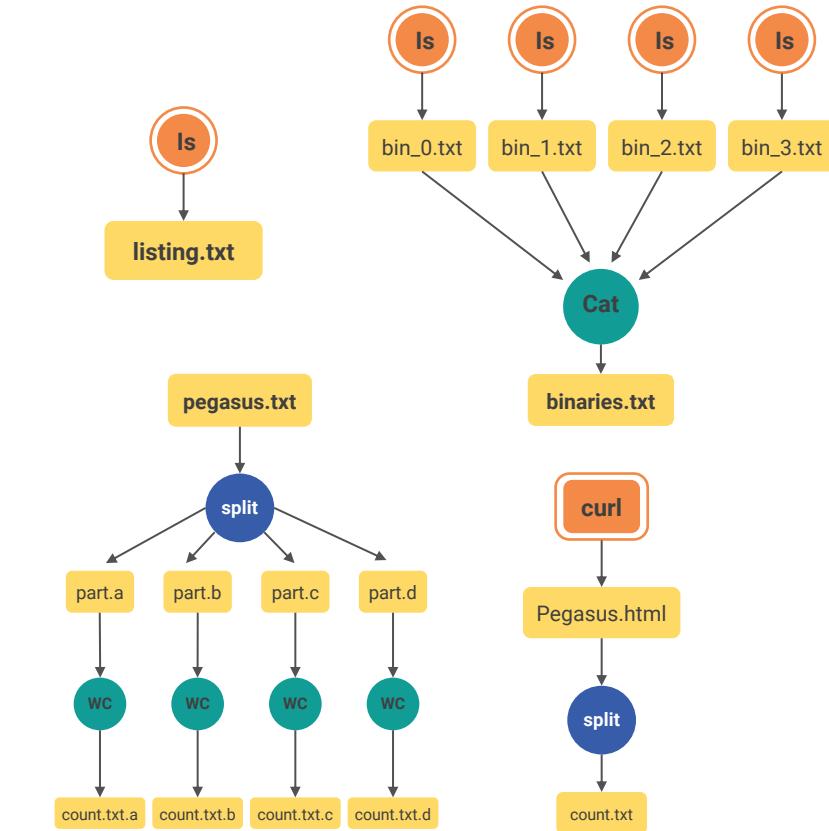


# What are Scientific Workflows



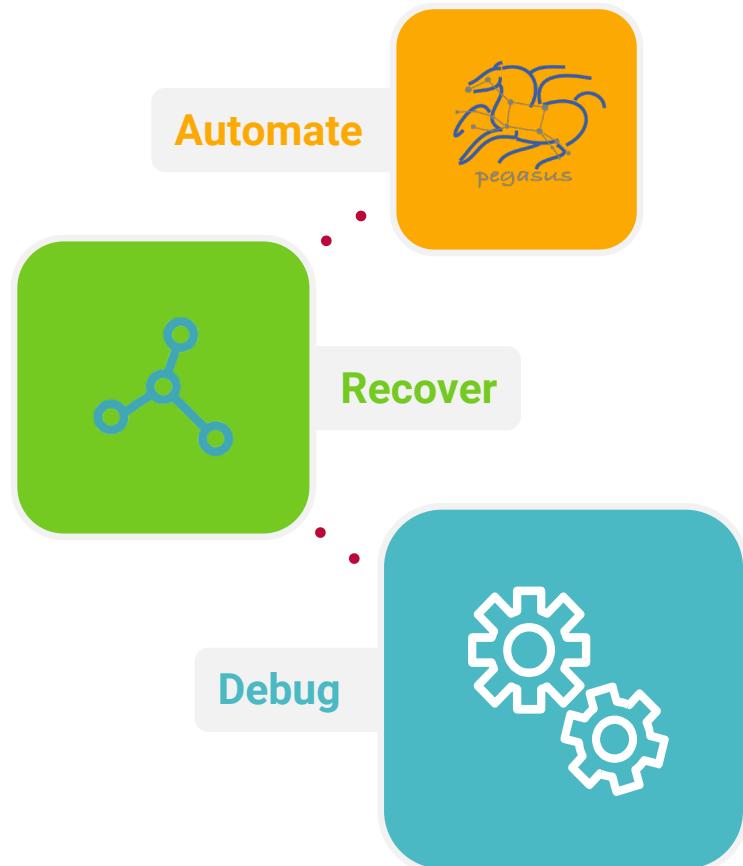
- ▶ **Conducts a series of computational tasks.**
  - Resources distributed across Internet.
- ▶ **Chaining (outputs become inputs) replaces manual hand-offs.**
  - Accelerated creation of products.
- ▶ **Ease of use - gives non-developers access to sophisticated codes.**
  - Resources distributed across 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, SCEC, USC

# Why Pegasus?



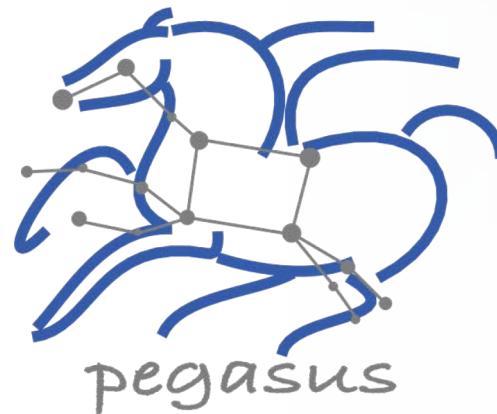
- ▶ **Automates Complex**, Multi-stage Processing Pipelines
- ▶ Enables Parallel, **Distributed Computations**
- ▶ **Automatically Executes** Data Transfers
- ▶ Reusable, Aids **Reproducibility**
- ▶ Records How Data was Produced (**Provenance**)
- ▶ Handles **Failures** with to Provide Reliability
- ▶ 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

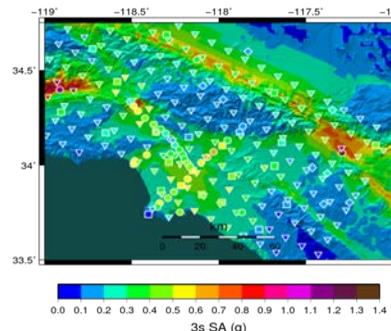
## Our Focus

- ▶ 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.



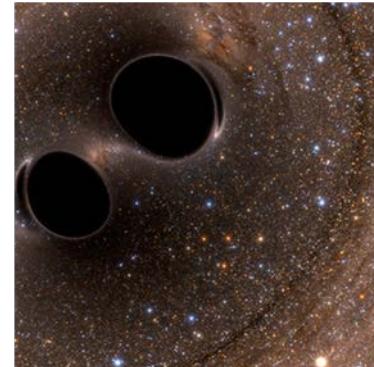
## Some of The Success Stories...

# Southern California Earthquake Center's CyberShake



*First Physics-Based "Shake map" of Southern California*

# Laser Interferometer Gravitational-Wave Observatory (LIGO)



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

# XENONnT - Dark Matter Search



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

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

Supported since 2001, distributed data, opportunistic computing resources

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



**CPU jobs**  
(Mesh generation, seismogram synthesis)



1,094,000 node-hours

**GPU jobs:**



439,000 node-hours

AWP-ODC finite-difference code

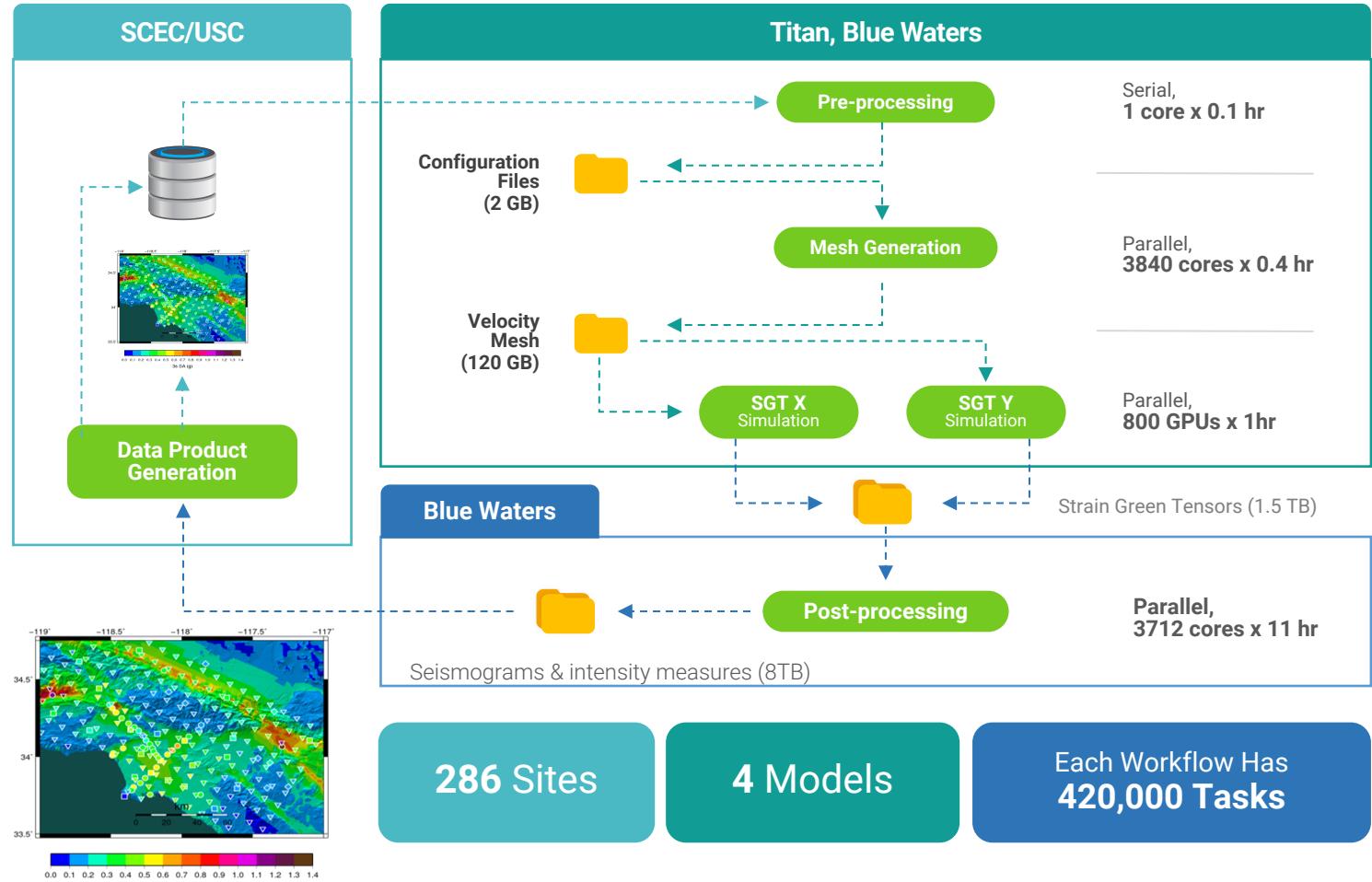
5 billion points per volume, 23,000 timesteps

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





# Data Flow for LIGO Pegasus Workflows in OSG

## Advanced LIGO

Laser Interferometer  
Gravitational Wave Observatory



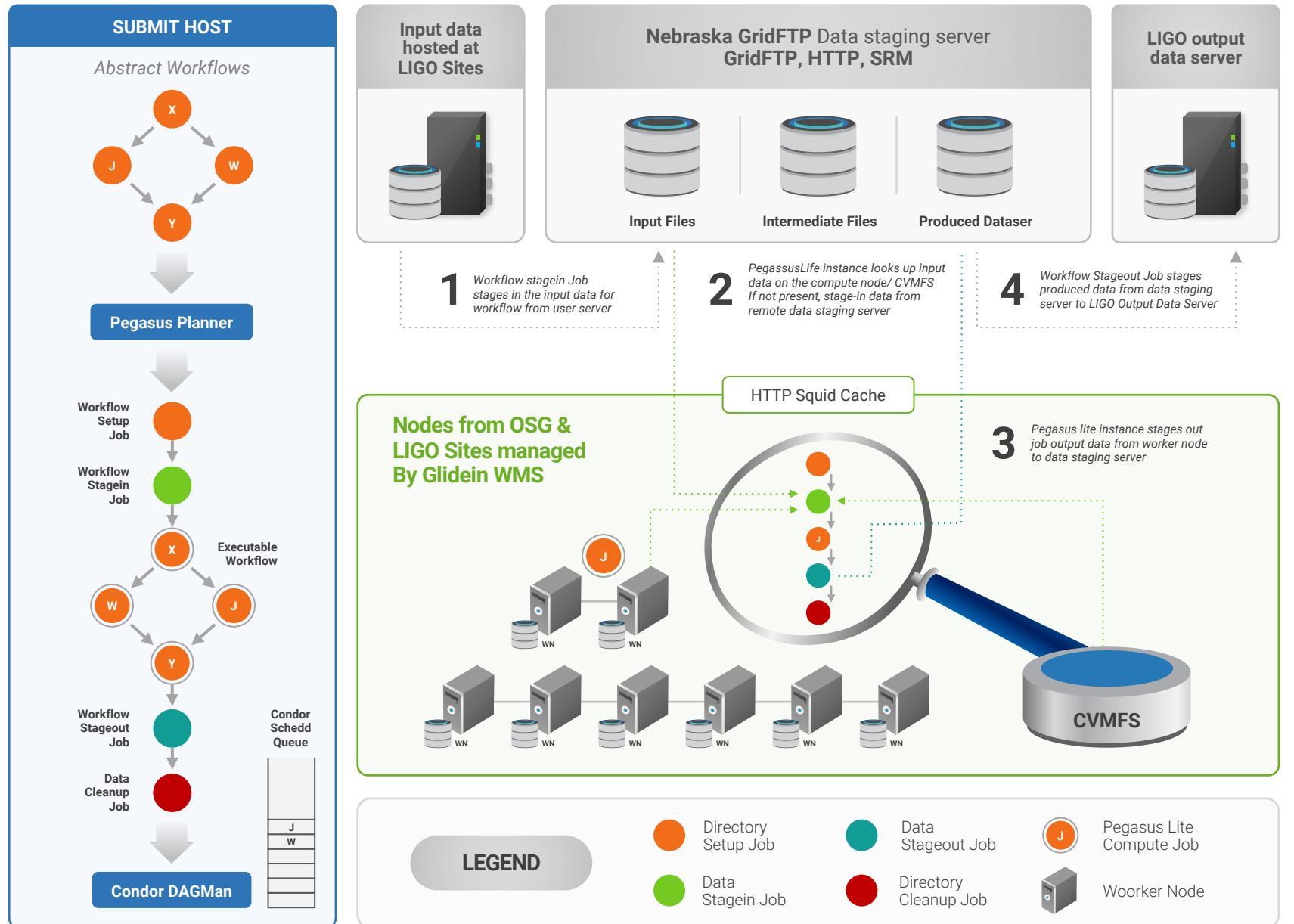
60,000 Compute Tasks

Input Data: 5000 files (10GB total)

Output Data: 60,000 files (60GB total)

Processed Data: 725 GB

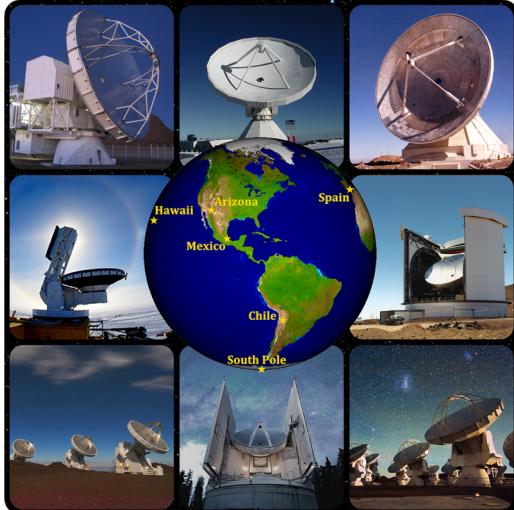
Executed on LIGO Data Grid, EGI,  
Open Science Grid and XSEDE



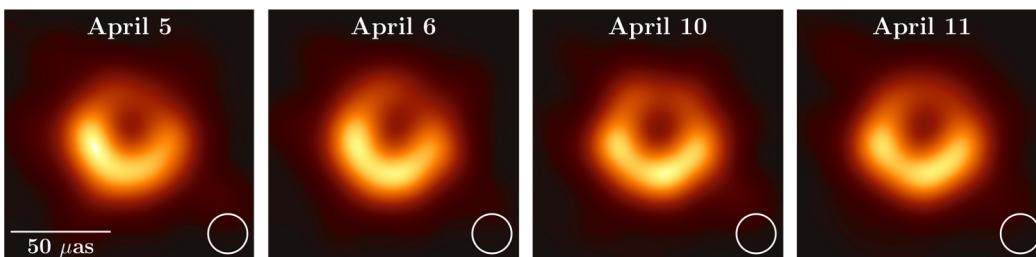
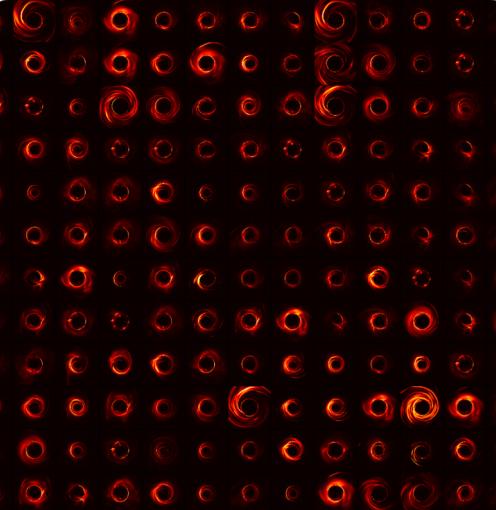
# Event Horizon Telescope

## Bringing Black Holes into Focus

8 telescopes: 5 PB of data



60 simulations: 35 TB data



First images of black hole at the center of the M87 galaxy

**Improve constraints on Einstein's theory of general relativity by 500x**

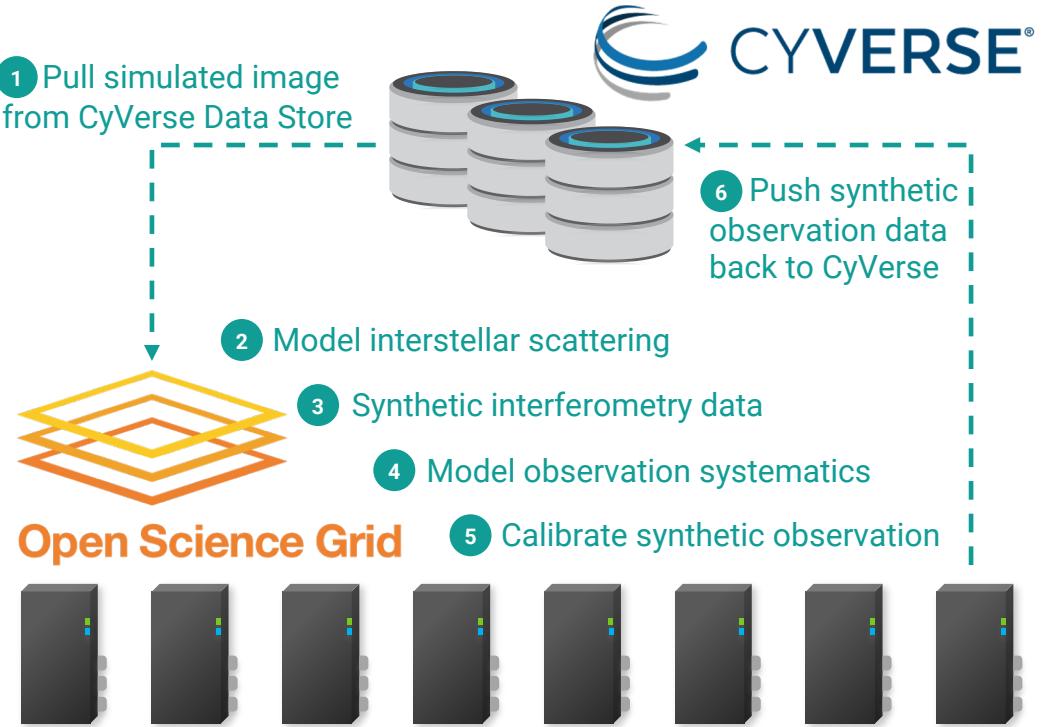
480,000 jobs - 2,600,000 core hours

#15 in all OSG projects in last 6 months

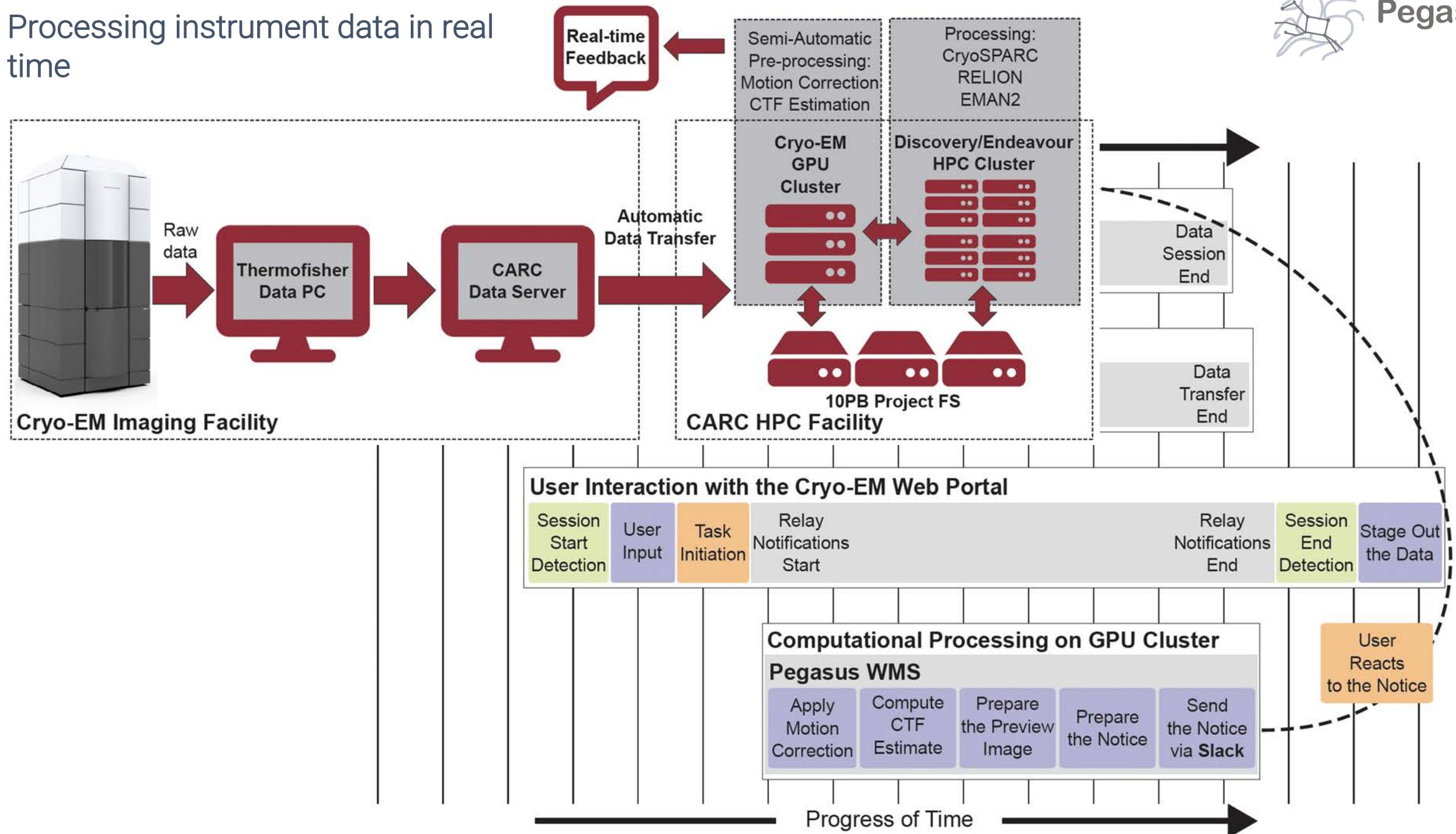
#2 in all OSG astronomy projects in the last 6 months

### Pegasus-SYMBIA Pipeline

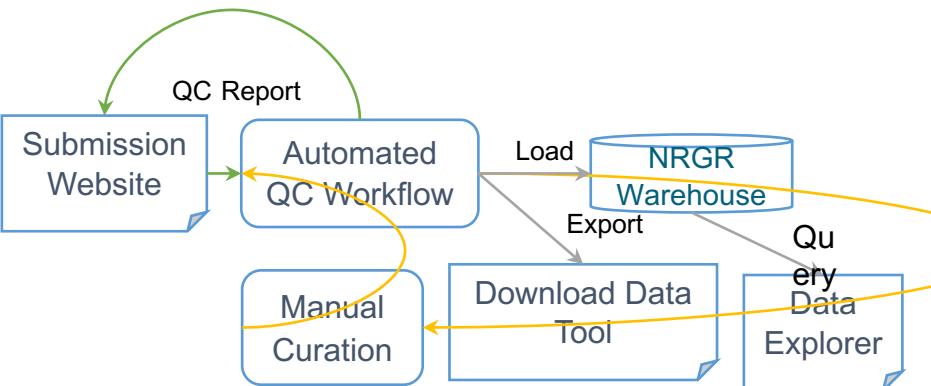
Physically accurate synthetic observation data from simulations are keys to develop calibration and imaging algorithms, as well as comparing the observation with theory and interpreting the results.



# 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.



## 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**

Flowchart of validation steps:

```

graph TD
    id_validation[id_validation] --> extended_diagnosis_validation[extended_diagnosis_validation]
    race_ethnicity_validation[race_ethnicity_validation] --> extended_diagnosis_validation
    extended_diagnosis_validation --> advanced_qc[advanced_qc]
    phenotypic_validation[phenotypic_validation] --> advanced_qc
    submission_validation[submission_validation] --> advanced_qc
    pedigree_validation[pedigree_validation] --> advanced_qc
  
```

- 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

<https://pegasus.isi.edu>

## Auto QC Status

◀ Back to Previous Validations

Successful: 100%

### Summary

UID	5e6a6ddd95f6e
Disorder	Depression
Study Id	149
File	shaptest7.zip
User	JaclynVitanza
Email	jv607@dis.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	Log
study_149_sub.csv	✓ Standardized File	✓ Log	✓ Log
study_149_id.csv	✓ Standardized File	✓ Log	✓ Log
shaps01.phen.csv	✓ Standardized File	✓ Log	✓ Log
Advanced QC			
study_149_sub.canon.csv	✓ Corrected Submission File		
study_149_id.canon.csv	✓ Corrected ID File		
Corrections Log	✓ Corrections Log		
Advanced QC Report	✓ Advanced QC Report		

Basic Concepts...



# 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

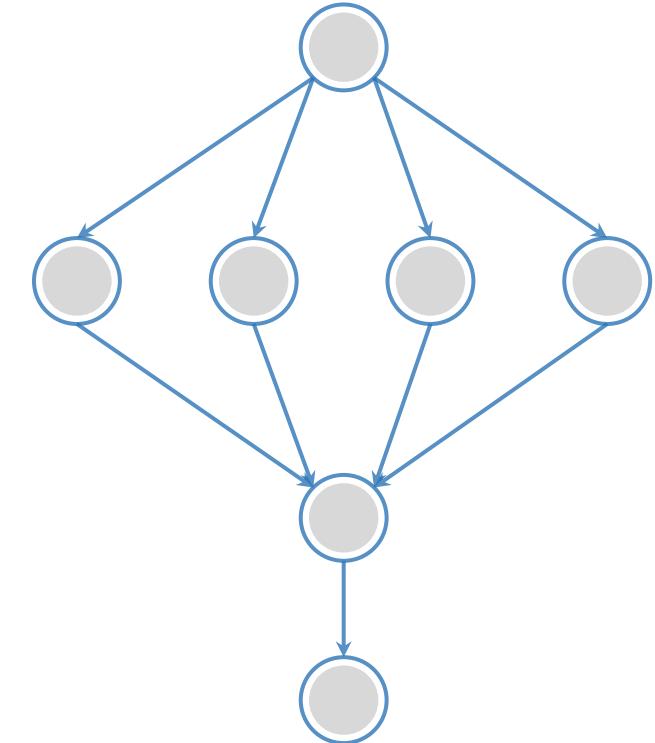
## ► Workflows are DAGs

- 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

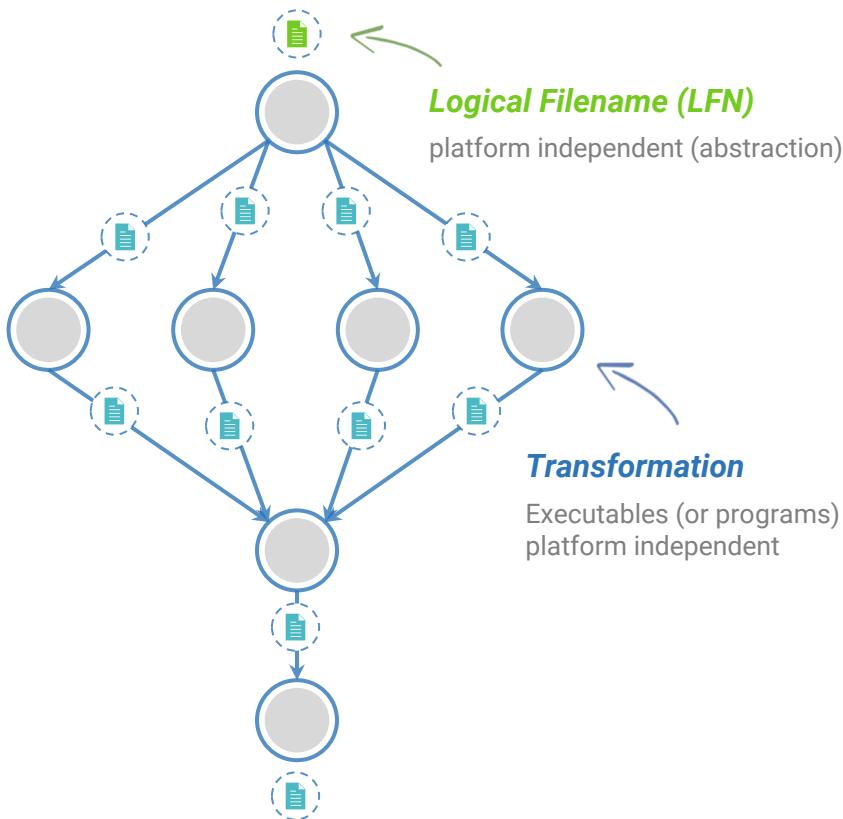


## Input Workflow Specification YAML formatted

### Portable Description

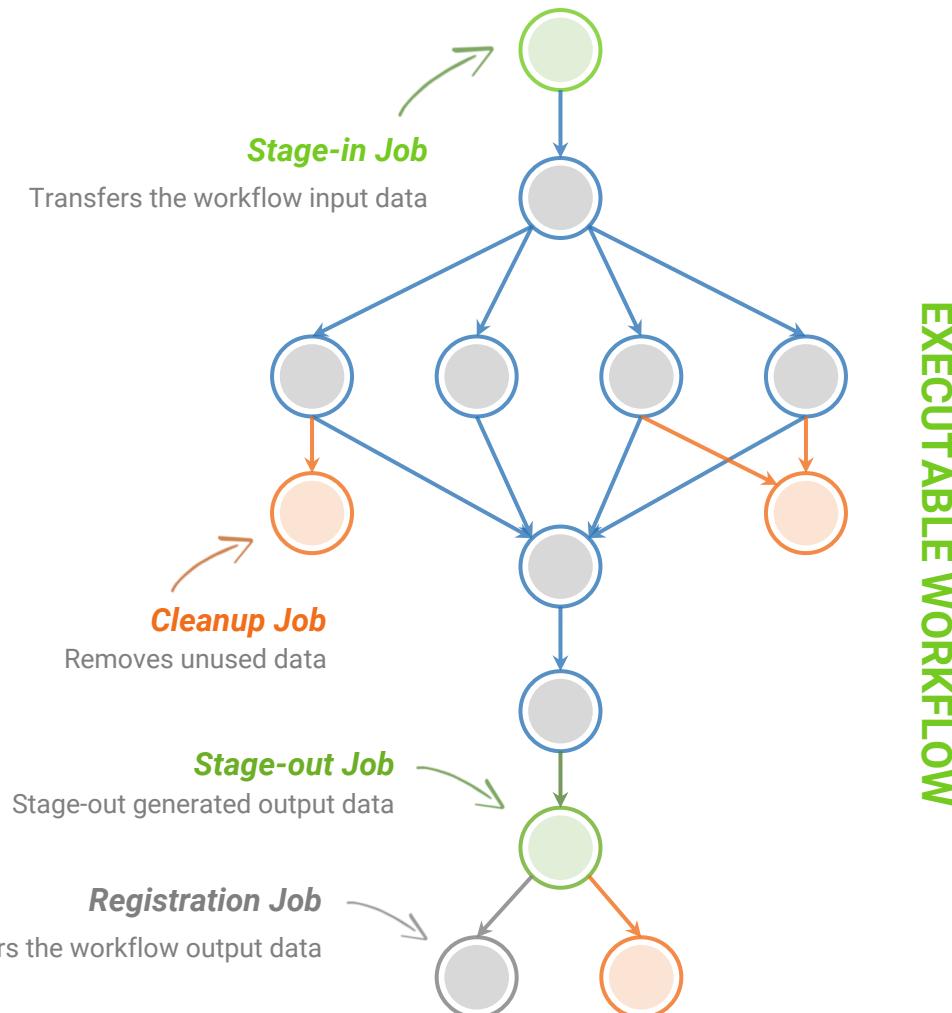
Users do not worry about low level execution details

### ABSTRACT WORKFLOW

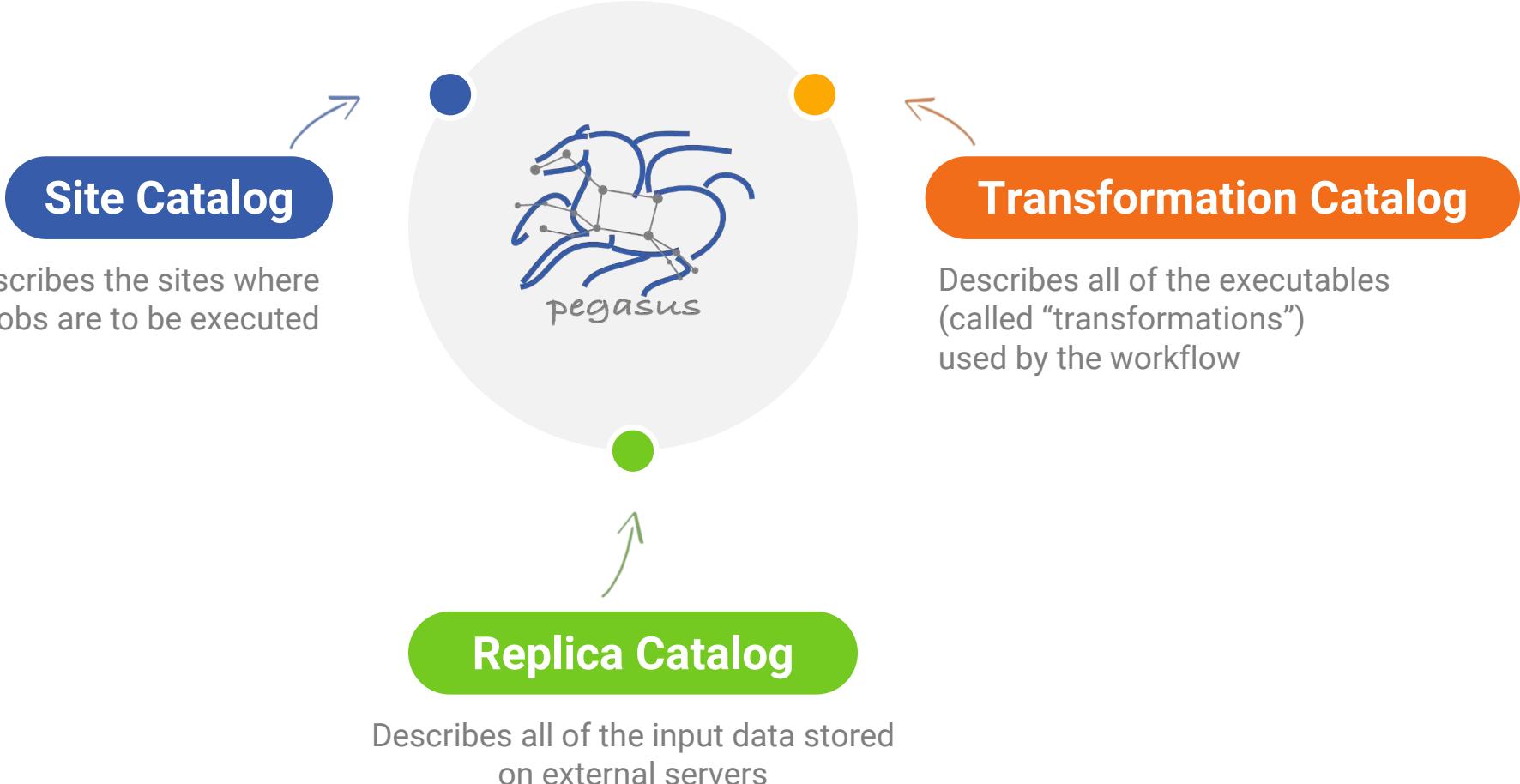


directed-acyclic graphs

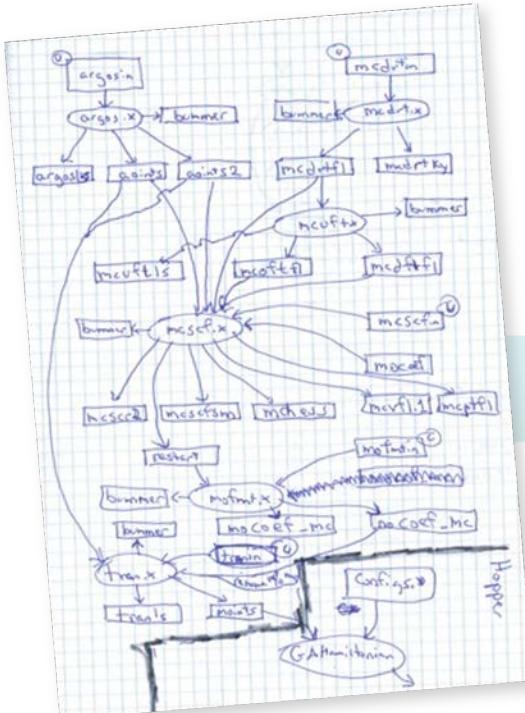
## Output Workflow



# So, what information does Pegasus need?



# Pegasus also provides tools to generate the Abstract Workflow



```

#!/usr/bin/env python3

import os
import logging
from pathlib import Path
from argparse import ArgumentParser

logging.basicConfig(level=logging.DEBUG)

# --- Import Pegasus API -----
from Pegasus.api import *

# --- Create Abstract Workflow -----
wf = Workflow("pipeline")

webpage = File("pegasus.html")

# --- Create Parent Job -----
curl_job = (
    Job("curl")
    .add_args("-o", webpage, "http://pegasus.isi.edu")
    .add_outputs(webpage, stage_out=False, register_replica=False)
)

count = File("count.txt")

# --- Create Dependent Job -----
wc_job = (
    Job("wc")
    .add_args("-l", webpage)
    .add_inputs(webpage)
    .set_stdout(count, stage_out=True, register_replica=True)
)

# --- Add jobs to the Abstract Workflow -----
wf.add_jobs(curl_job, wc_job)

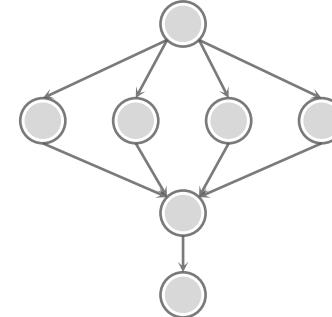
# --- Add control flow dependency -----
wf.add_dependency(wc_job, parents=[curl_job])

# --- Write out the Abstract Workflow -----
wf.write()

```



```
x-pegasus:  
  apilang: python  
  createdBy: vahi  
  createdOn: 11-19-20T14:57:58Z  
pegasus: '5.0'  
name: pipeline  
jobs:  
- type: job  
  name: curl  
  id: ID00000001  
  arguments:  
    - -o  
    - pegasus.html  
    - http://pegasus.isi.edu  
  uses:  
    - lfn: pegasus.html  
      type: output  
      stageOut: false  
      registerReplica: false  
- type: job  
  name: wc  
  id: ID00000002  
  stdout: count.txt  
  arguments:  
    - -l  
    - pegasus.html  
  uses:  
    - lfn: count.txt  
      type: output  
      stageOut: true  
      registerReplica: true  
    - lfn: pegasus.html  
      type: input  
jobDependencies:  
- id: ID00000001  
  children:  
  - ID00000002
```



# YAML Formatted

# Abstract Workflow

# 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()
```

# Pegasus Deployment



## Workflow Submit Node

- Pegasus WMS
- HTCondor

## One or more Compute Sites

- Compute Clusters
- Cloud
- OSG

## Input Sites

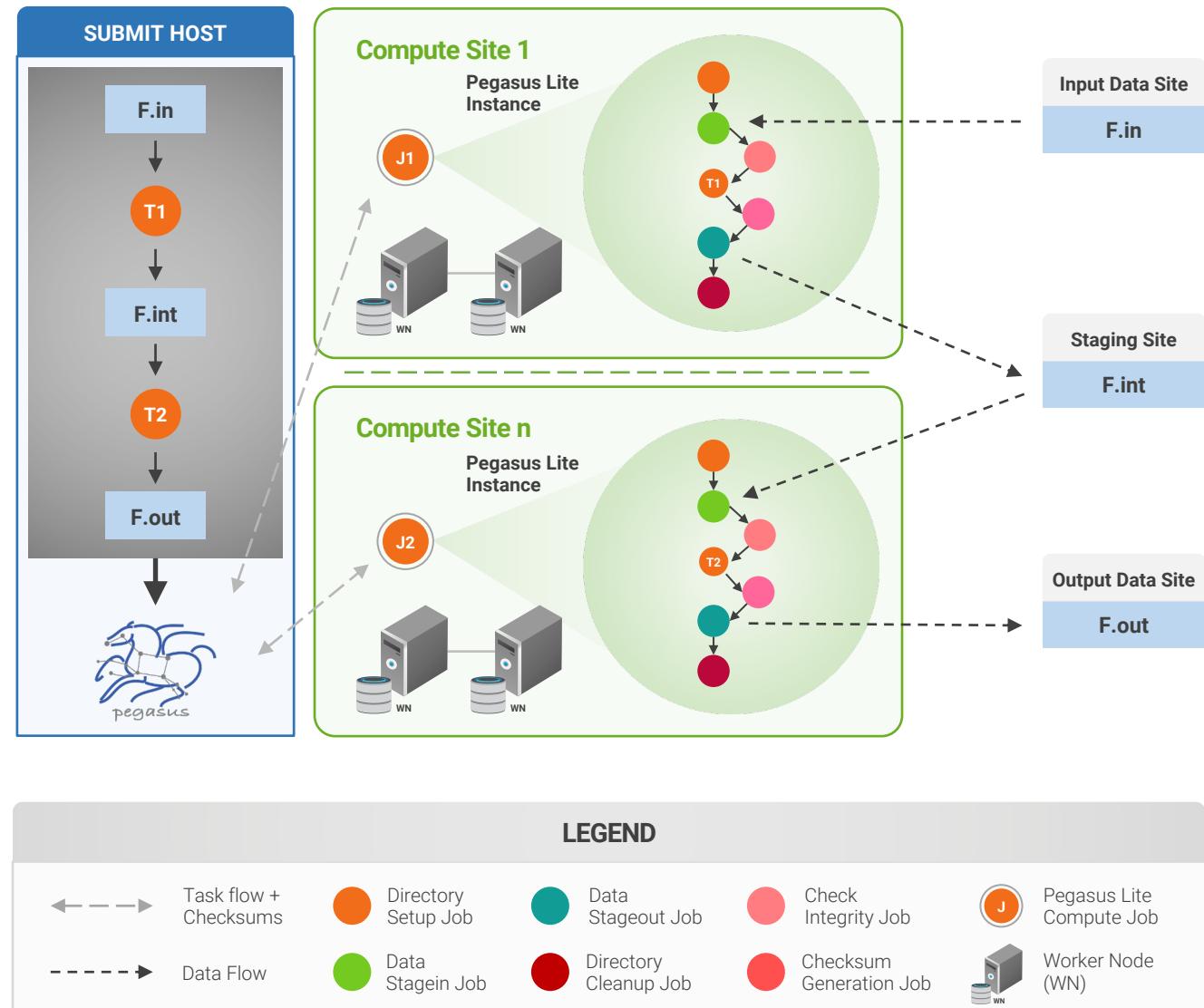
- Host Input Data

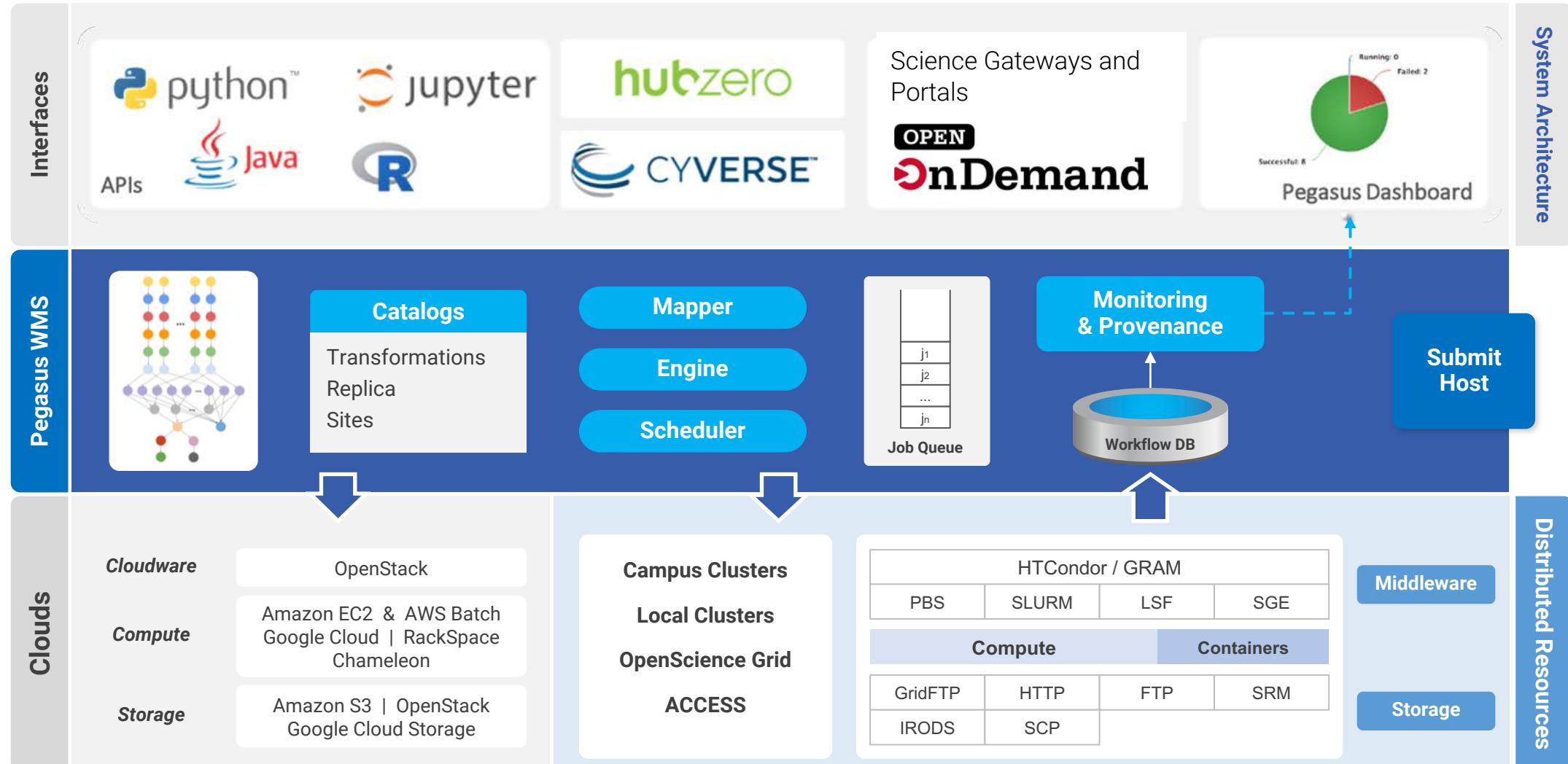
## Data Staging Site

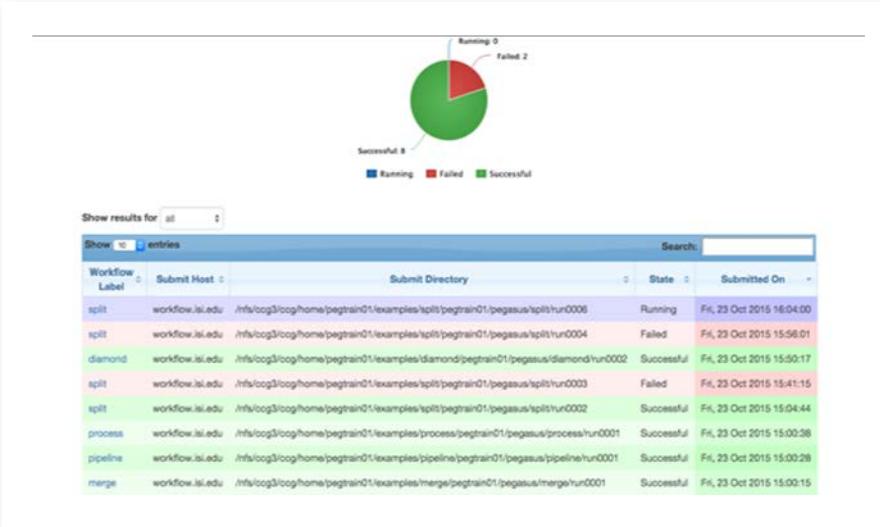
- Coordinate data movement for workflow

## Output Site

- Where output data is placed







Workflow Statistics:

Workflow Label	Submit Host	Submit Directory	State	Submitted On
split	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/split/pegtrain01/pegasus/split/run0006	Running	Fri, 23 Oct 2015 16:04:00
split	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/split/pegtrain01/pegasus/split/run0004	Failed	Fri, 23 Oct 2015 15:56:01
diamond	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/diamond/pegtrain01/pegasus/diamond/run0002	Successful	Fri, 23 Oct 2015 15:50:17
split	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/split/pegtrain01/pegasus/split/run0003	Failed	Fri, 23 Oct 2015 15:41:15
split	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/split/pegtrain01/pegasus/split/run0002	Successful	Fri, 23 Oct 2015 15:04:44
process	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/process/pegtrain01/pegasus/process/run0001	Successful	Fri, 23 Oct 2015 15:00:36
pipeline	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/pipeline/pegtrain01/pegasus/pipeline/run0001	Successful	Fri, 23 Oct 2015 15:00:28
merge	workflow.isi.edu	/nfs/cog3/cog/home/pegtrain01/examples/merge/pegtrain01/pegasus/merge/run0001	Successful	Fri, 23 Oct 2015 15:00:15

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

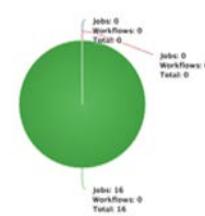
Workflow ID: 5bb4de1d-e986-42b8-9160-ab9488494ecf

Label	split
Type	root-wf
Progress	Successful
Submit Host	workflow.isi.edu
User	pegtrain01
Submit Directory	/nfs/cog3/cog/home/pegtrain01/examples/split/split/run0006
DAGMan Out File	split-0.dag.dagman.out
Wall Time	12 mins 23 secs
Cumulative Wall Time	9 mins 34 secs

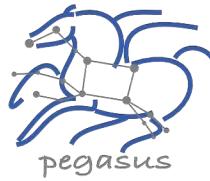
Job Status (Entire Workflow)



Job Status (Per Workflow)



<https://pegasus.isi.edu>



# 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 Waittime as seen from Submit Side	9 mins 35 secs
Workflow Cumulative Badput Time	9 mins 23 secs
Cumulative Job Badput Waittime as seen from Submit Side	9 mins 20 secs
Workflow Retries	1

### Workflow Statistics

This Workflow						
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
Entire Workflow						
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

- Job Breakdown Statistics
- Job Statistics

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?



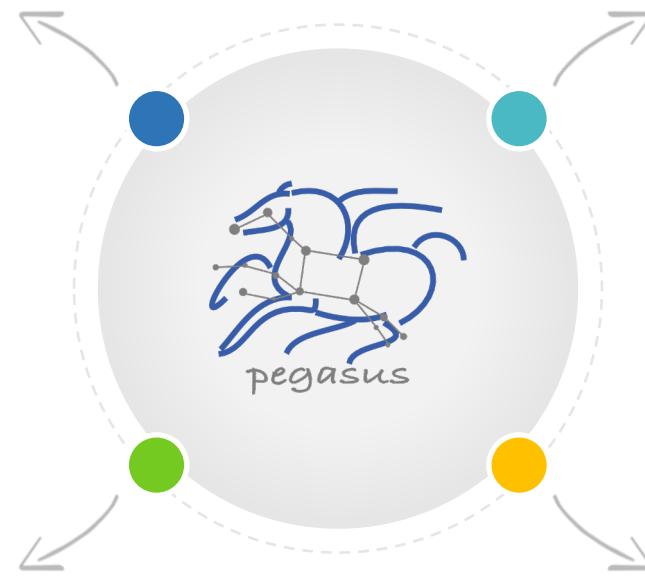
## Postscript

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



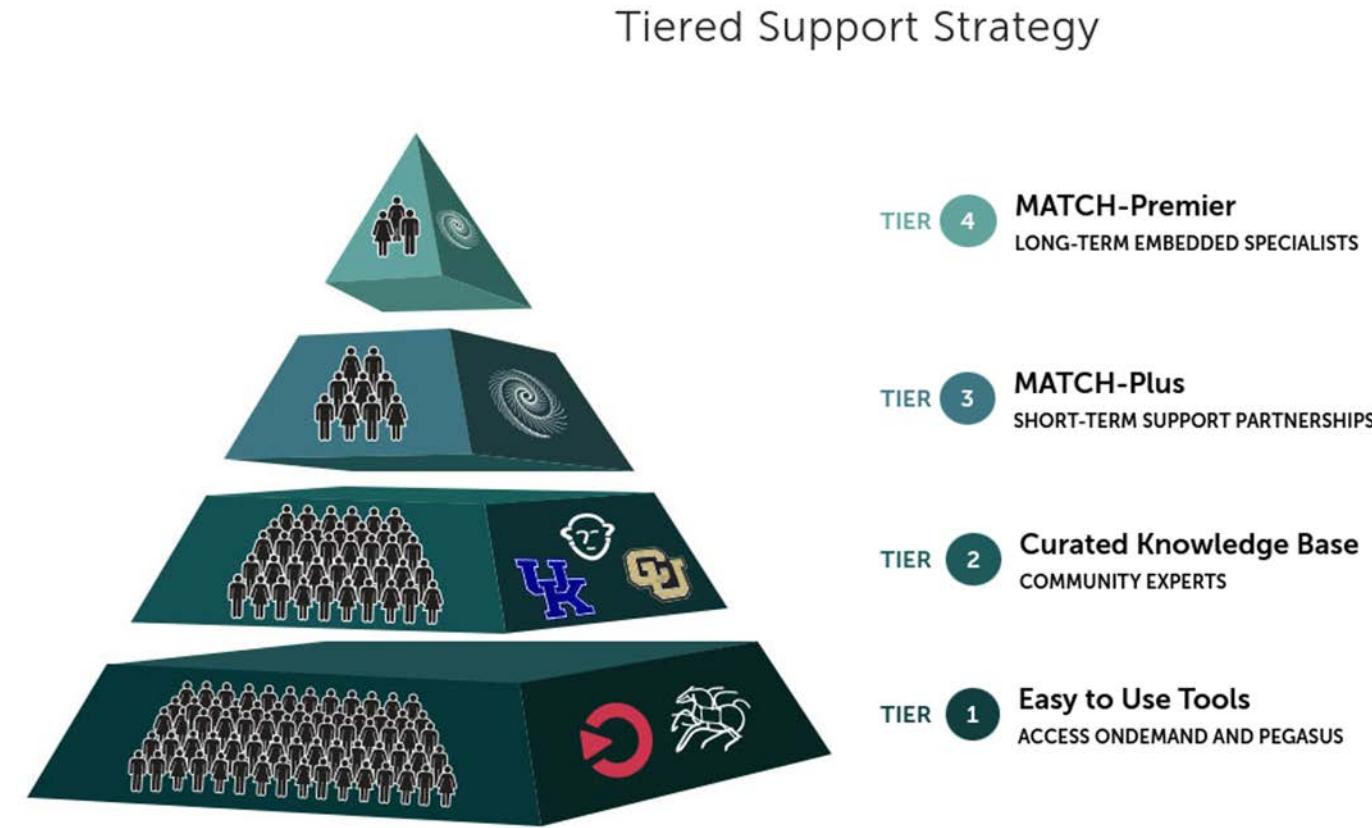
ACCESS...

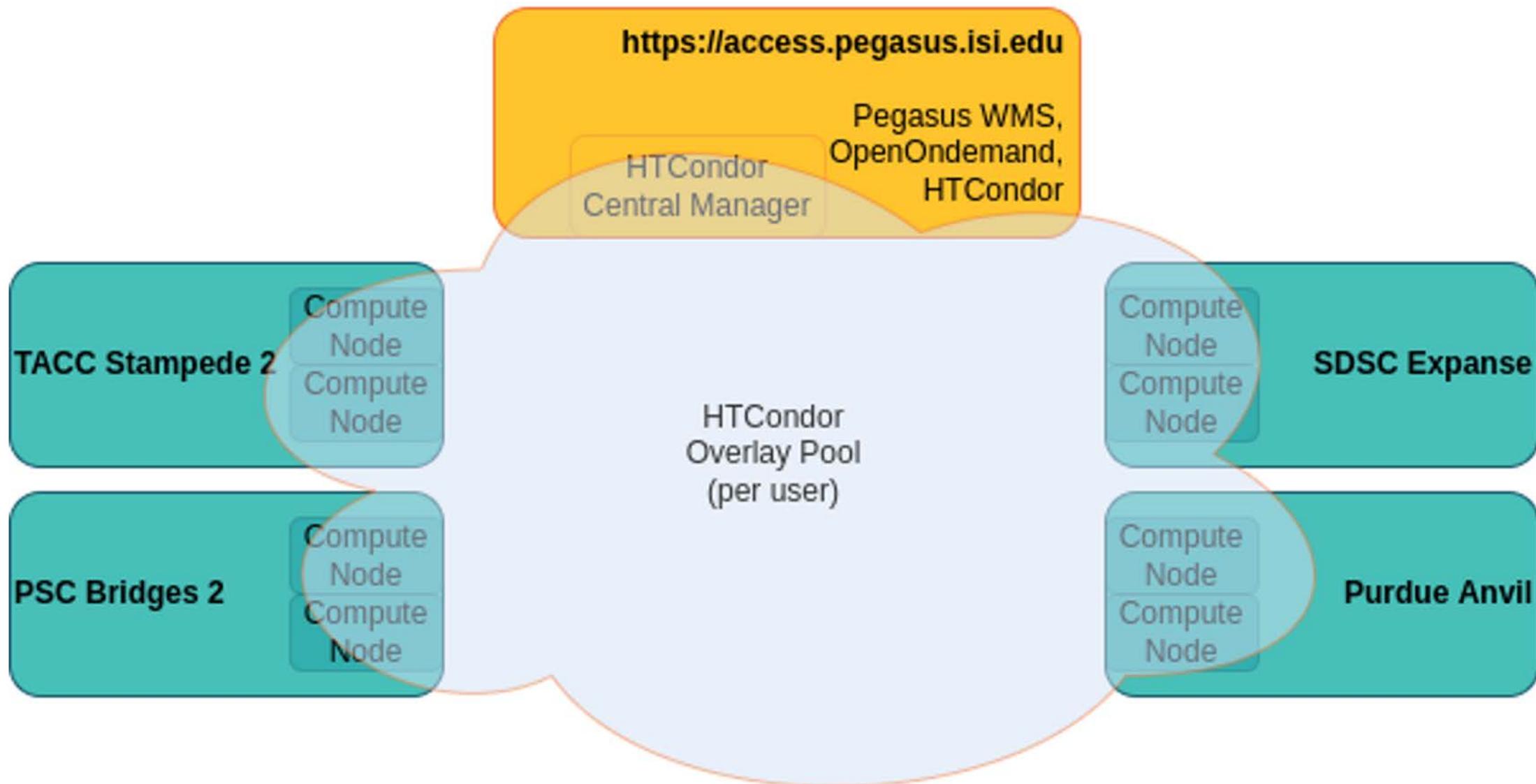
# Pegasus is part of the ACCESS support strategy

Pegasus will be used as a tier 1 tool

**Central Open OnDemand instance with Pegasus, HTCondor and Jupyter**

It will be easy to run HTC workflows across ACCESS sites





# HTCondor Annex / Pilot Jobs

- **A pilot can run multiple user jobs** - it stays active until no more user jobs are available or until end of life has been reached, whichever comes first.
- **A pilot is partitionable** - job slots will dynamically be created based on the resource requirements in the user jobs. This means you can fit multiple user jobs on a compute node at the same time.
- **A pilot will only run jobs for the user who started it.**

## Understanding Pegasus Features...



# 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



# Data Staging Configurations

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

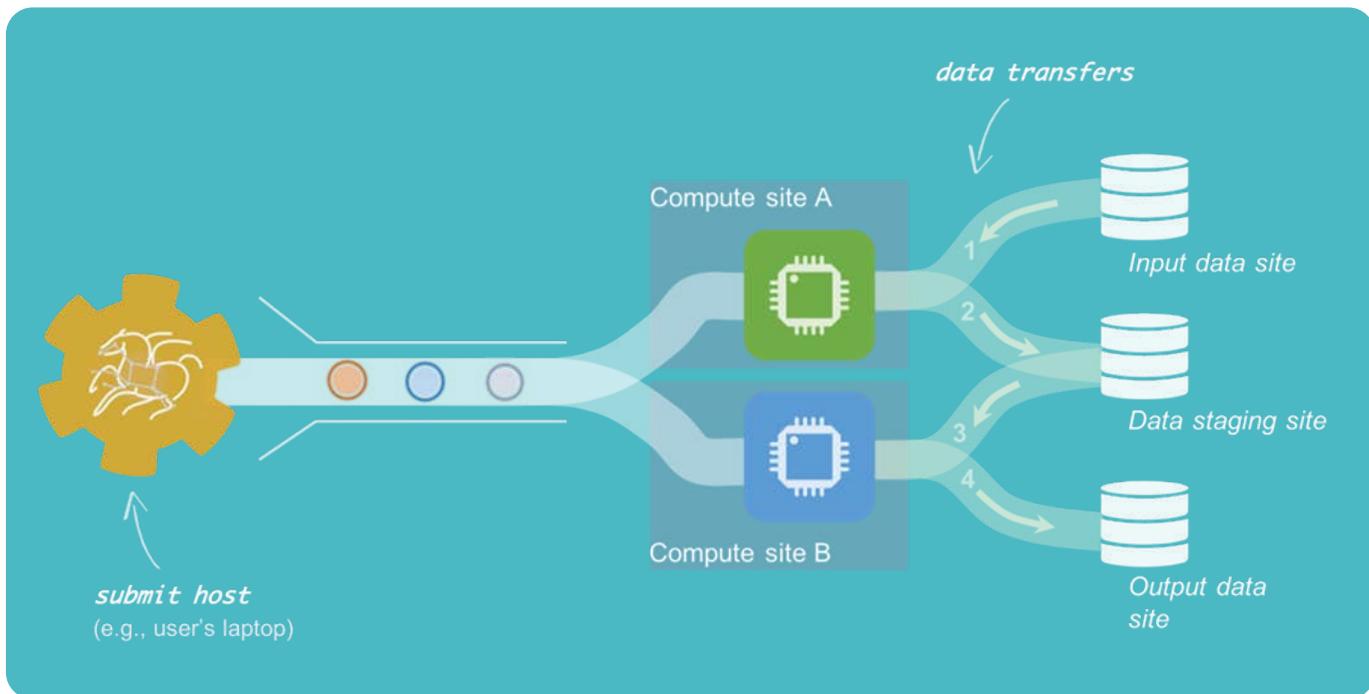
- Worker nodes do not share a file system
- Data is pulled from / pushed to the submit host via HTCondor file transfers
- Staging site is the submit host

## 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, ...)

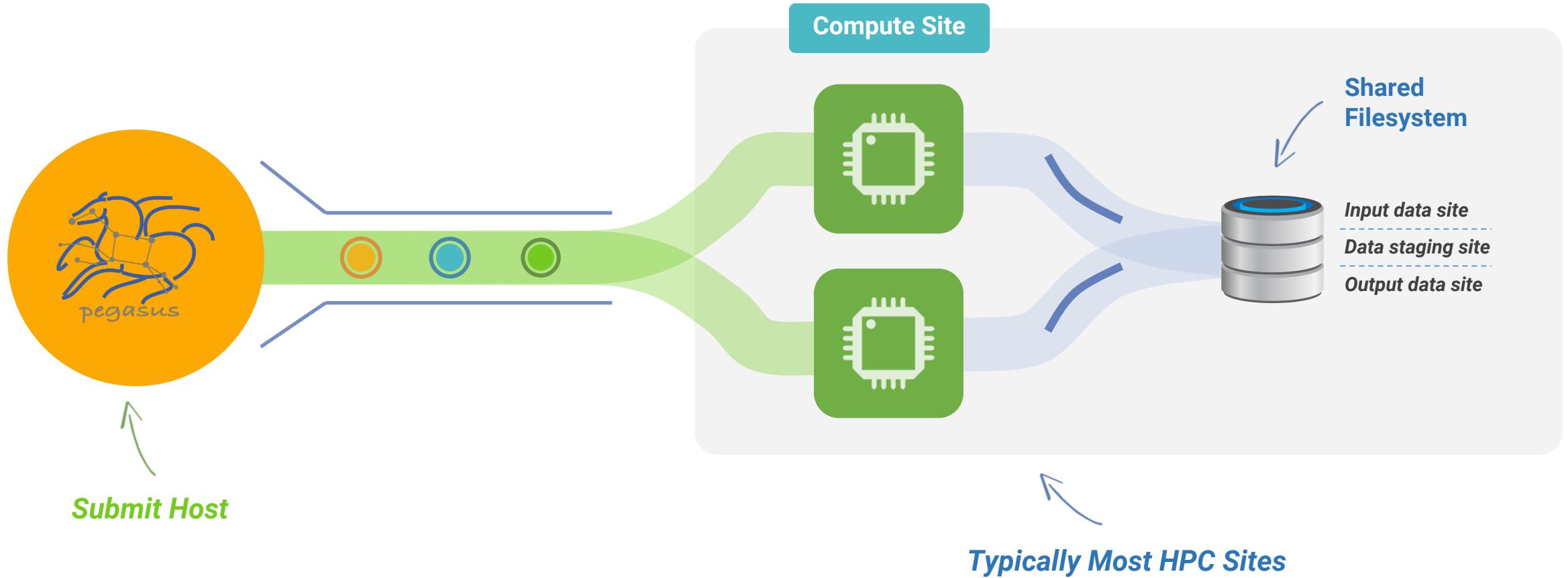
- I/O is directly against the shared file system





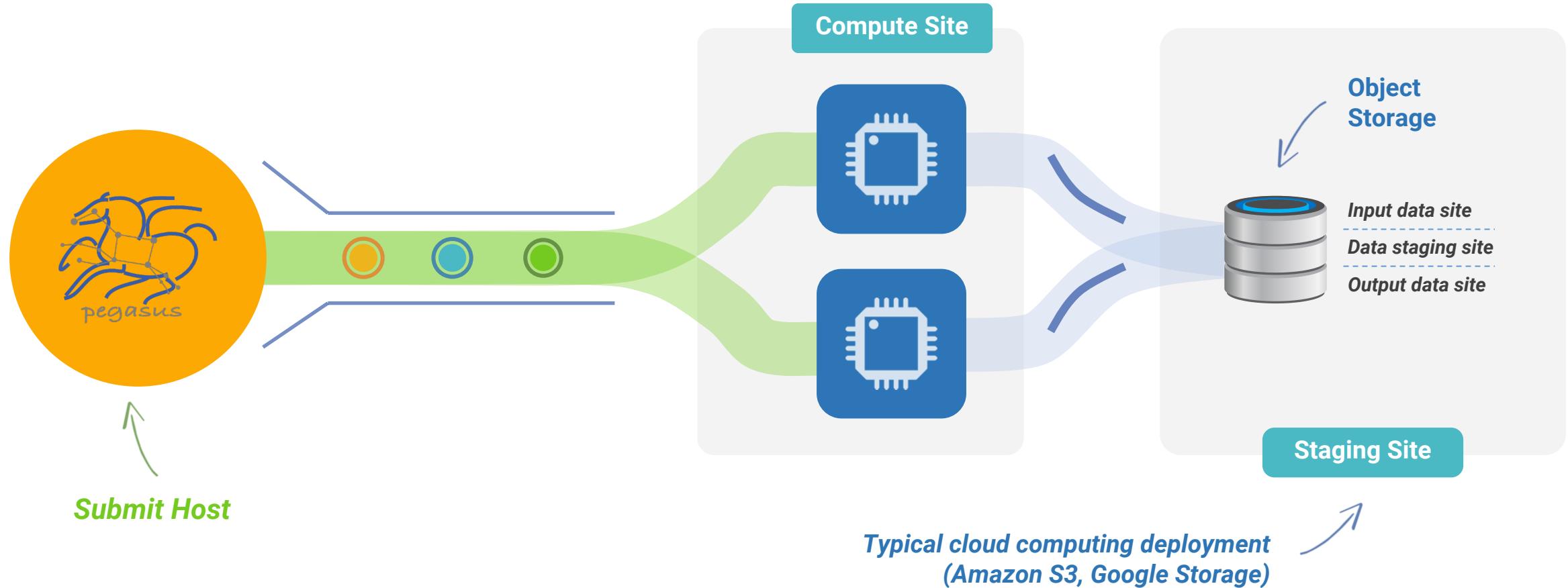
# High Performance Computing

There are several possible configurations...



# Cloud Computing

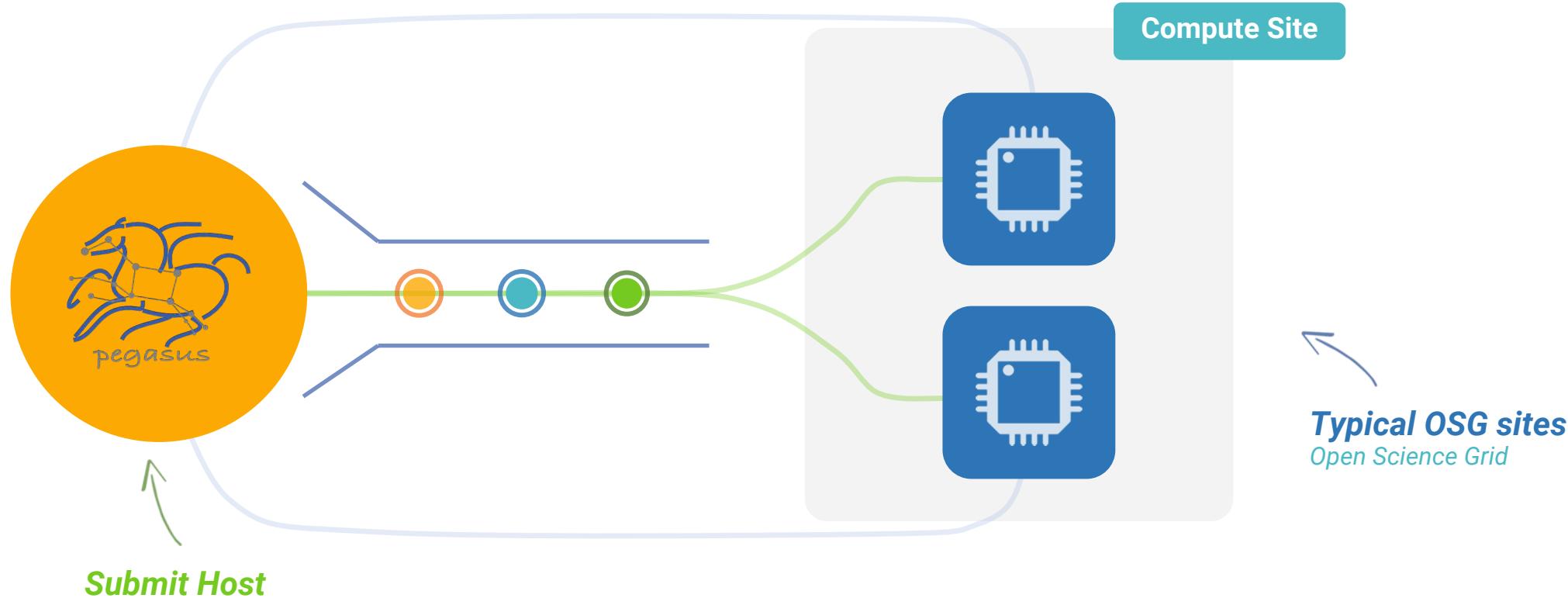
High-scalable object storages



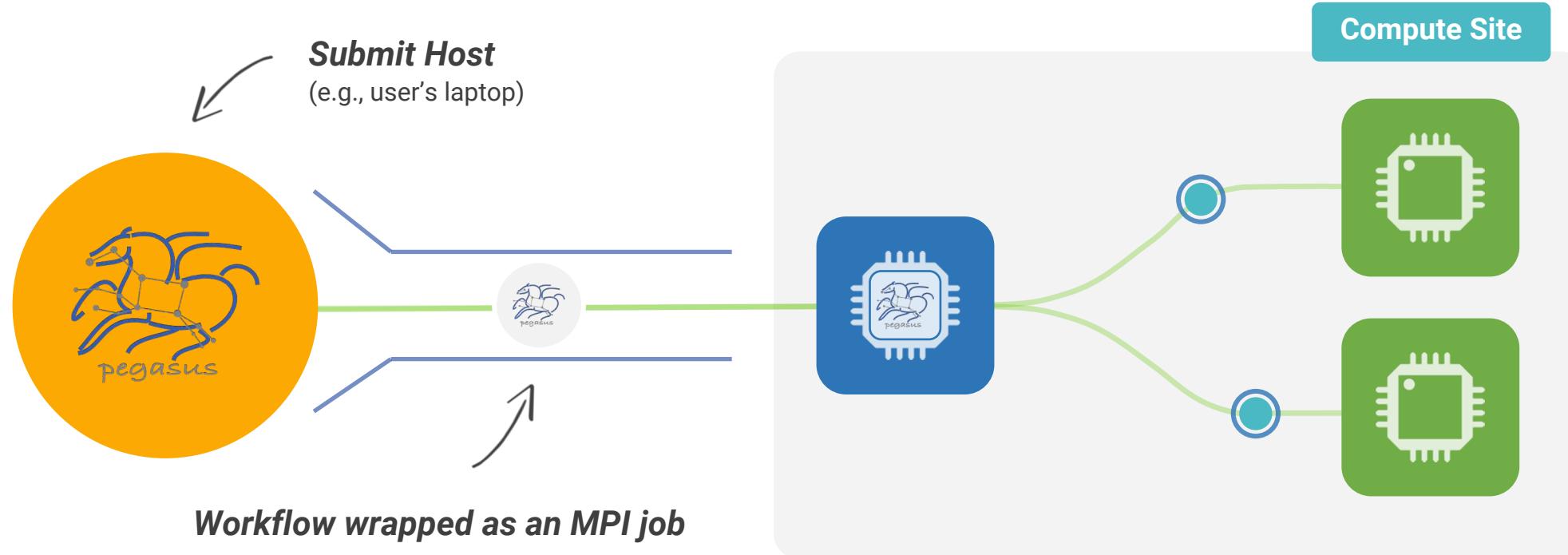


# Grid Computing

## Local data management



# Running fine-grained workflows on HPC systems...



# Pegasus Container Support



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



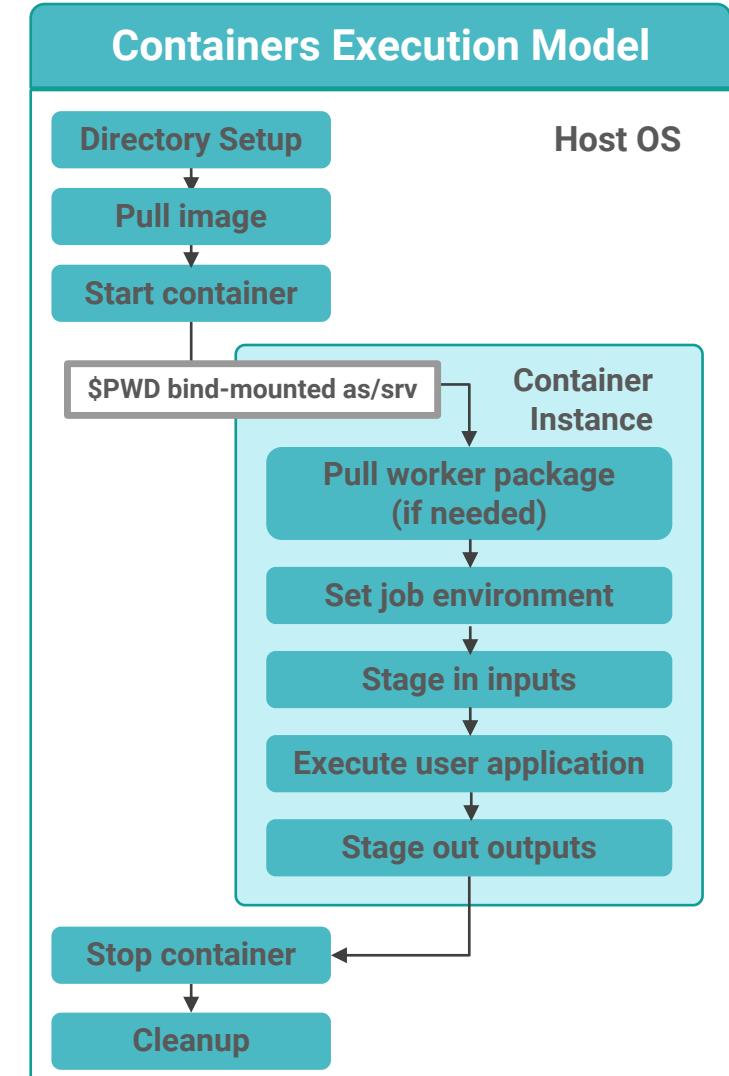
Users can **refer** to a **container** they want to **use – Pegasus stages** their executables and containers to the node

- 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

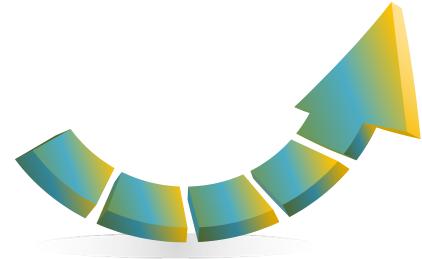
- Staged to compute node if not present
- Docker or Singularity Hub URL's
- Docker Image exported as a TAR file and available at a server, just like any other input dataset

## 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

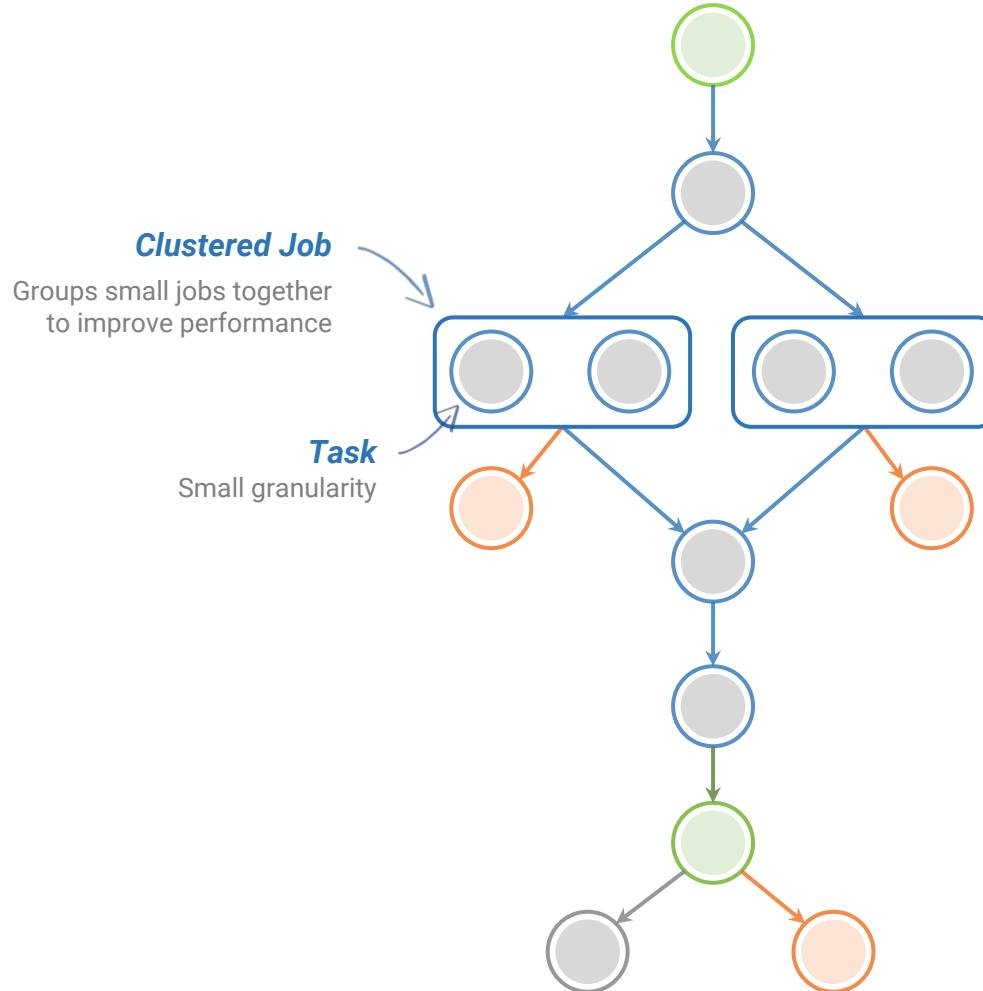
## Other Optimizations

- 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

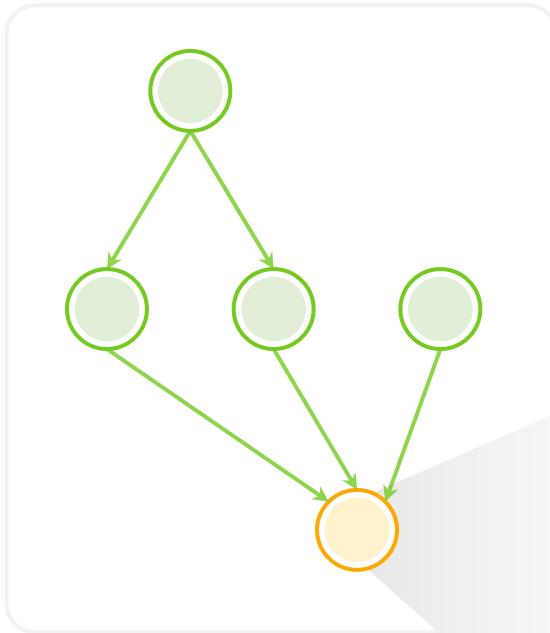
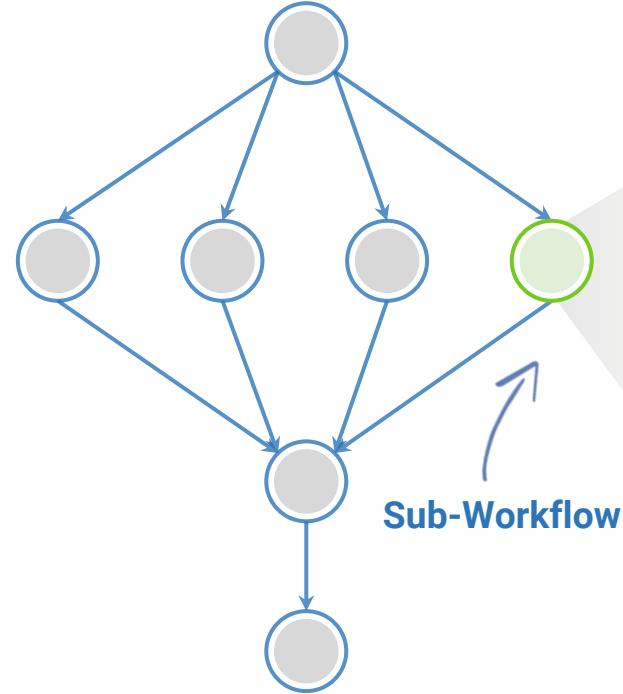


# Performance.

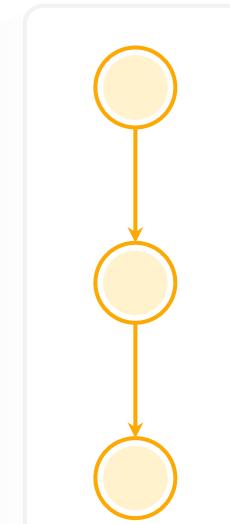
## Why not improve it?



# Pegasus also handles large-scale workflows

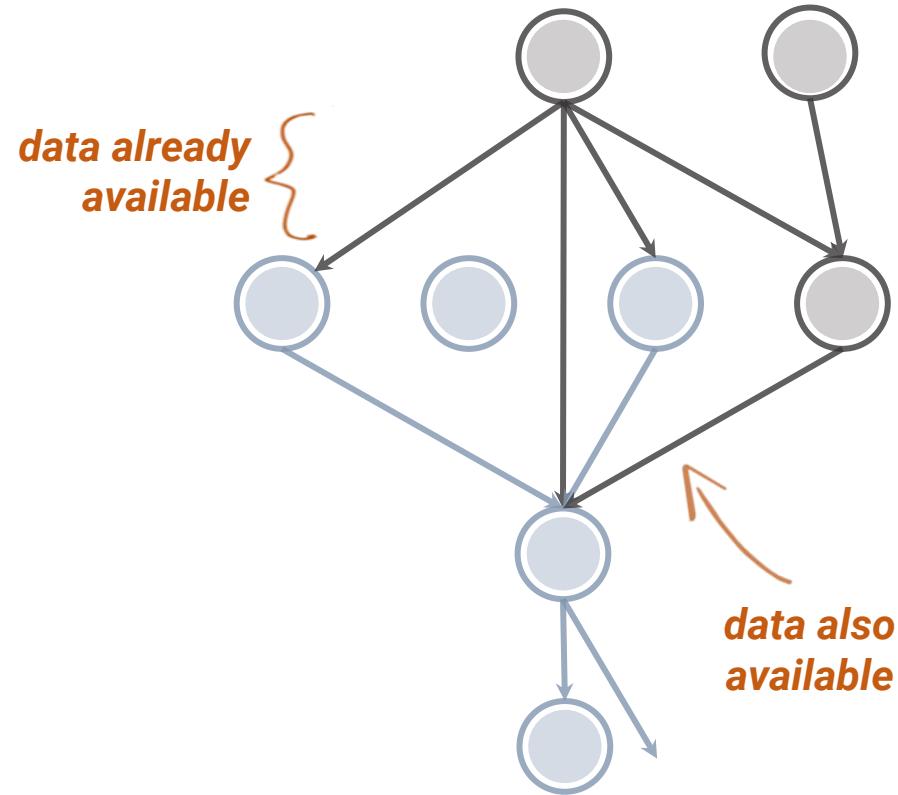


Sub-Workflow

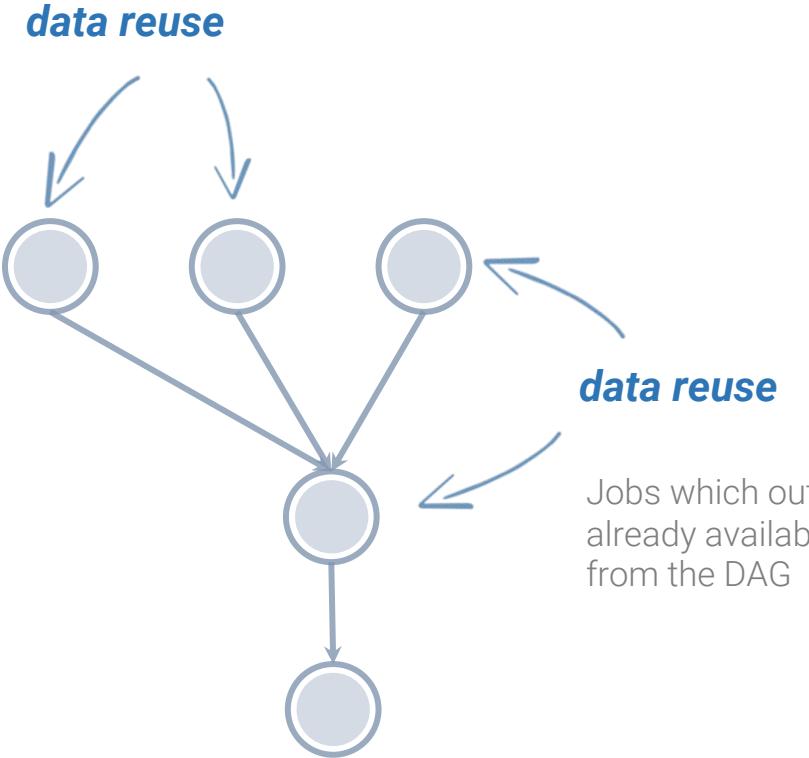


Recursion ends  
When abstract  
workflow with  
only compute jobs  
is encountered

# Data Reuse prune jobs if output data already exists



workflow  
reduction



# And if a job fails?



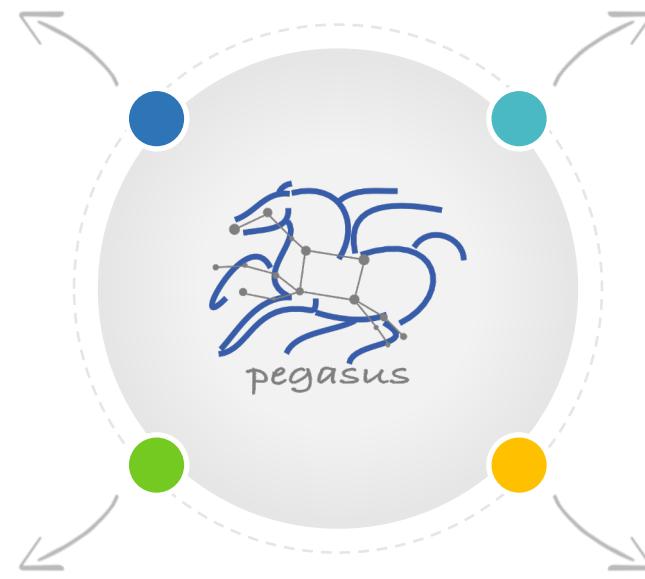
## Postscript

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

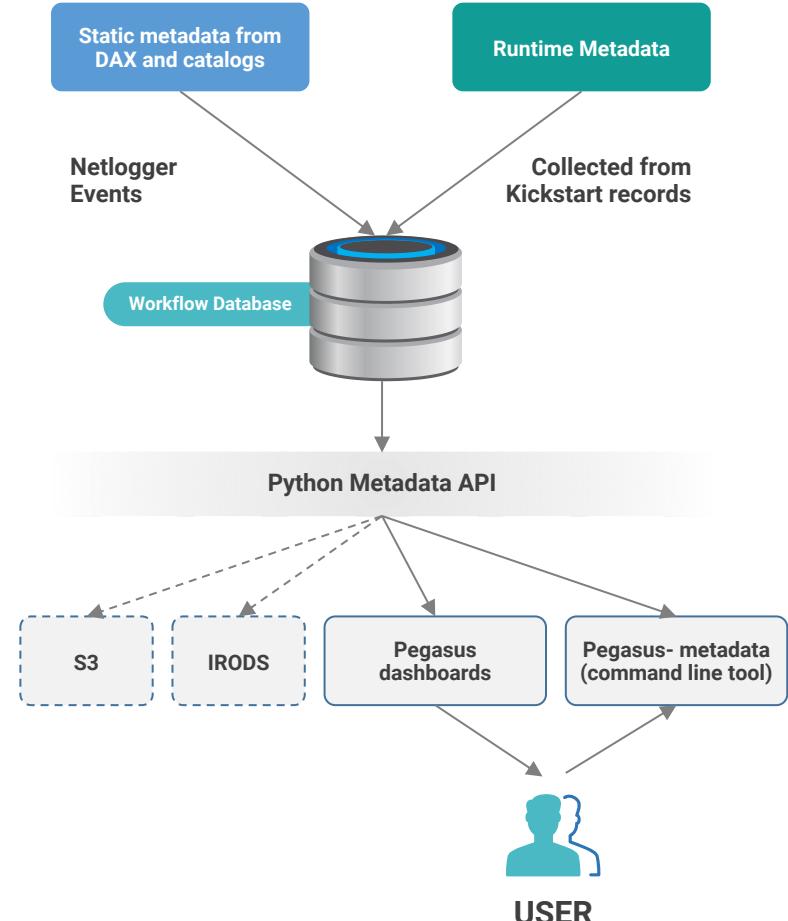
```
x-pegasus:  
apiLang: python  
createdBy: vahi  
createdOn: 12-08-20T10:08:48Z  
pegasus: "5.0"  
name: diamond  
metadata:  
    experiment:"par_all27_prot_lipid"  
jobs:  
- type: "job"  
  name: "namd"  
  id: "ID0000001"  
  arguments: ["equilibrate.conf"]  
  uses:  
  - lfn: "Q42.psf"  
    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?”

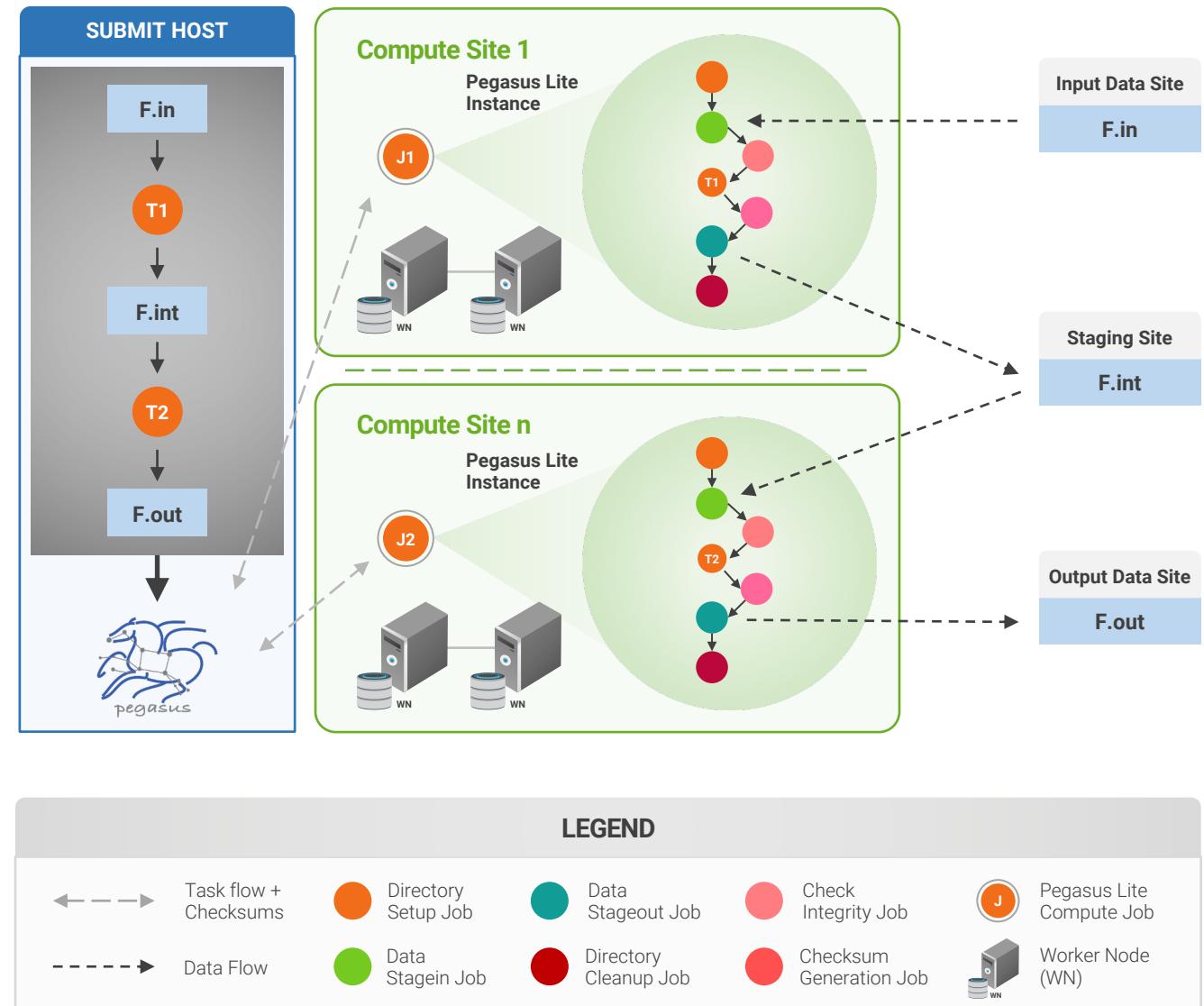
# Automatic Integrity Checking in Pegasus

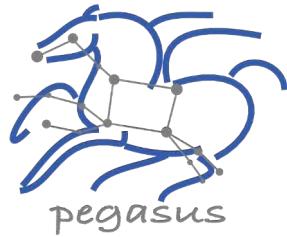


Pegasus performs integrity checksums on input files right before a job starts on the remote node.

- 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





# Pegasus

est. 2001

Automate, recover, and debug scientific computations.

## ► Get Started

### ► Pegasus Website

<https://pegasus.isi.edu>

### ► Users Mailing List

[pegasus-users@isi.edu](mailto:pegasus-users@isi.edu)

### ► Support

[pegasus-support@isi.edu](mailto:pegasus-support@isi.edu)

### ► Slack

Ask for an invite by trying to join [pegasus-users.slack.com](https://pegasus-users.slack.com) in the Slack app

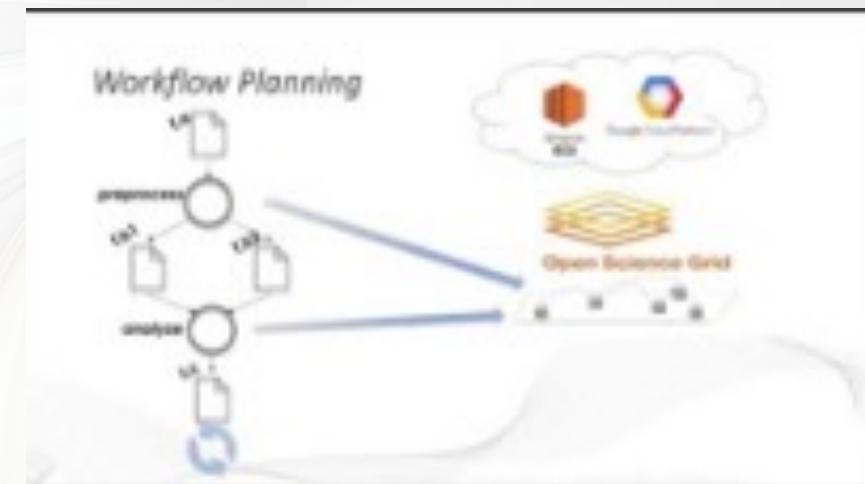
### ► Pegasus Online Office Hours

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



### YouTube Channel

<https://www.youtube.com/channel/UCwJQIn1CqBvTJqiNr9X9F1Q/featured>



[Pegasus in 5 Minutes](#)