

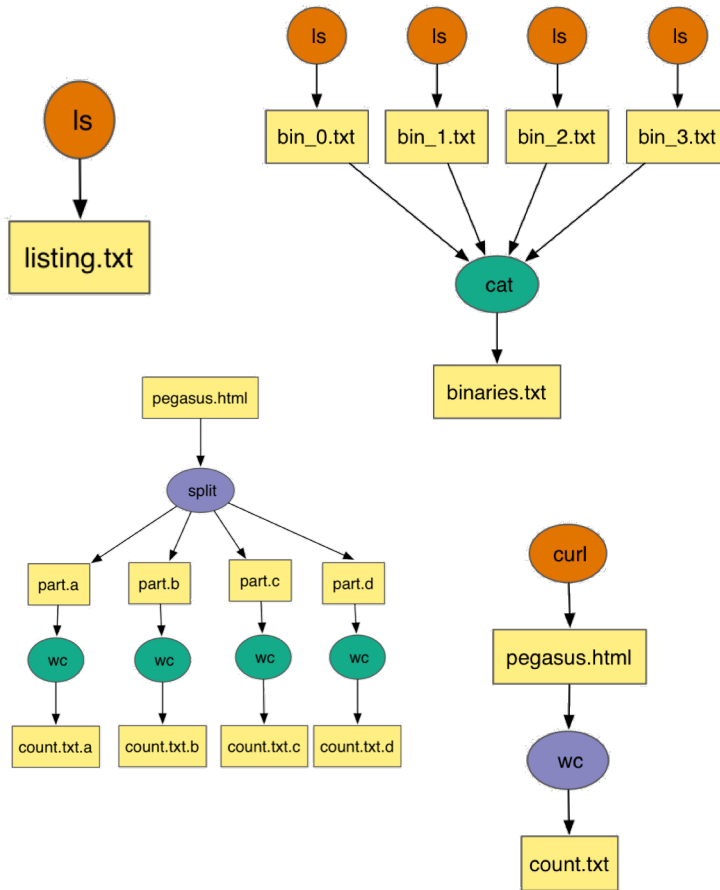
# ***Pegasus Workflow Management System***

*Karan Vahi*

*USC Information Sciences Institute*

# Benefits of Scientific Workflows

(from the point of view of an application scientist)



- 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.
  - Avoids need to download-install-learn how to use someone else's code.
- 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.

Slide courtesy of David Okaya, SCEC, USC

# Challenges of Workflow Management

## Challenges across domains

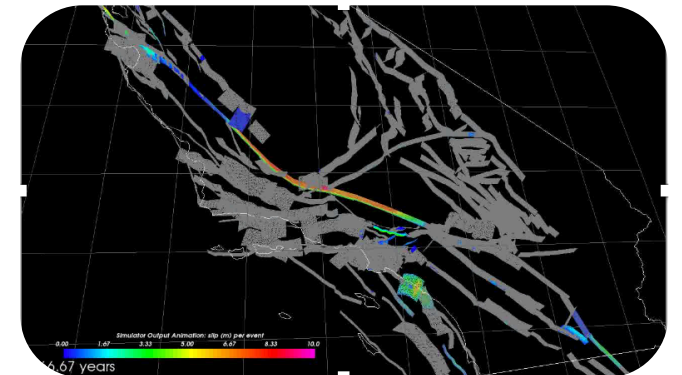
- Need to describe complex workflows in a simple way
- Need to access distributed, heterogeneous data and resources (heterogeneous interfaces)
- Need to deal with resources/software that change over time

## Our focus

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



Sky mosaic, IPAC, Caltech



Earthquake simulation, SCEC, USC

# Pegasus Workflow Management System

- Operates at the level of files and individual applications
- Allows scientists to describe their computational processes (workflows) at a logical level
- Without including details of target heterogeneous CI (portability)
- Scalable to  $O(10^6)$  tasks, TBs of data
- Captures provenance and supports reproducibility
- Includes monitoring and debugging tools

## Workflow Listing



Show results for all 1

Workflow Label	Submit Host	Submit Directory	State	Submitted On
split	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/run0008	Running	Fri, 23 Oct 2015 16:04:00
split	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/run0004	Failed	Fri, 23 Oct 2015 16:06:01
diamond	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/diamond/run0002	Successful	Fri, 23 Oct 2015 16:05:17
split	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/run0003	Failed	Fri, 23 Oct 2015 16:04:15
split	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/run0002	Successful	Fri, 23 Oct 2015 16:04:44
process	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/process/run0001	Successful	Fri, 23 Oct 2015 16:00:28
pipeline	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/pipeline/run0001	Successful	Fri, 23 Oct 2015 16:00:15
merge	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/merge/run0001	Successful	Fri, 23 Oct 2015 16:00:06
diamond	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/diamond/run0001	Successful	Fri, 23 Oct 2015 16:00:06
split	workflow.isl.edu	/ifs/cog3/cog/home/pegasus01/learningsplit/pegasus01/pegasus01/run0001	Successful	Fri, 23 Oct 2015 16:00:06

Showing 1 to 10 of 10 entries

## 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

## Job Breakdown Statistics

## Job Statistics

Composition in Python, R, Java, Perl, Jupyter Notebook

hubzero

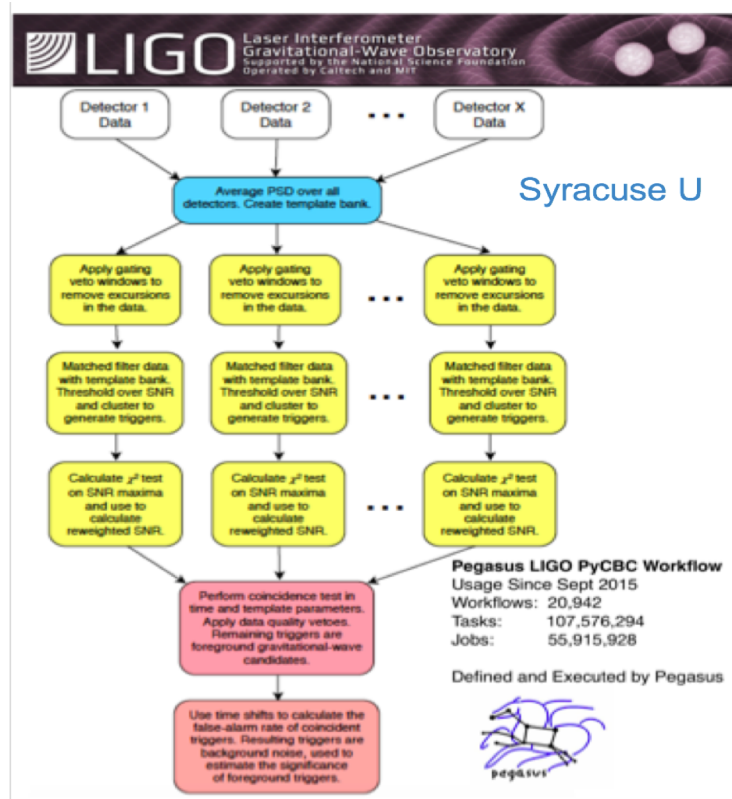


# Pegasus Workflow Management System, Production Use

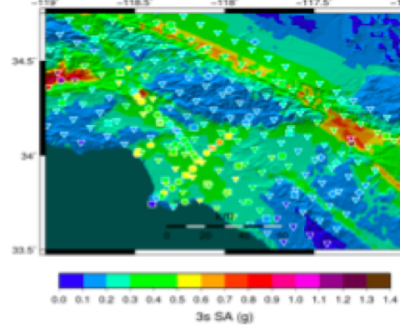


Last 12 months: Pegasus users ran **240K** workflows, **145M** jobs

Majority of these include data transfers, using LAN, the Internet, local and remote storage



Southern California Earthquake Center, USC



Nek2 Kinase

PDB ID: 2DKA

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Binding Site: A

Other Proteins in Crystal Structure (Chain)

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure

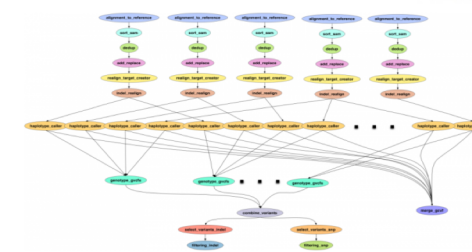
No other chains found for this PDB structure

No other chains found for this PDB structure

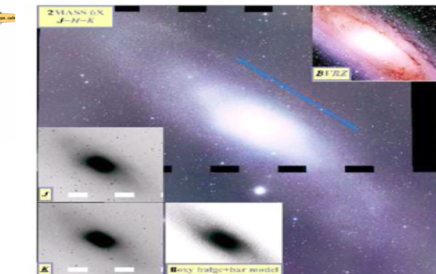
No other chains found for this PDB structure

No other chains found for this PDB structure

No other chains found for this PDB structure



Bioinformatics: Protein interactions, IU



Montage, Caltech

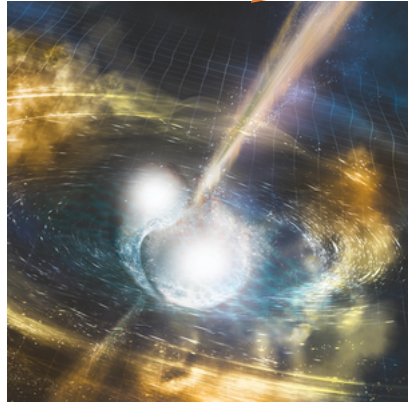
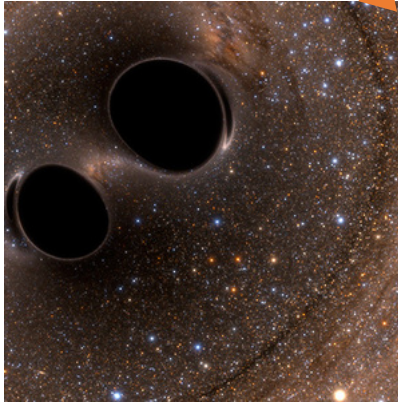
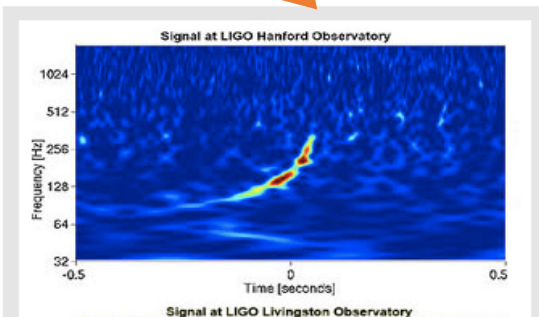
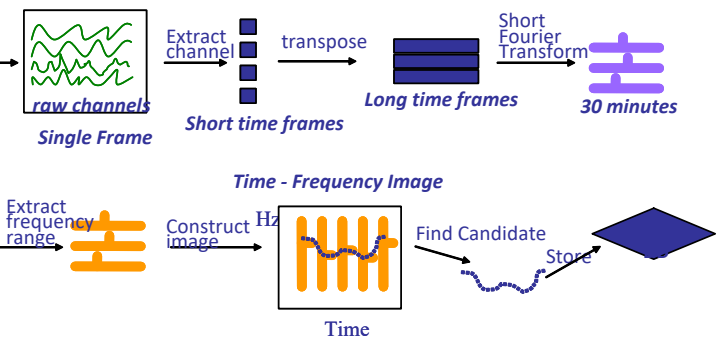
Bioinformatics: SoyKB, University of Arizona

# Pegasus: Grounding Research and Development

**Nobel Prize**



## Working with LIGO



**First Pegasus  
prototype for  
LIGO Pulsar  
Searches**

**Blind injection detection**

**First detection of  
black hole collision**

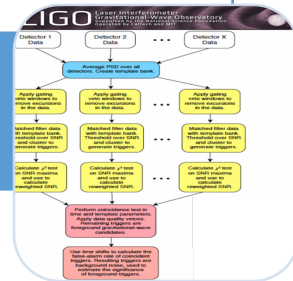
**Multi-messenger  
neutron star  
merger  
observation**

# Complexity of LIGO Workflows

First GW detection:  $\sim 21\text{K}$  Pegasus workflows,  $\sim 107\text{M}$  tasks

Analysis measures the statistical significance of collected data

# Science Workflow



Efficient,  
scalable, and  
robust execution  
of tasks and  
data access

# Automation



LIGO, Open  
Science Grid,  
XSEDE, Blue  
Waters

# Distributed Power



2015/16



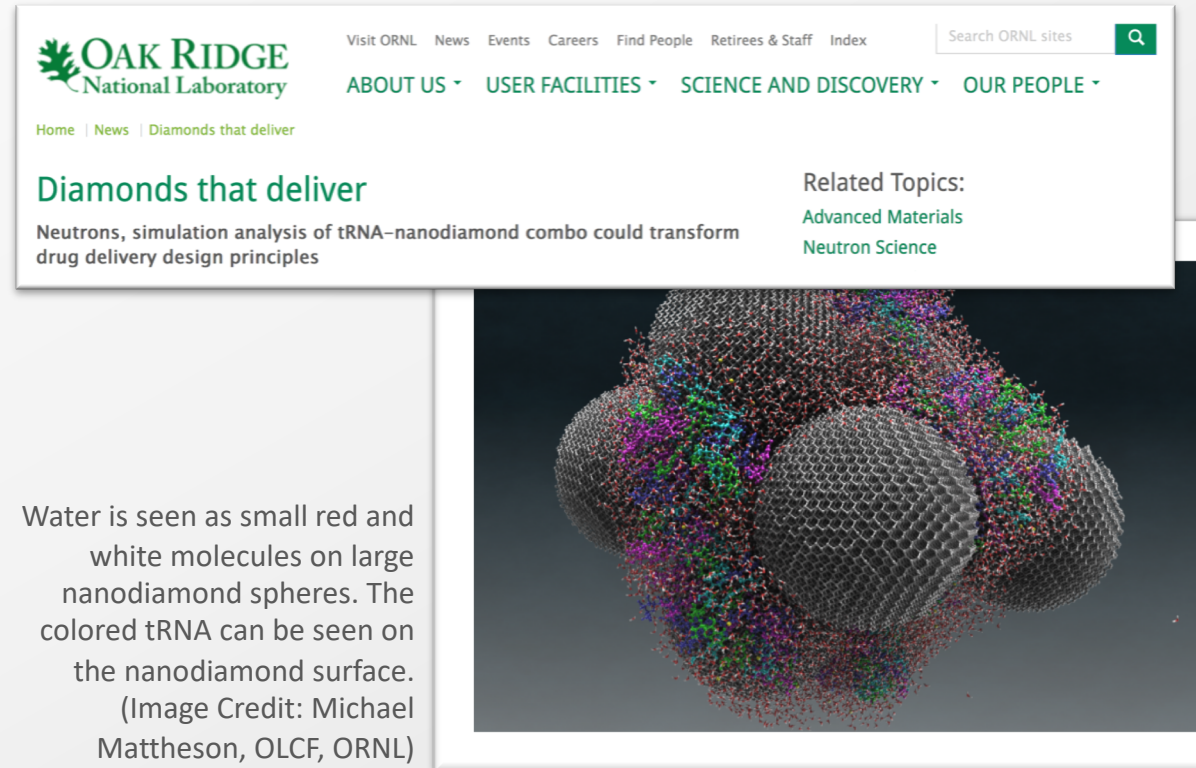
# Impact on DOE Science

Enabled cutting-edge domain science (e.g., drug delivery) through collaboration with scientists at the DoE **Spallation Neutron Source (SNS)** facility

A Pegasus workflow was developed that confirmed that **nanodiamonds** can enhance the dynamics of tRNA

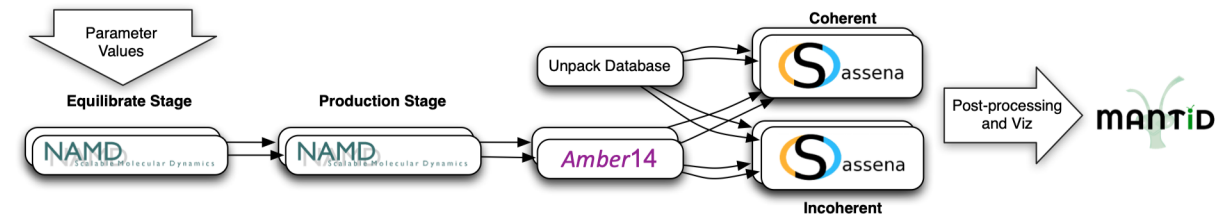
It compared SNS neutron scattering data with MD simulations by calculating the epsilon that best matches experimental data

Ran on a Cray XE6 at NERSC using 400,000 CPU hours, and generated 3TB of data.



The screenshot shows the Oak Ridge National Laboratory website. The header includes the ORNL logo and navigation links: Visit ORNL, News, Events, Careers, Find People, Retirees & Staff, Index, and a search bar. Below the header, there are links for ABOUT US, USER FACILITIES, SCIENCE AND DISCOVERY, and OUR PEOPLE. The main content area features a news article titled "Diamonds that deliver" with the subtext "Neutrons, simulation analysis of tRNA-nanodiamond combo could transform drug delivery design principles". To the right of the article title, there are "Related Topics" listed: "Advanced Materials" and "Neutron Science". Below the text, there is a large 3D visualization of a nanodiamond sphere with water molecules (small red and white spheres) and tRNA molecules (colored spheres) on its surface.

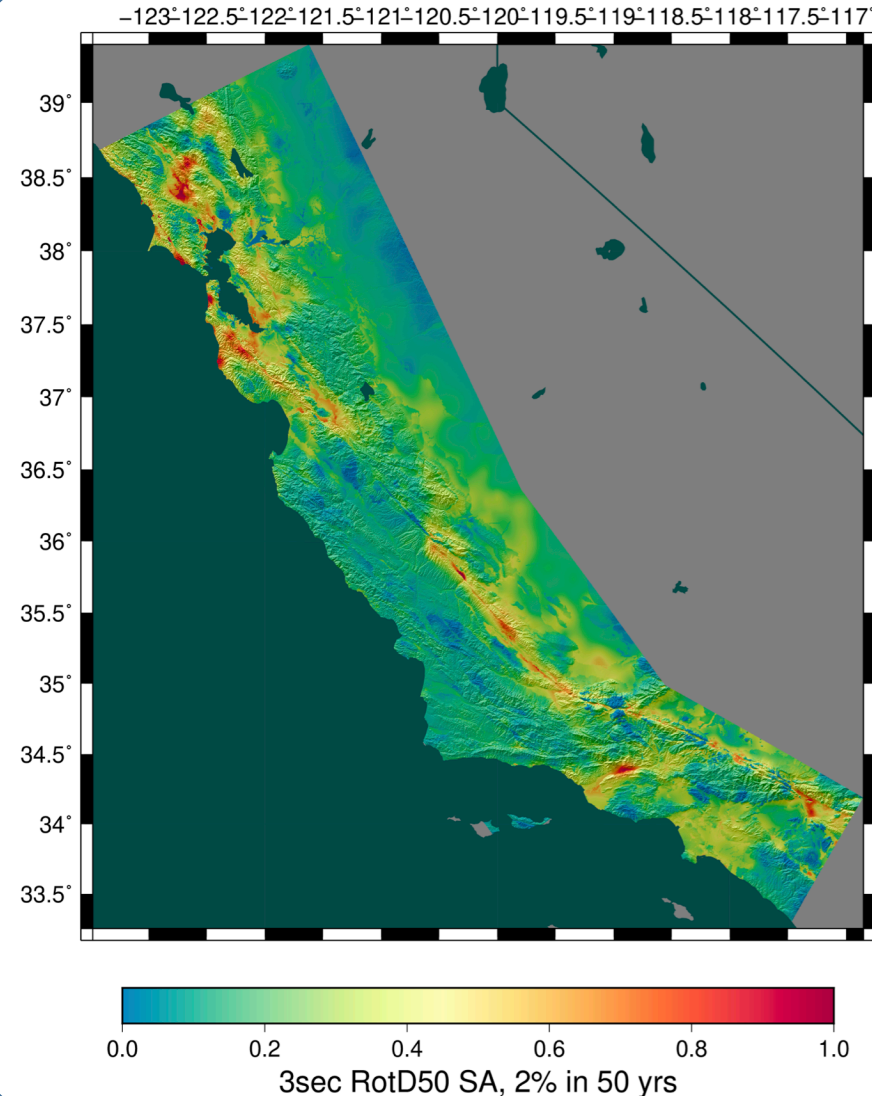
Water is seen as small red and white molecules on large nanodiamond spheres. The colored tRNA can be seen on the nanodiamond surface.  
(Image Credit: Michael Mattheson, OLCF, ORNL)



*An automated analysis workflow for optimization of force-field parameters using neutron scattering data. V. E. Lynch, J. M. Borreguero, D. Bhowmik, P. Ganesh, B. G. Sumpter, T. E. Proffen, M. Goswami, Journal of Computational Physics, July 2017.*

# Supporting Heterogeneous Workflows

**SCEC's  
CyberShake:  
What will the  
peak  
earthquake  
motion be  
over the next  
50 years?**



Useful information for:

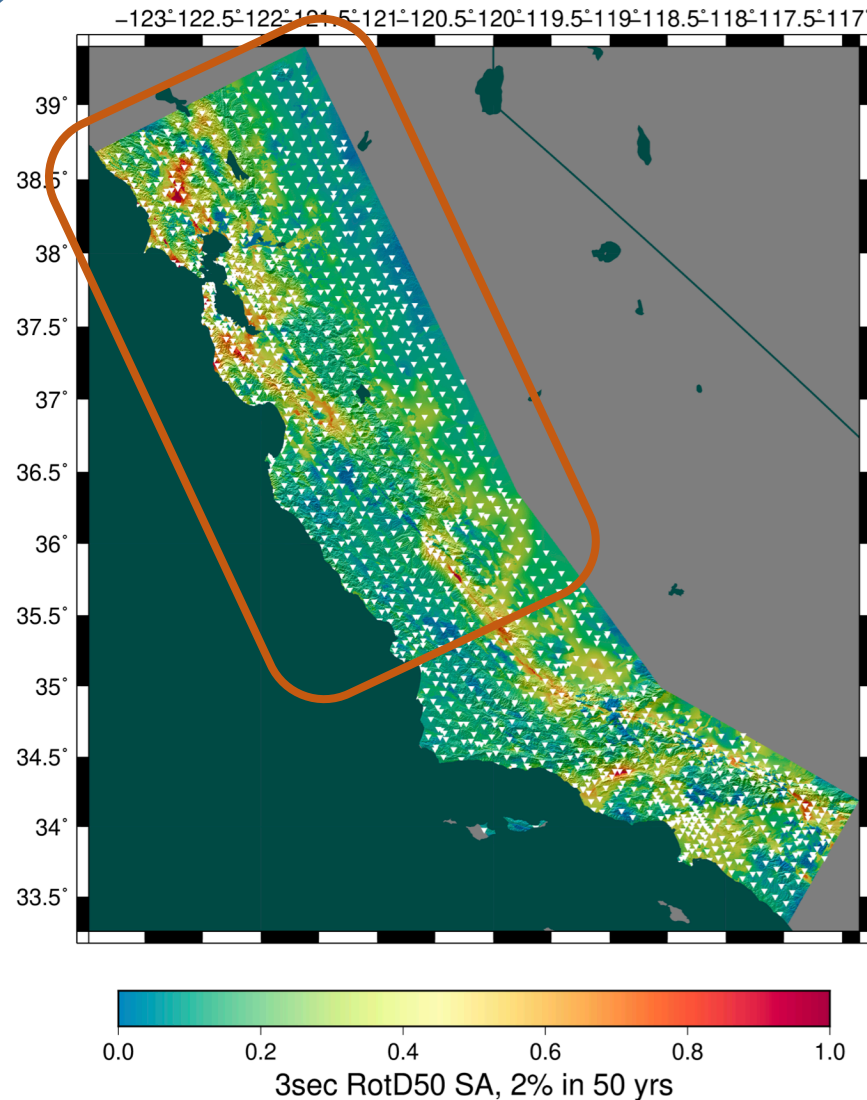
- Building engineers
- Disaster planners
- Insurance agencies

Slide credit: Southern California  
Earthquake Center



# Supporting Heterogeneous Workflows

## 2018-2019 Mapping Northern California



- 120 million core-hours
- 39,285 jobs
- 1.2 PB of data managed
- 157 TB of data automatically transferred
- 14.4 TB of output data archived

- NCSA *Blue Waters*
- OLCF *Titan*

Total map:  
170 million core hours  
> 19,407 core years

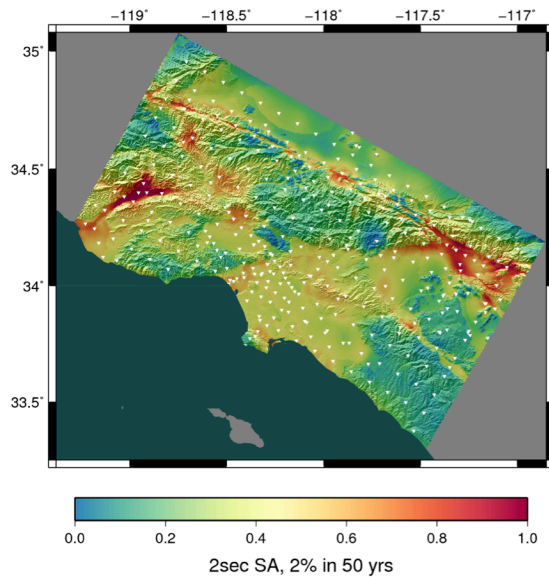
# Mix Workloads on Heterogeneous/ Changing CI

Since 2007: 215 million core-hours (24,543 years)

9 different supercomputers

## Pegasus Optimizations:

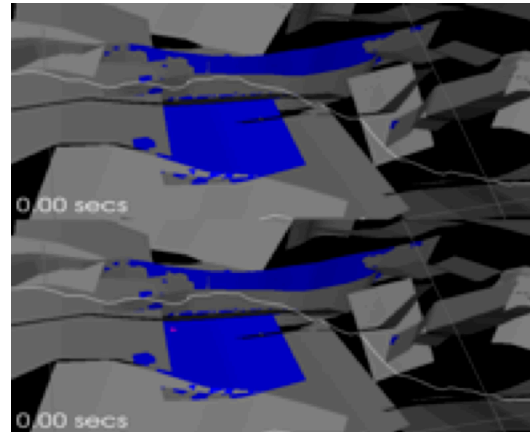
- Task clustering
- MPI-based workflow engine



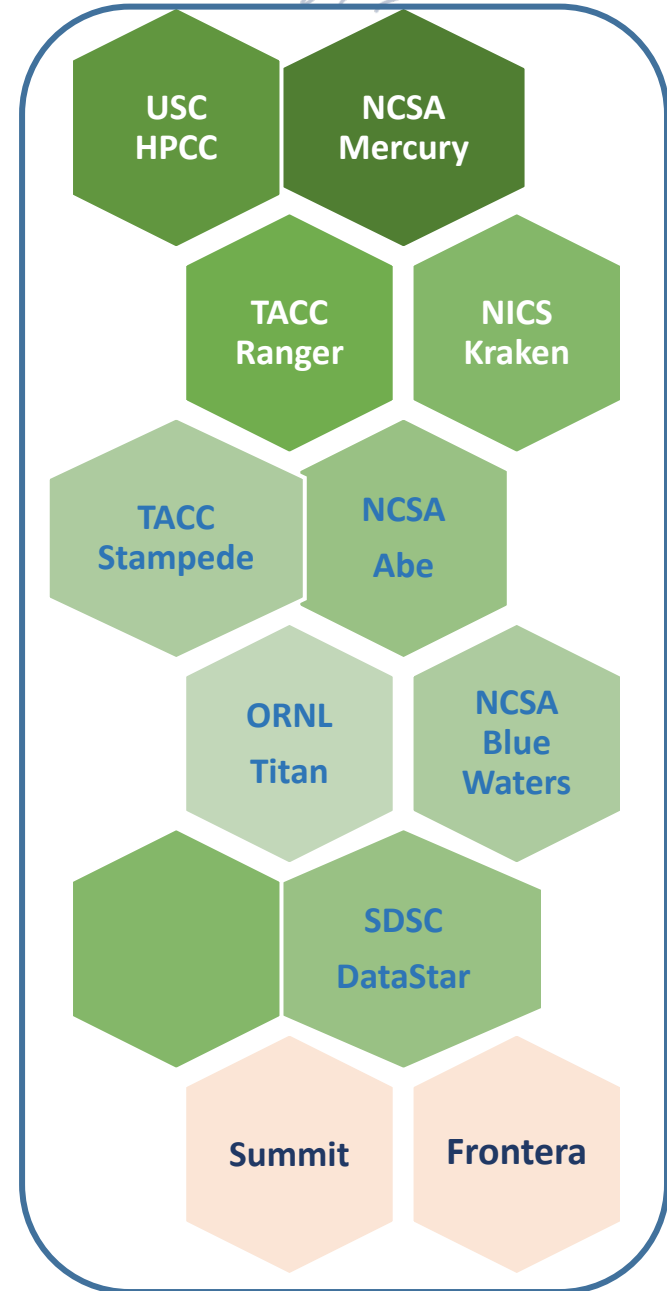
2010: World's first physics-based probabilistic seismic hazard map

## Application Optimizations:

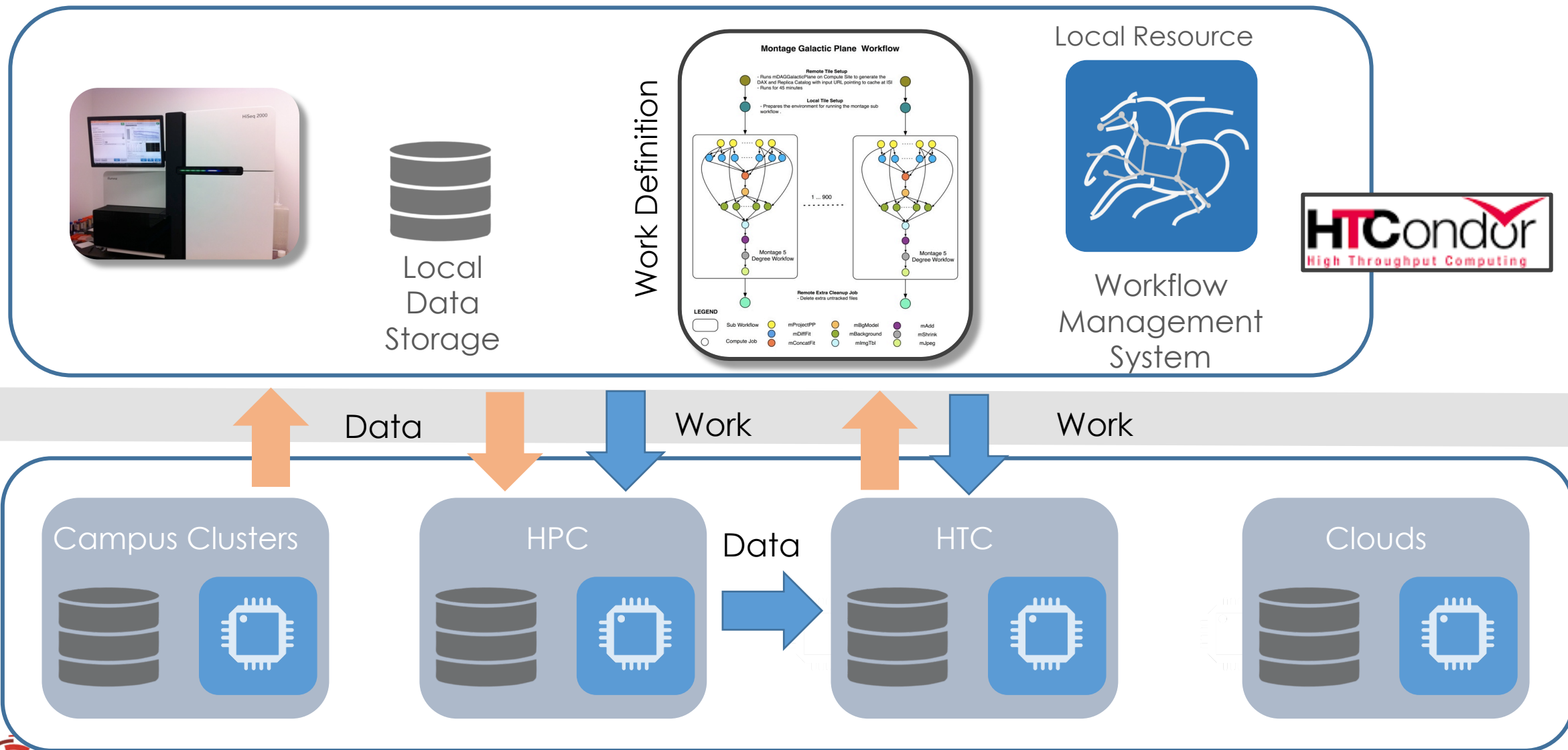
- Workflow restructuring
- MPI/code tuning
- Porting to GPUs



2018: Incorporating earthquake simulator with a 1 million-year catalog of California seismicity



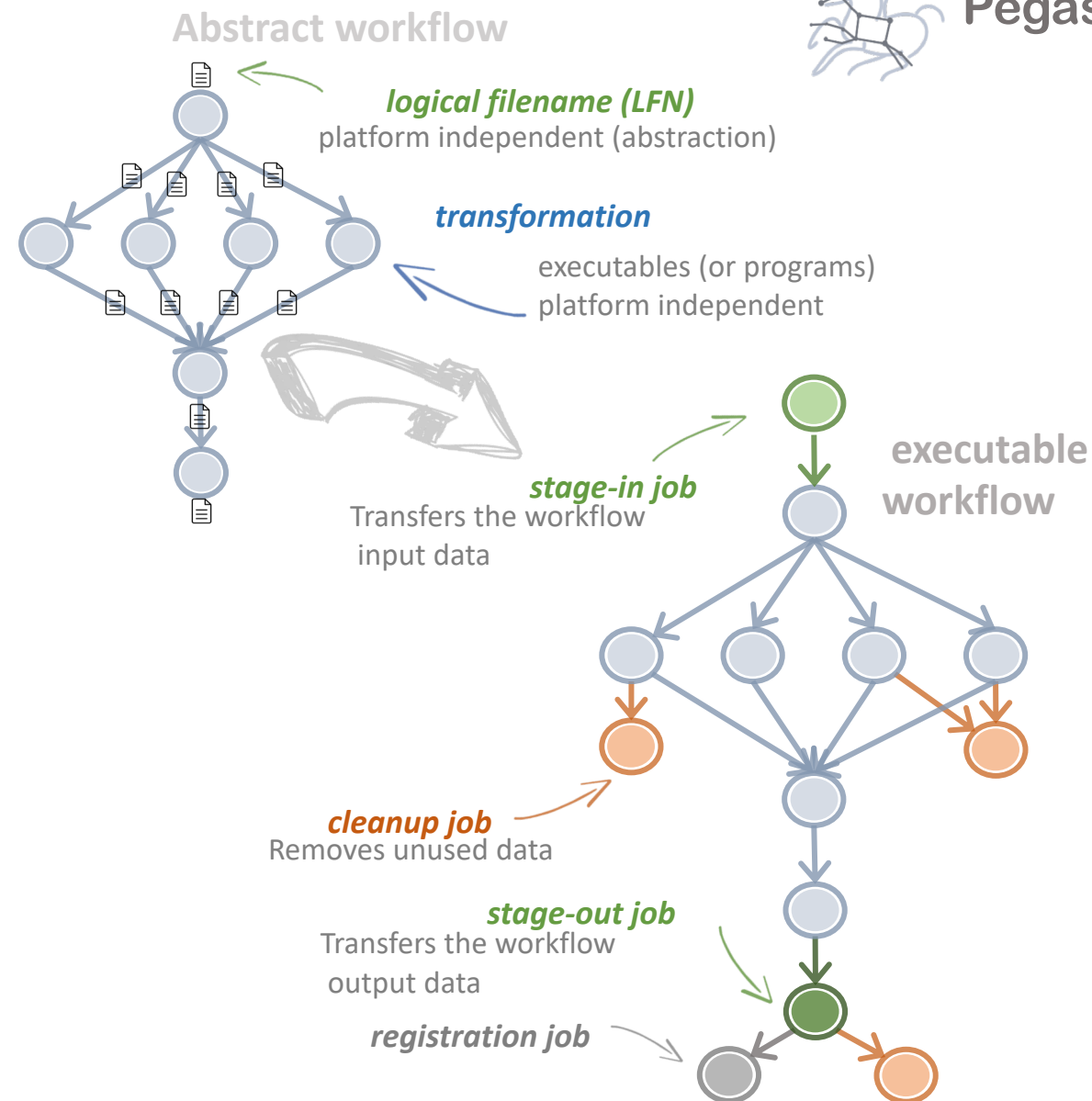
# Submit locally run globally



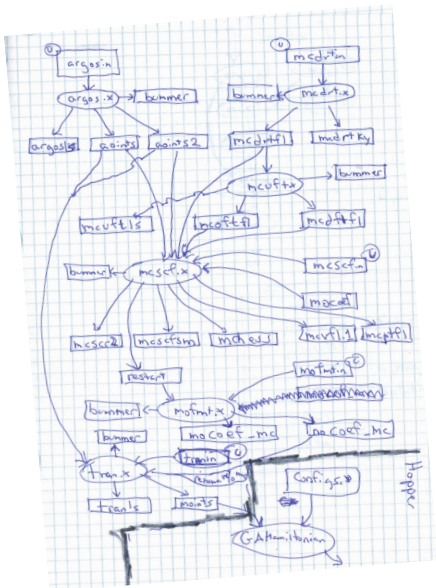
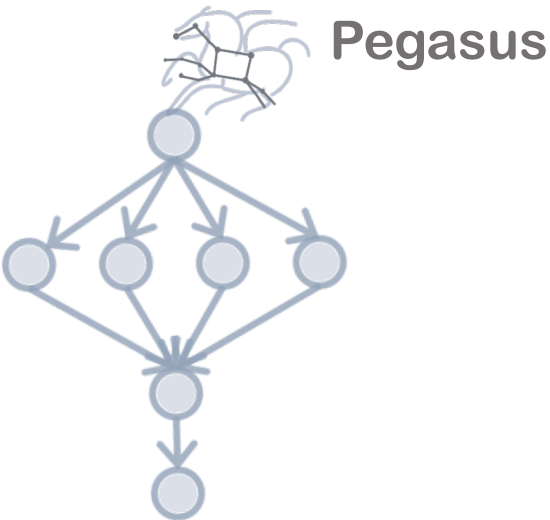


# Pegasus Concepts

- Users describe their pipelines in a **portable** format called Abstract Workflow, **without worrying** about **low level execution** details.
- Workflows are DAGs
  - Nodes: jobs, edges: dependencies
  - No while loops, no conditional branches
  - Jobs are standalone executables
- Pegasus takes this and **generates an executable workflow** that
  - has **data management** tasks added
  - **transforms** the workflow for **performance** and **reliability**



# Pegasus also provides tools to generate the workflow descriptions



```
#!/usr/bin/env python
from Pegasus.DAX3 import *
import sys
import os

# Create a abstract dag
dax = ADAG("hello_world")

# Add the hello job
hello = Job(namespace="hello_world",
            name="hello", version="1.0")
b = File("f.b")
hello.uses(a, link=Link.INPUT)
hello.uses(b, link=Link.OUTPUT)
dax.addJob(hello)

# Add the world job (depends on the hello job)
world = Job(namespace="hello_world",
            name="world", version="1.0")
c = File("f.c")
world.uses(b, link=Link.INPUT)
world.uses(c, link=Link.OUTPUT)
dax.addJob(world)

# Add control-flow dependencies
dax.addDependency(Dependency(parent=hello,
                              child=world))

# Write the DAX to stdout
dax.writeXML(sys.stdout)
```



```
<?xml version="1.0" encoding="UTF-8"?>
<!-- generator: python -->
<adag xmlns="http://pegasus.isi.edu/schema/DAX"
      version="3.4" name="hello_world">

  <!-- describe the jobs making
  up the hello world pipeline -->
  <job id="ID0000001" namespace="hello_world"
       name="hello" version="1.0">

    <uses name="f.b" link="output"/>
    <uses name="f.a" link="input"/>
  </job>

  <job id="ID0000002" namespace="hello_world"
       name="world" version="1.0">

    <uses name="f.b" link="input"/>
    <uses name="f.c" link="output"/>
  </job>

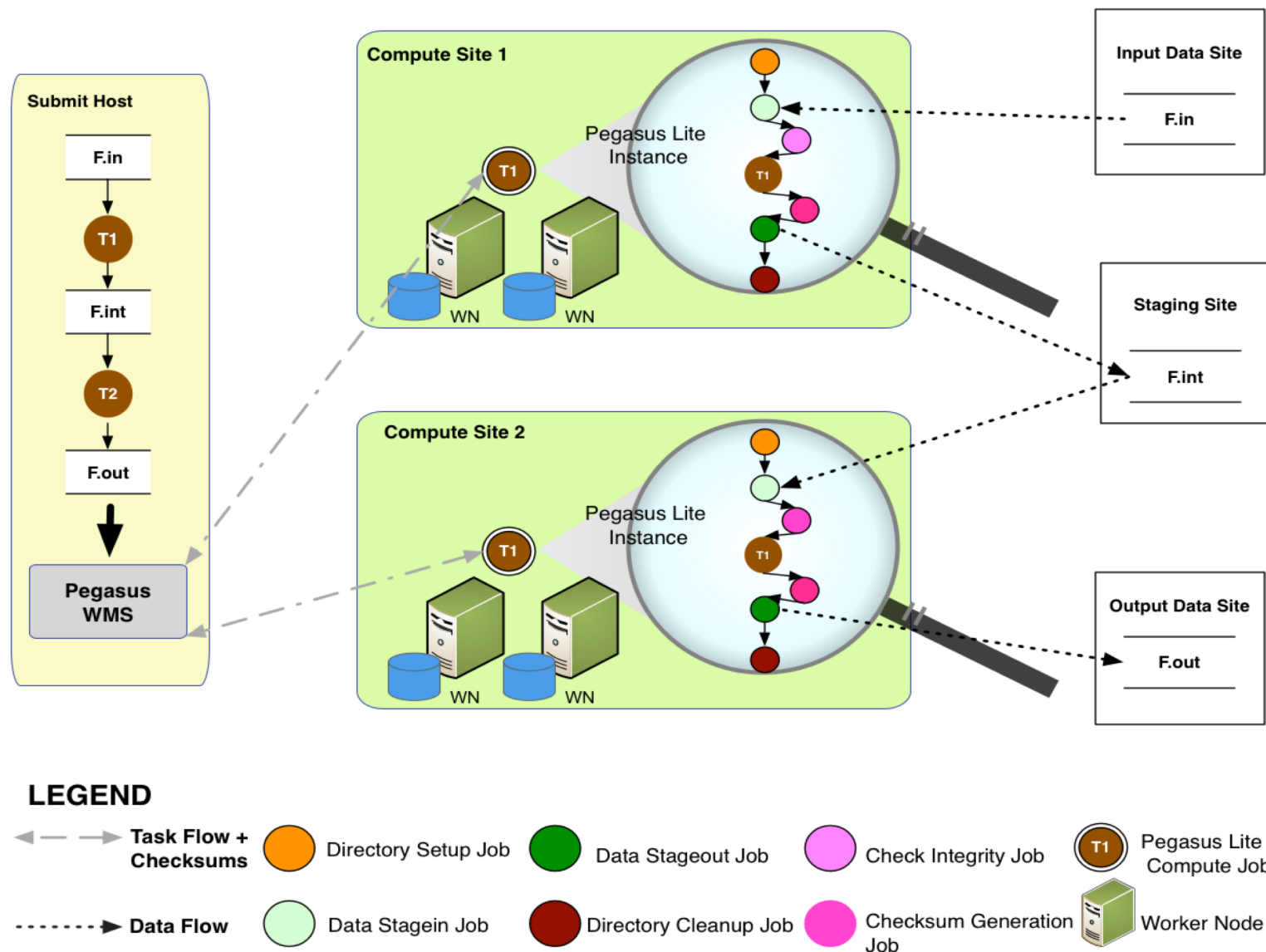
  <!-- describe the edges in the DAG -->
  <child ref="ID0000002">
    <parent ref="ID0000001"/>
  </child>
</adag>
```

DAX



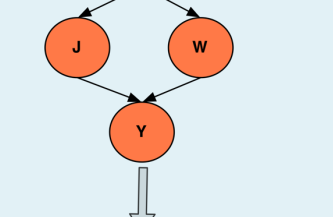
# Pegasus Deployment

- **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

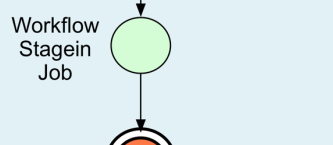
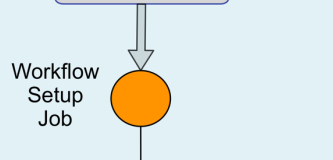


# Data Flow for LIGO Pegasus Workflows in OSG

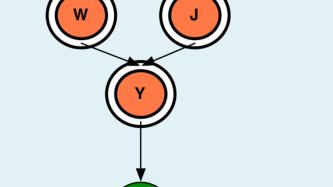
**SUBMIT HOST** Abstract Workflow



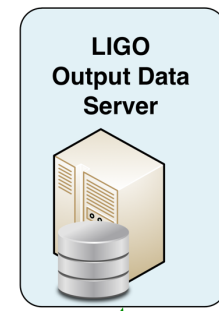
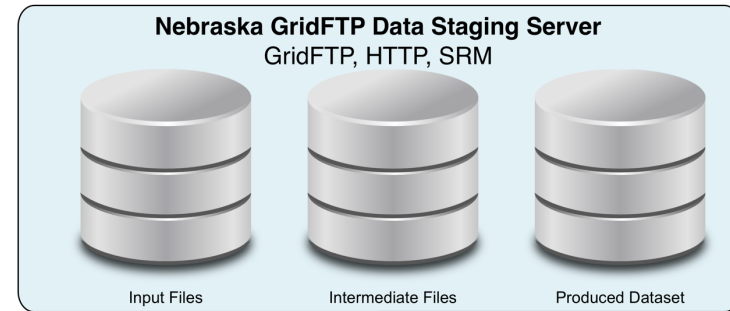
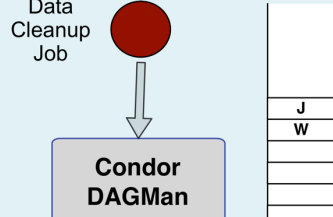
**Pegasus Planner**



**Condor Schedd Queue**



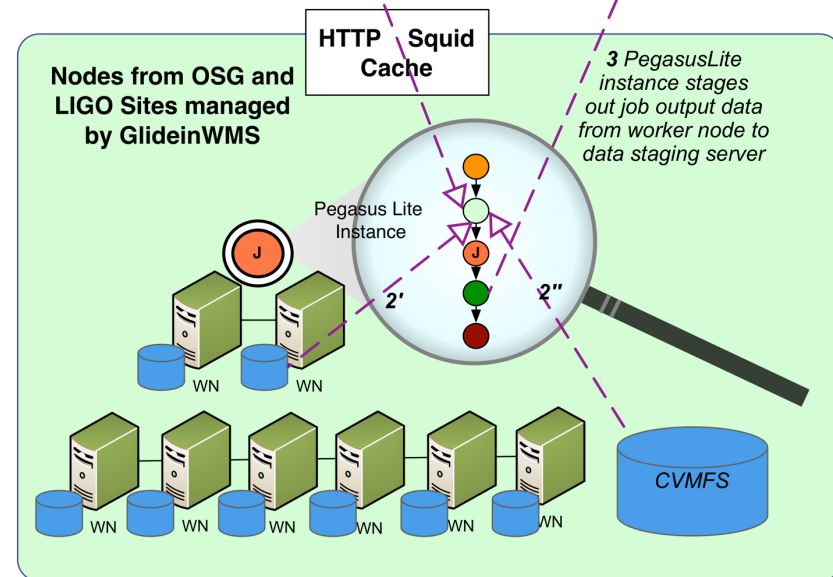
**Condor DAGMan**



1 Workflow Stagein Job stages in the input data for workflow from user server

2 PegasusLite instance looks up input data on the compute node/ CVMFS If not present, stage-in data from remote data staging server

4 Workflow Stageout Job stages produced data from data staging server to LIGO Output Data Server



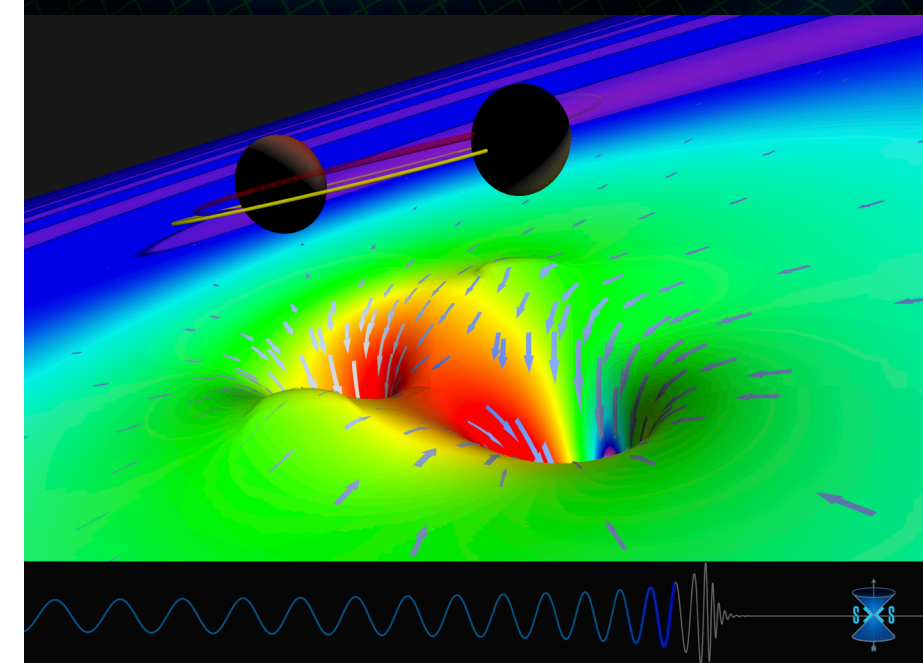
## LEGEND

- Orange circle: Directory Setup Job
- Green circle: Data Stageout Job
- Circle with J: Pegasus Lite Compute Job
- Light green circle: Data Stagein Job
- Red circle: Directory Cleanup Job
- Server rack icon: Worker Node

# Advanced LIGO – Laser Interferometer Gravitational Wave Observatory

60,000 compute tasks  
Input Data: 5,000 files (10GB total)  
Output Data: 60,000 files (60GB total)

Executed on LIGO Data Grid, Open Science Grid and XSEDE



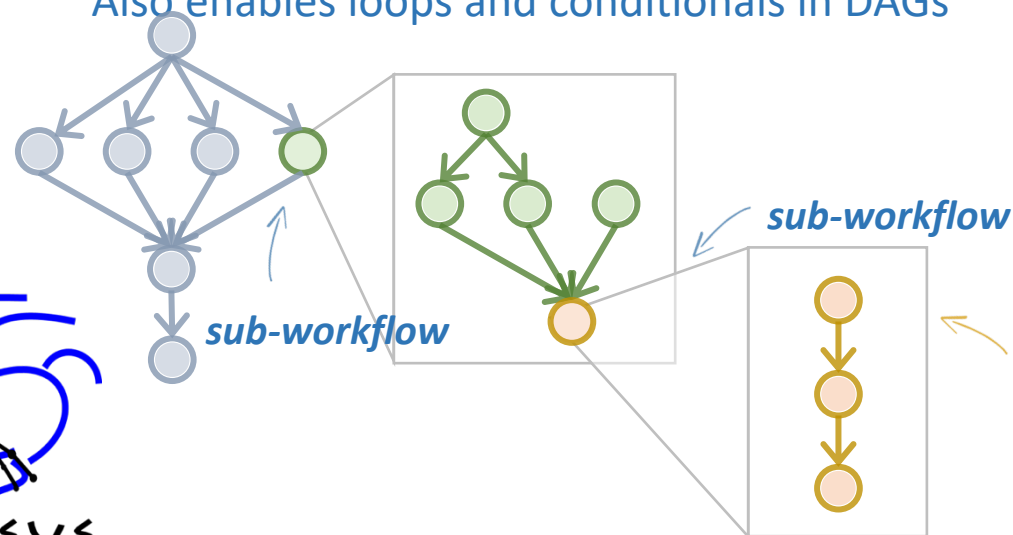
# Pegasus Optimizations



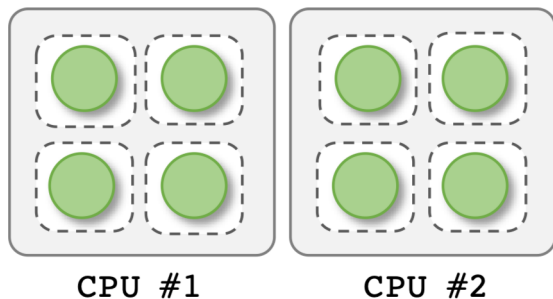
Pegasus

## Hierarchical workflows

Enacts the execution of **millions of tasks**  
Also enables loops and conditionals in DAGs

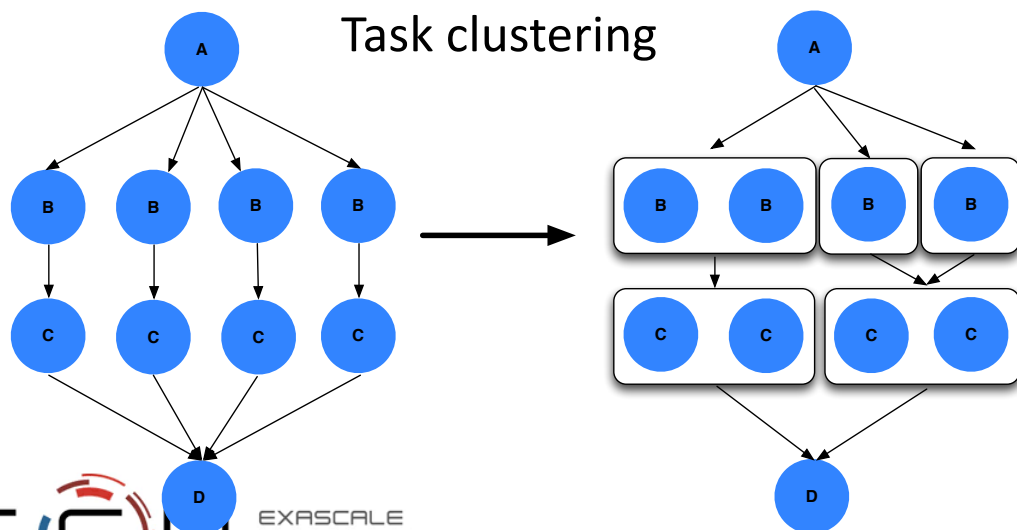


## Task-resource co-allocation

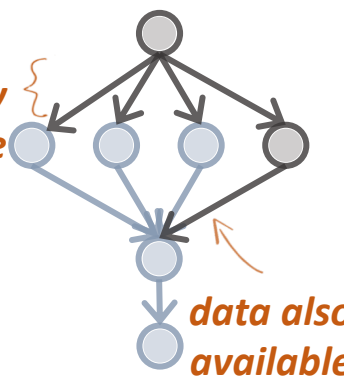


pegasus

## Task clustering



data already available

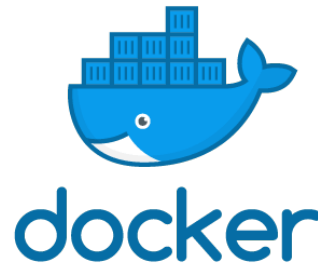


workf  
low  
reduc  
tion

modern workflow optimizations

# Pegasus Container Support

- Support for
  - Docker
  - Singularity
  - Shifter (coming soon)



- Users can refer to **containers** in the **Transformation Catalog** with their executable preinstalled.
- Users can **refer** to a **container** they want to **use**. However, they let **Pegasus stage** their executable 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 at tar file and ship them around ( planned for 4.8.X )*

# Pegasus: Containers Data Management

- Treat containers as input data dependency
  - Needs to be staged to compute node if not present
- Users can refer to container images as
  - Docker Hub or Singularity Library URL's
  - Docker Image exported as a TAR file and available at a server , just like any other input dataset.
- If an image is specified to be residing in a hub
  - 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
  - Motivation: Avoid hitting Docker Hub/Singularity Library repeatedly for large workflows
- Symlink against a container image if available on shared filesystem
  - For e.g. CVMFS hosted images on Open Science Grid



# 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.

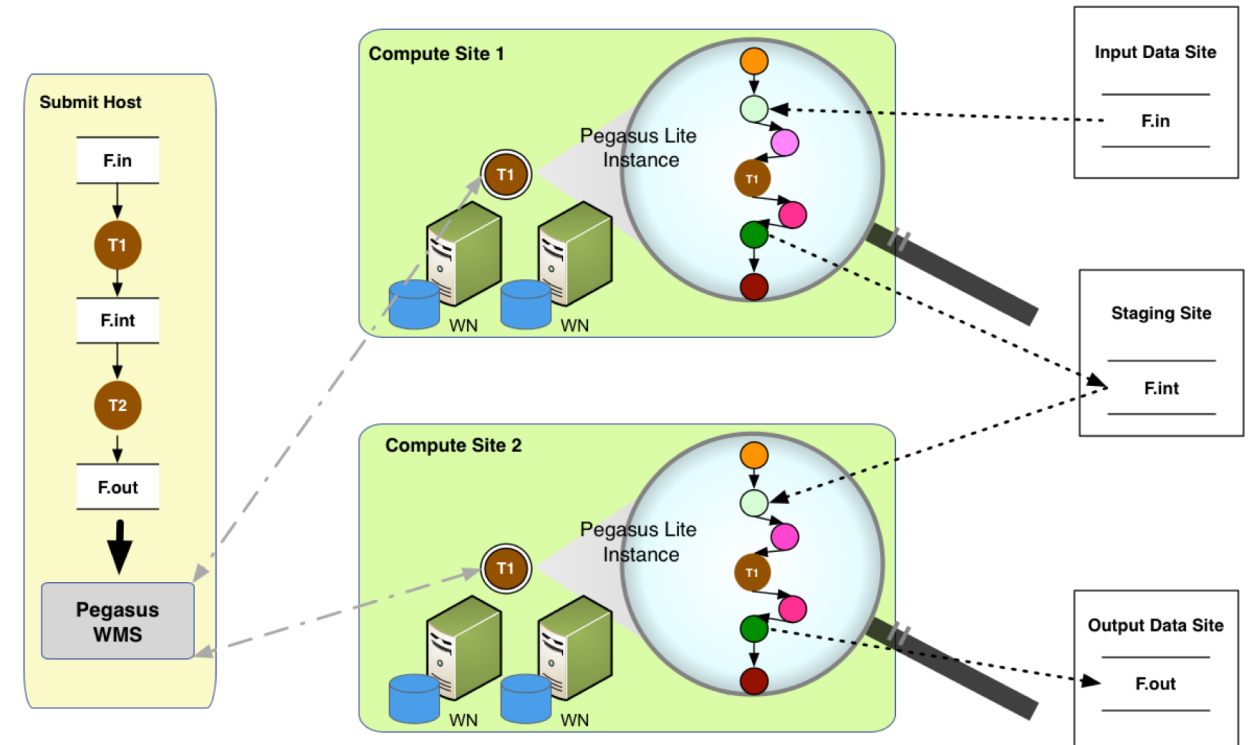
Plus there is the threat of intentional changes: malicious attackers, insider threats, etc.

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?”

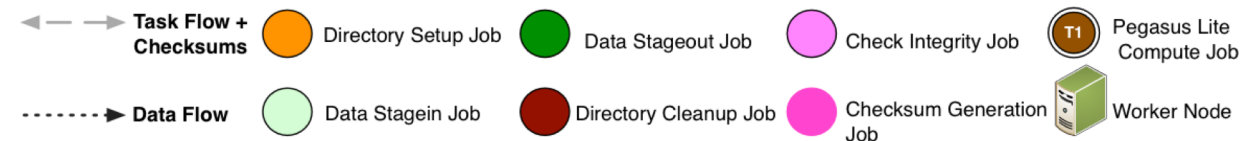
# Automatic Integrity Checking in Pegasus

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



## LEGEND



**Job failure** is triggered if checksums fail

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

Worker nodes do not share a file system

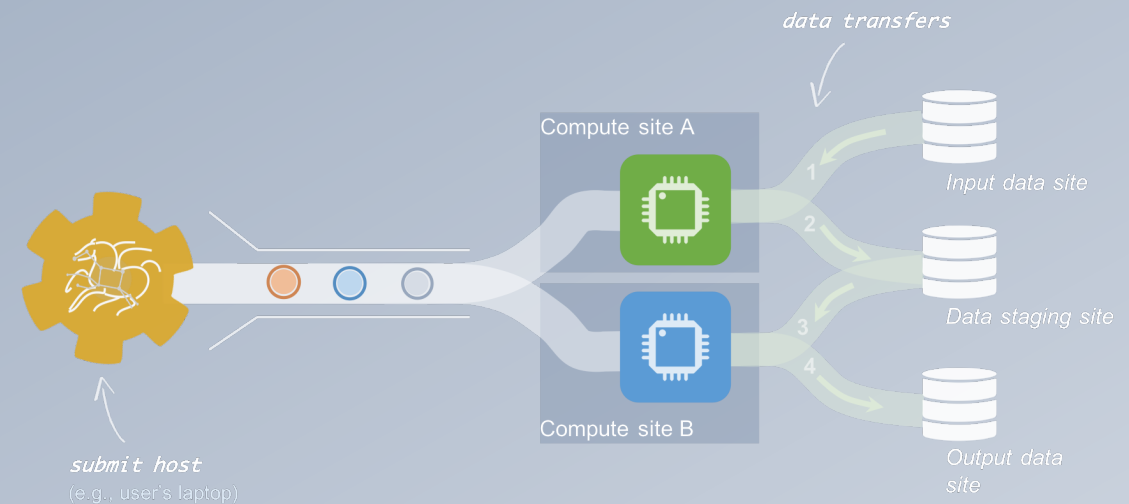
Data is pulled from / pushed to the submit host via HTCondor file transfers

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

# pegasus-transfer

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



Pegasus

## 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

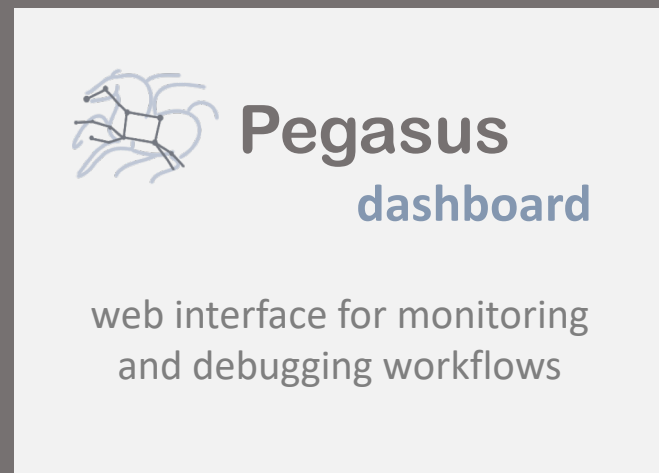
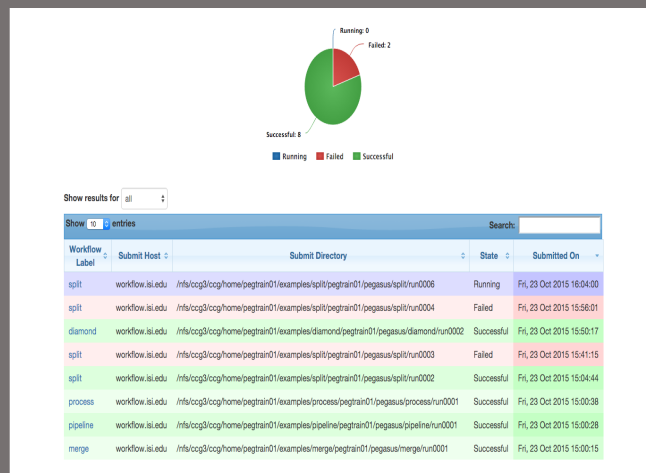
## Parallel transfers

## Automatic retries

## Credential management

Uses the appropriate credential for each site and each protocol (even 3<sup>rd</sup> party transfers)

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



### 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

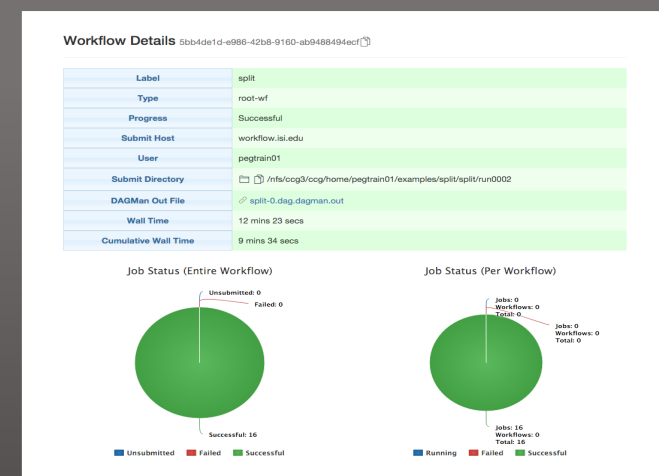
### Entire Workflow

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 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.



Real-time Monitoring  
Reporting  
Debugging  
Troubleshooting  
RESTful API





# Pegasus dashboard

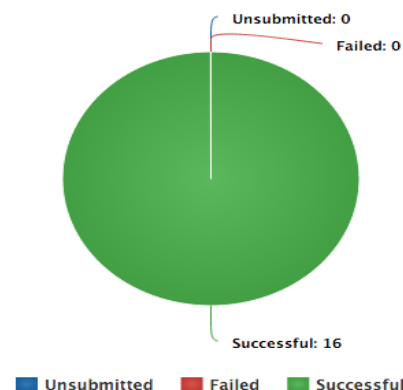
web interface for monitoring  
and debugging workflows

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.

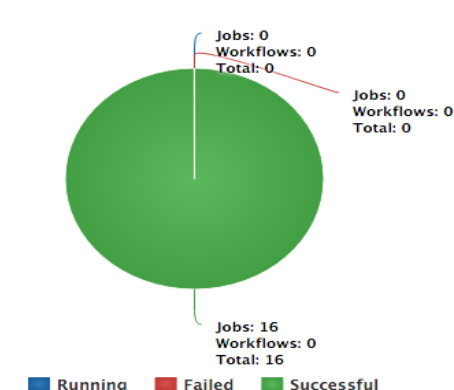
## Workflow Details 5bb4de1d-e986-42b8-9160-ab9488494ecf

Label	split
Type	root-wf
Progress	Successful
Submit Host	workflow.isi.edu
User	pegtrain01
Submit Directory	/nfs/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

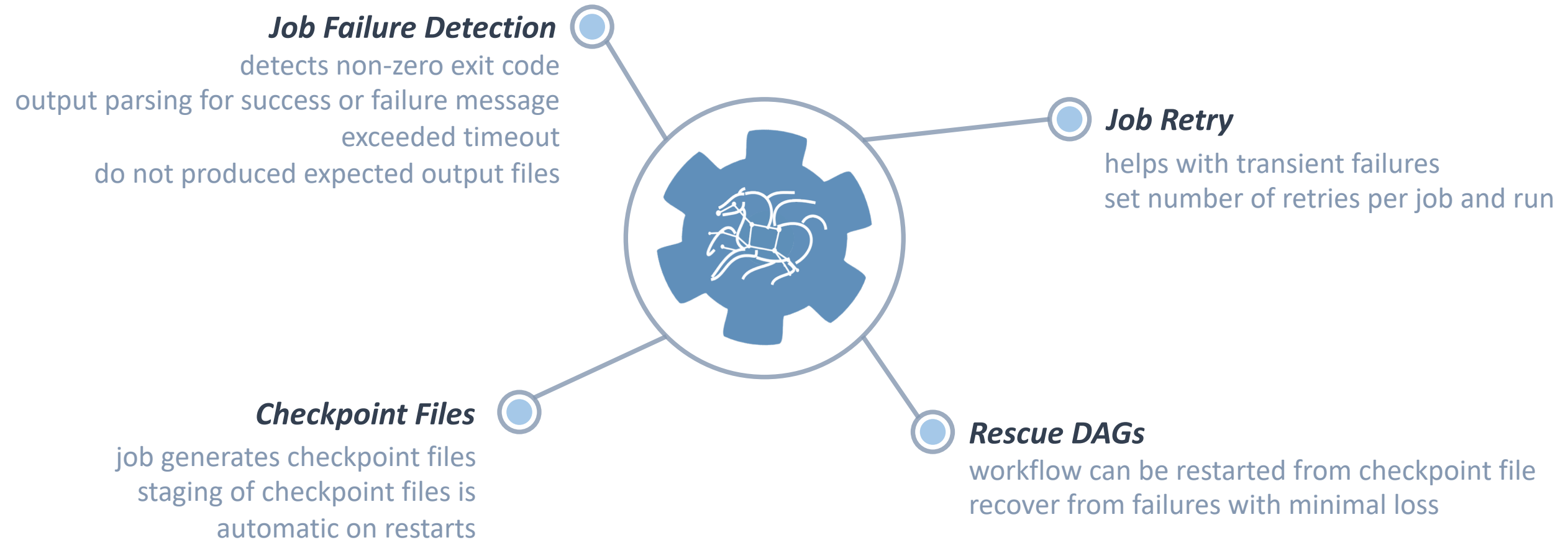
Job Status (Entire Workflow)



Job Status (Per Workflow)



# And if a job fails?



# 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*

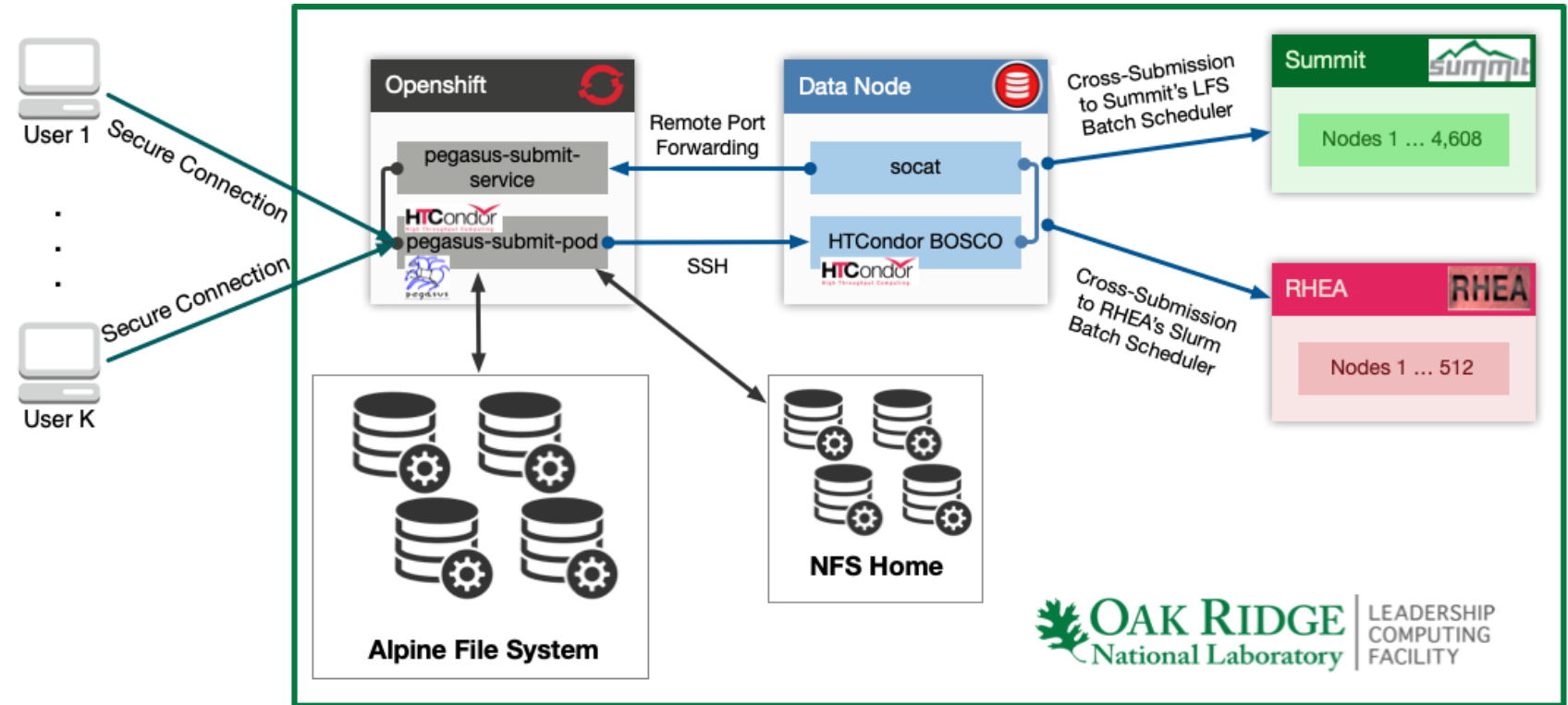
### **Globus GRAM**

*Uses CondorG*

# Pegasus at OLCF: Kubernetes Deployment



- Pegasus workflow **environments** at OLCF have been **simplified**.
- Using the Kubernetes cluster at OLCF, we can deploy Pegasus submit nodes as services.
- This solution uses HTCondor's BOSCO SSH style submissions on the DTNs and achieves submissions to the SLURM and LSF batch schedulers.



- This approach is powerful because a single workflow can be configured to use **all** of OLCF's resources. Execute transfers on the DTNs, run simulations and heavy processing on Summit and then do lightweight post processing steps on RHEA.

GitHub: <https://github.com/pegasus-isi/pegasus-olcf-kubernetes>

# Questions?





# 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)

### Pegasus Online Office Hours

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

*Bi-monthly basis on second Friday of the month, where we address user questions and also apprise the community of new developments*