



U.S. DEPARTMENT OF
ENERGY



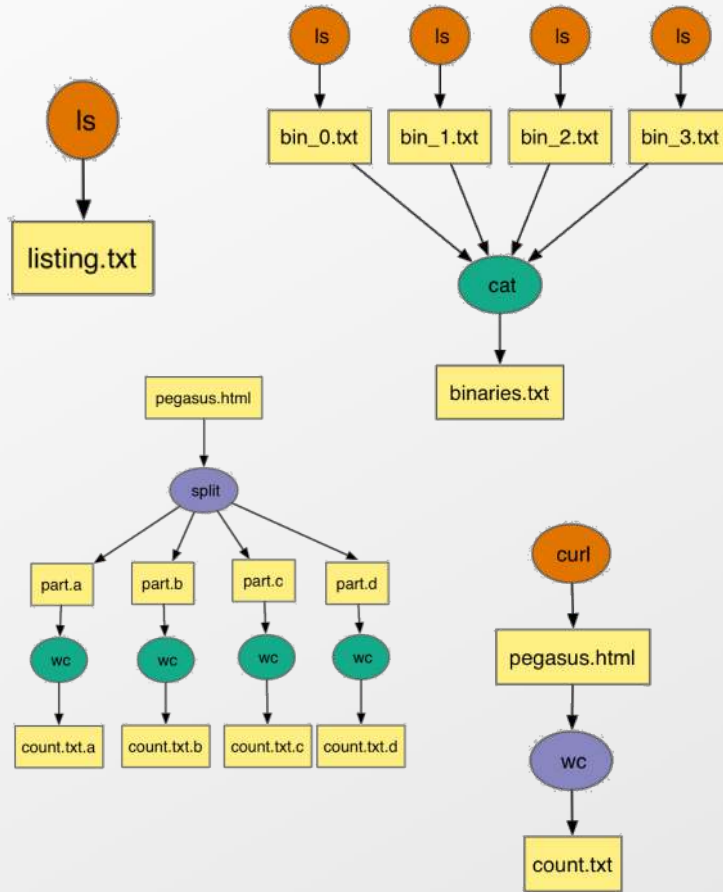
Workflows using Pegasus

Pegasus Workflow Management System

Karan Vahi



Compute Pipelines Building Blocks



Compute Pipelines

Allows scientists to connect different codes together and execute their analysis

Pipelines can be very simple (independent or parallel) jobs or complex represented as DAG's

Helps users to automate scale up

However, it is still up-to user to figure out

Data Management

How do you ship in the small/large amounts data required by your pipeline and protocols to use?

How best to leverage different infrastructure setups

OSG has no shared filesystem while XSEDE and your local campus cluster has one!

Debug and Monitor Computations

Correlate data across lots of log files

Need to know what host a job ran on and how it was invoked

Restructure Workflows for Improved Performance

Short running tasks? Data placement

Why Pegasus ?

Automates complex, multi-stage processing pipelines

Enables parallel, **distributed** computations

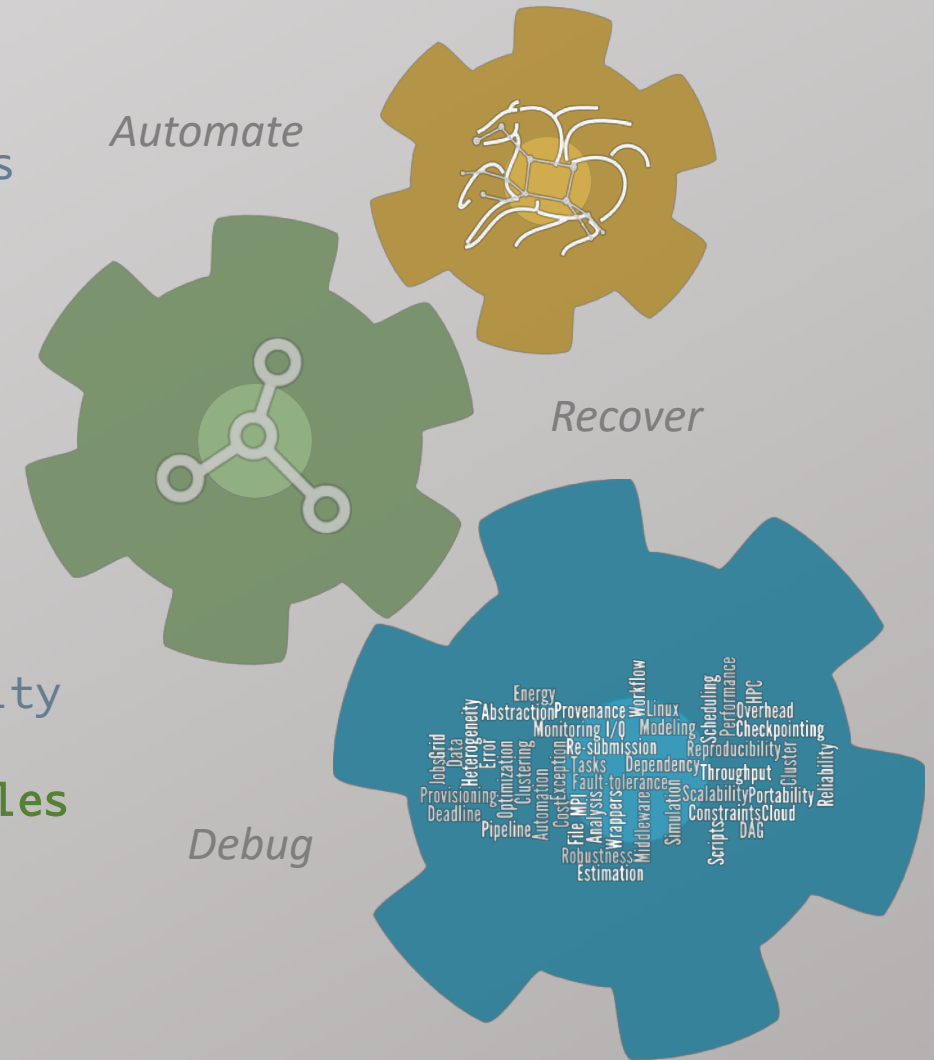
Automatically executes data transfers

Reusable, aids **reproducibility**

Records how data was produced (**provenance**)

Handles **failures** with to provide reliability

Keeps track of data and **files**



NSF funded project since 2001,
with close collaboration with
HTCondor team

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

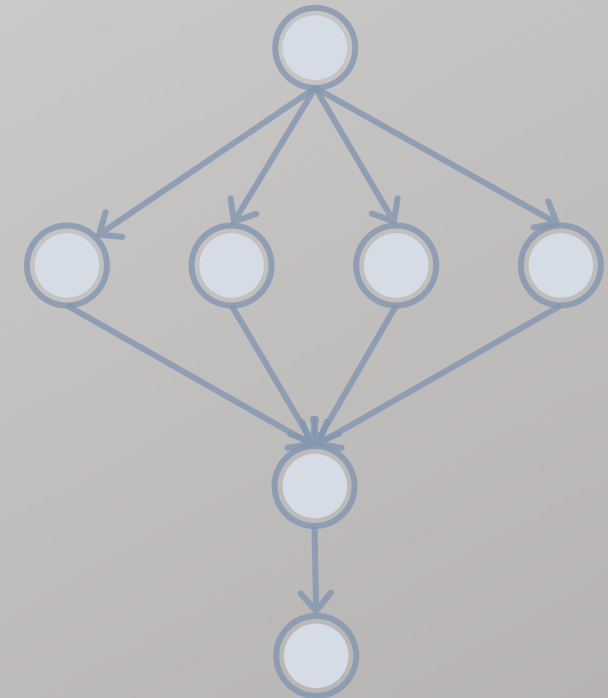
No while loops, no conditional branches

Jobs are standalone executables

Planning occurs ahead of execution

Planning converts an abstract workflow into a concrete, executable workflow

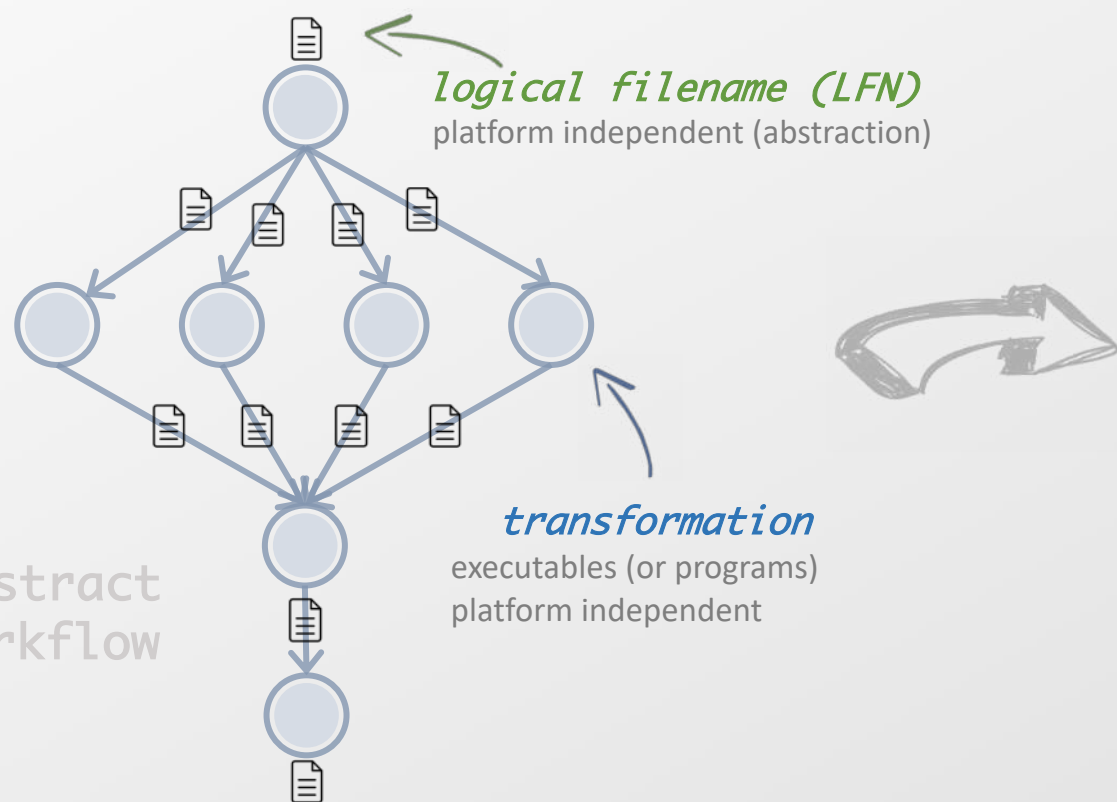
Planner is like a compiler



DAX

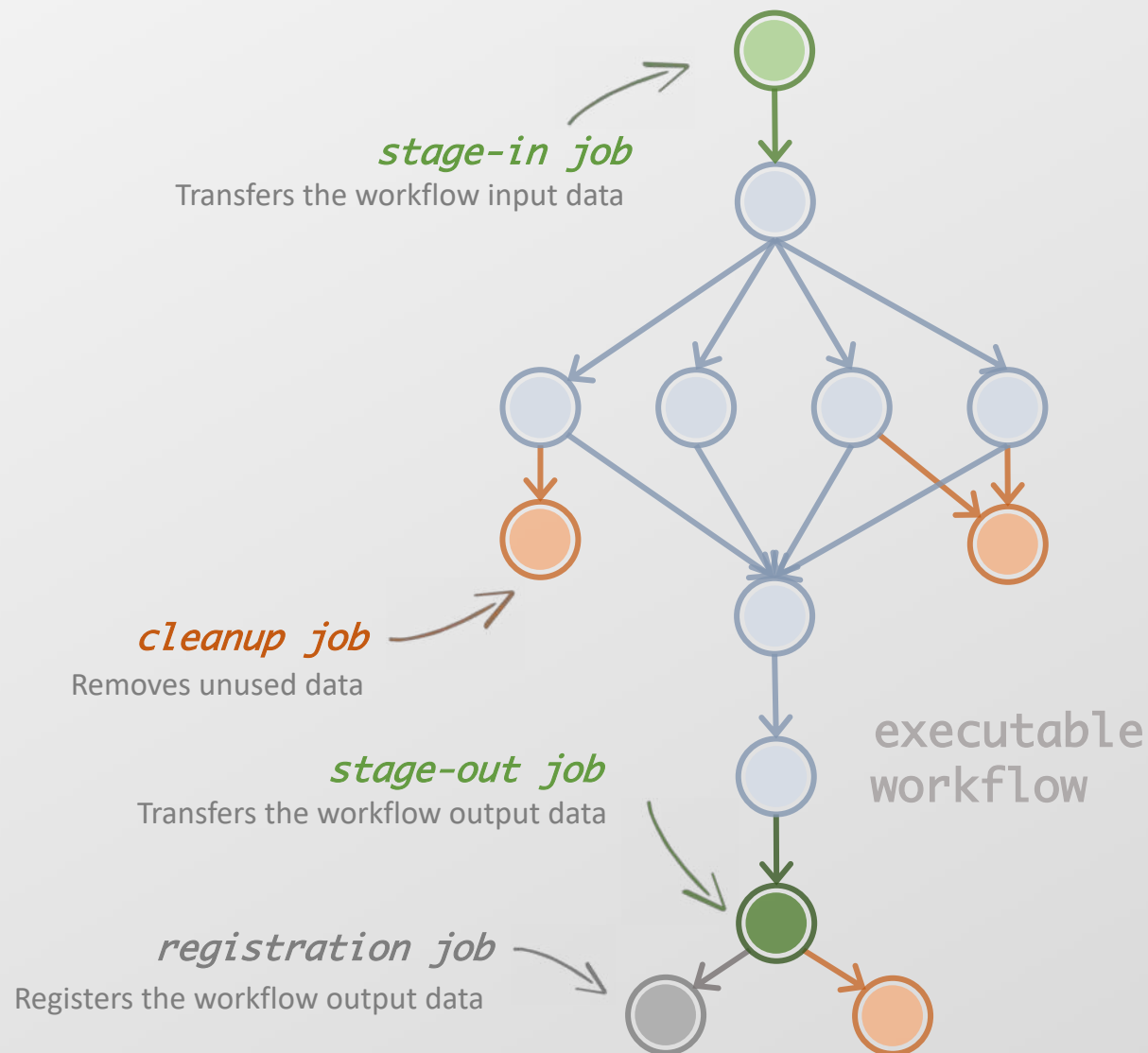
Portable Description

Users do not worry about
low level execution details



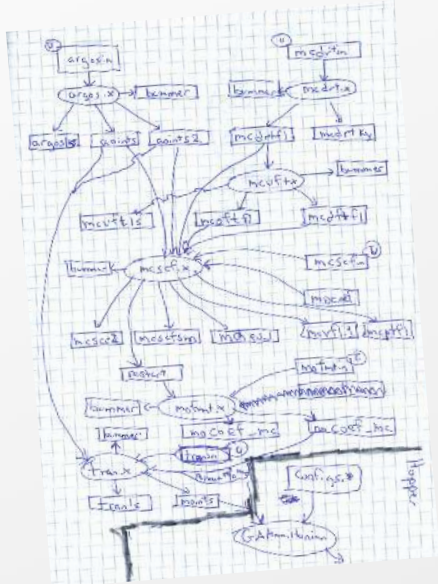
abstract
workflow

directed-acyclic graphs



executable
workflow

Pegasus also provides tools to generate the abstract workflow



```
#!/usr/bin/env python

from Pegasus.DAX3 import *
import sys
import os

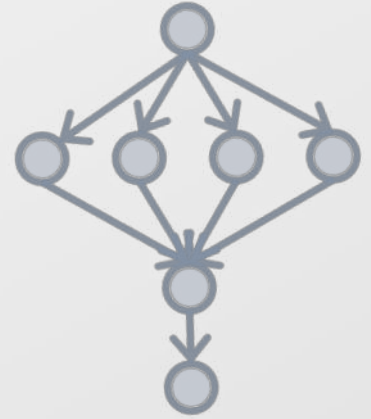
# Create an abstract dag
dax = ADAG("hello_world")

# 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">

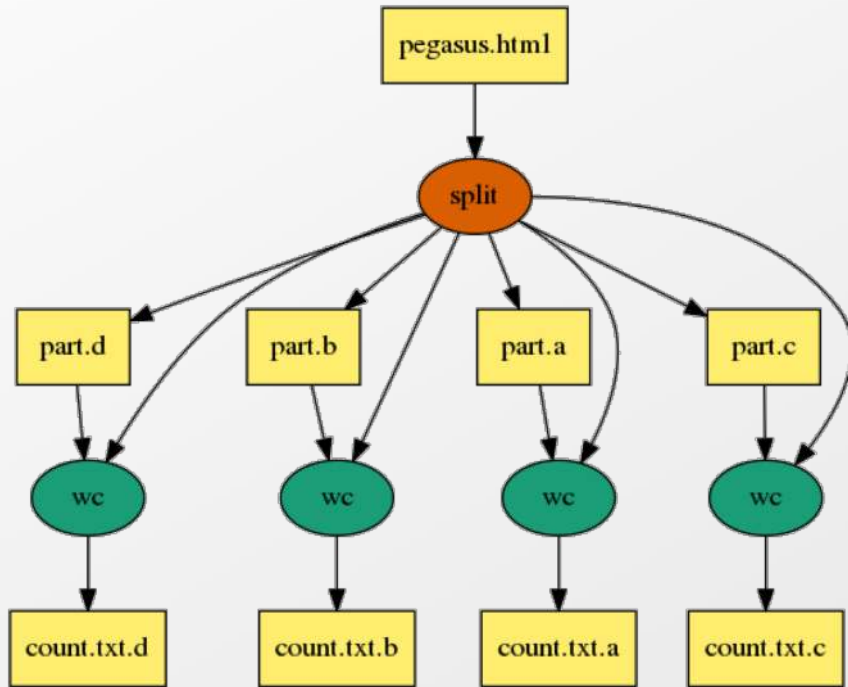
    <uses name="f.b" link="input"/>
    <uses name="f.c" link="output"/>
  </job>

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



DAG IN XML
DAX

An example Split Workflow



Visualization Tools:
pegasus-graphviz
pegasus-plots

https://pegasus.isi.edu/documentation/tutorial_submitting_wf.php

```
#!/usr/bin/env python
```

```
import os, pwd, sys, time
from Pegasus.DAX3 import *
```

```
# Create an abstract dag
dax = ADAG("split")
```

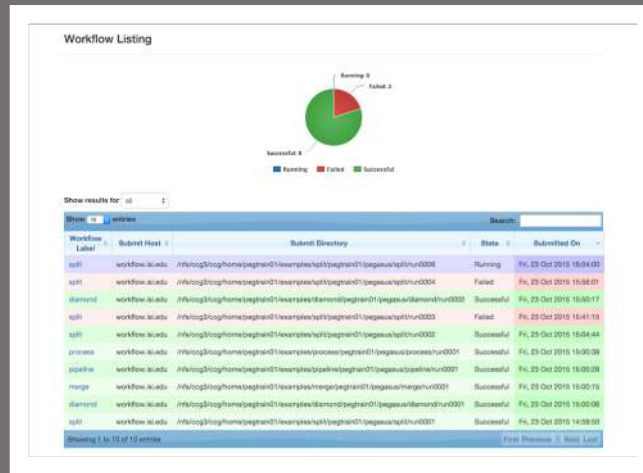
```
webpage = File("pegasus.html")
```

```
# the split job that splits the webpage into smaller chunks
split = Job("split")
split.addArguments("-l", "100", "-a", "1", webpage, "part.")
split.uses(webpage, link=Link.INPUT)
# associate the label with the job. all jobs with same label
# are run with PMC when doing job clustering
split.addProfile( Profile("pegasus", "label", "p1"))
dax.addJob(split)
```

```
# we do a parameter sweep on the first 4 chunks created
for c in "abcd":
    part = File("part.%s" % c)
    split.uses(part, link=Link.OUTPUT, transfer=False, register=False)
    count = File("count.txt.%s" % c)
    wc = Job("wc")
    wc.addProfile( Profile("pegasus", "label", "p1"))
    wc.addArguments("-l", part)
    wc.setStdout(count)
    wc.uses(part, link=Link.INPUT)
    wc.uses(count, link=Link.OUTPUT, transfer=True, register=True)
    dax.addJob(wc)
```

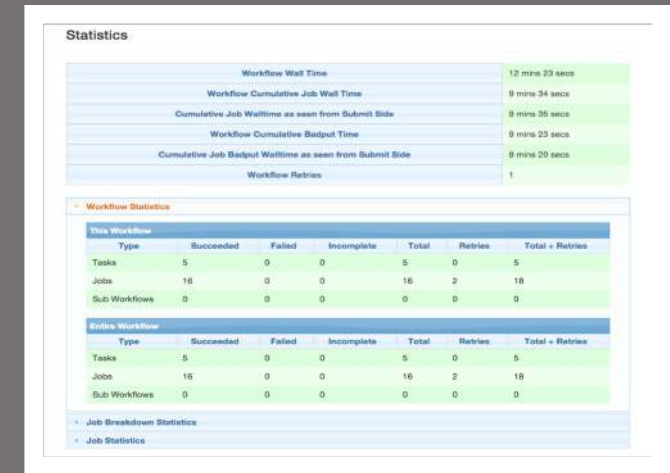
```
#adding dependency
dax.depends(wc, split)
```

```
f = open("split.dax", "w")
dax.writeXML(f)
f.close()
```

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.



Real-time Monitoring
Reporting
Debugging
Troubleshooting
RESTful API



Pegasus dashboard

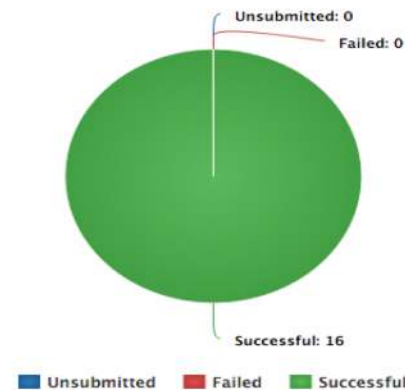
web interface for monitoring
and debugging workflows

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.

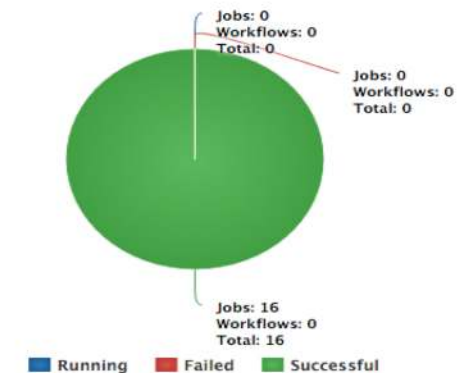
Workflow Details 5bb4de1d-e986-42b8-9160-ab9488494ecf

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)





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

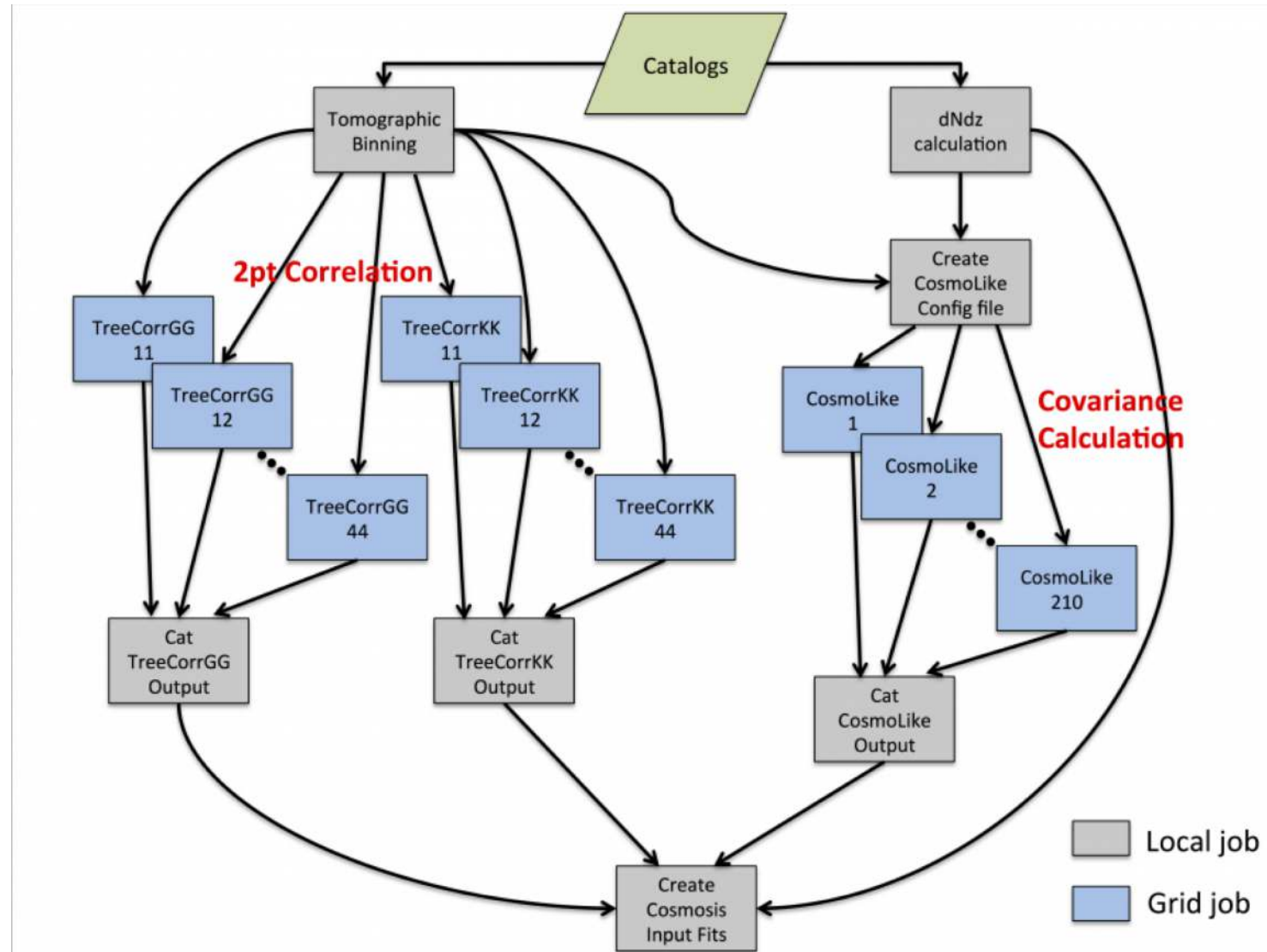
pegasus-analyzer



Hands-on Demo...

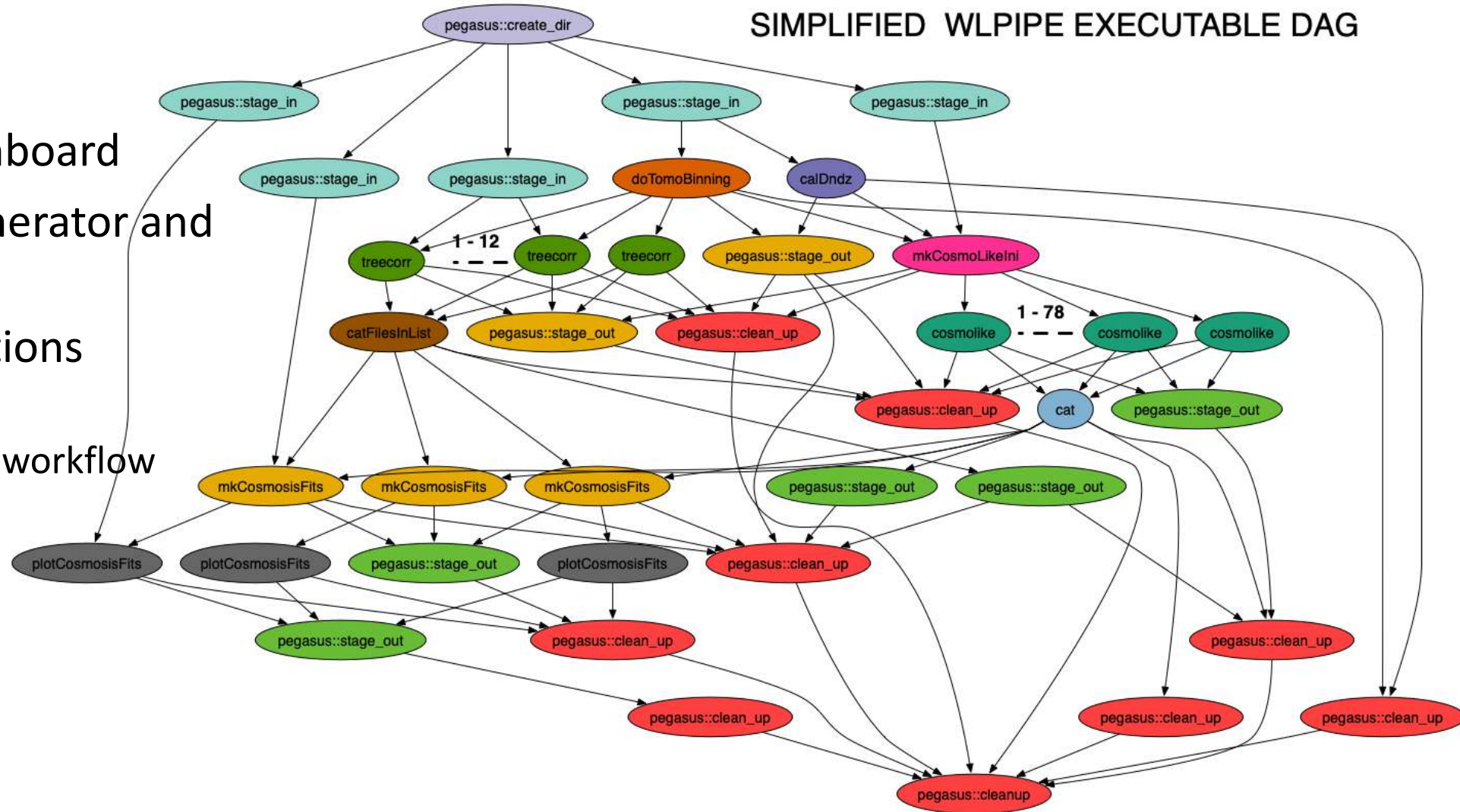
Hands On Demo

- Weak Lensing Pipeline
 - <https://github.com/pegasus-isi/pegasus-wlpipe>
- An example of a typical gravitational weak lensing analysis. It uses publicly available Science Verification catalogs the Dark Energy Survey (DES-SV).
- The pipeline is run currently at Fermi Grid
- We will run the example version at a cluster at ISI
- Science Codes are bundled into a **Singularity Container**
- **This is an interactive session. Please interrupt at anytime to ask questions.**

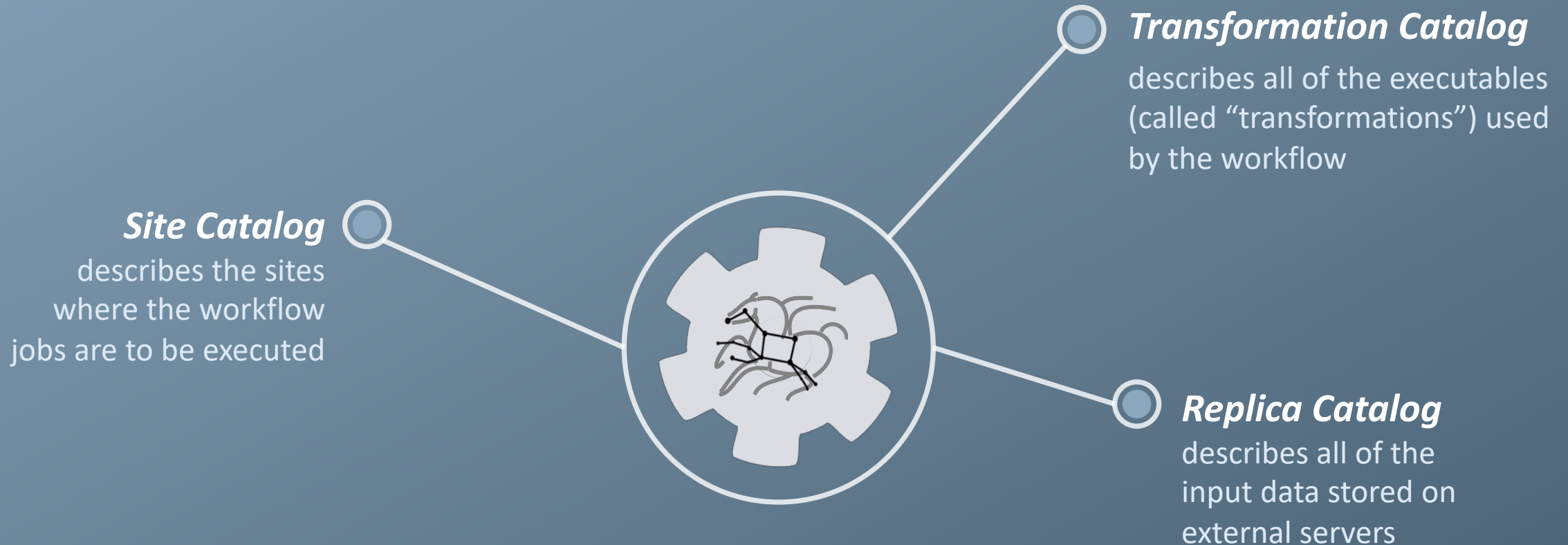


Outline

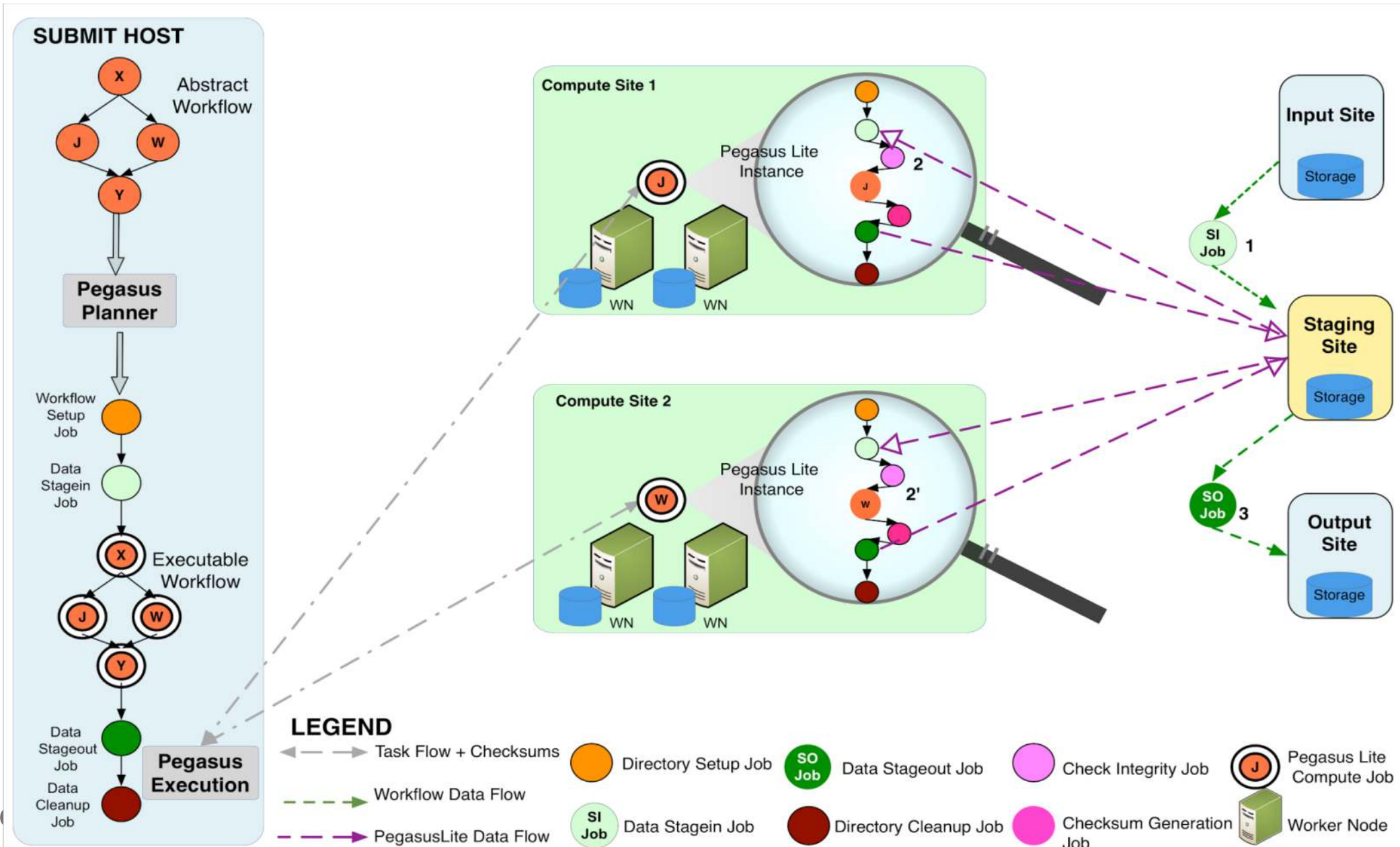
- Submit a workflow
- Go through the dashboard
- Go over the DAX generator and catalogs
- Show debugging options
 - pegasus-analyzer
 - Recover from failed workflow
- pegasus-statistics



So, what information does Pegasus need?



Distributed Execution



Job Submissions

Local

Submit Machine

Personal HTCondor

*Local Campus Cluster accessible via
Submit Machine **

HTCondor via Glite

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

Supported schedulers

PBS SGE SLURM 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

ICE Cube Glidein service

OSG using glideinWMS

CREAMCE

Uses CondorG

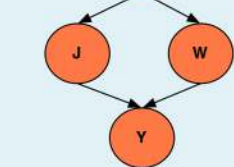
Globus GRAM

Uses CondorG

Some of the successful stories...

Data Flow for LIGO Pegasus Workflows in OSG

SUBMIT HOST Abstract Workflow



Pegasus Planner

Workflow Setup Job

Workflow Stagein Job

Executable Workflow

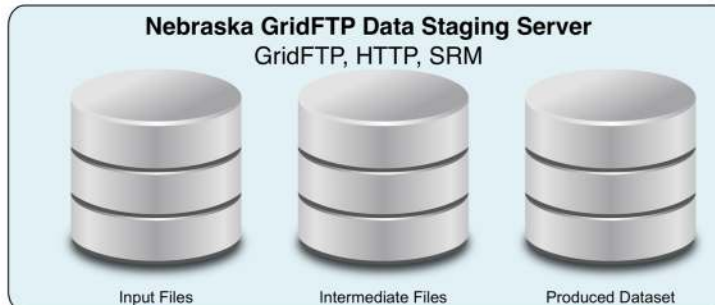
Workflow Stageout Job
Data Cleanup Job

Condor Schedd Queue

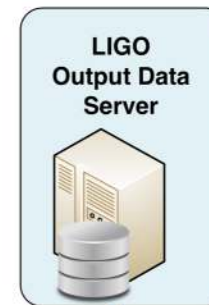
Condor DAGMan



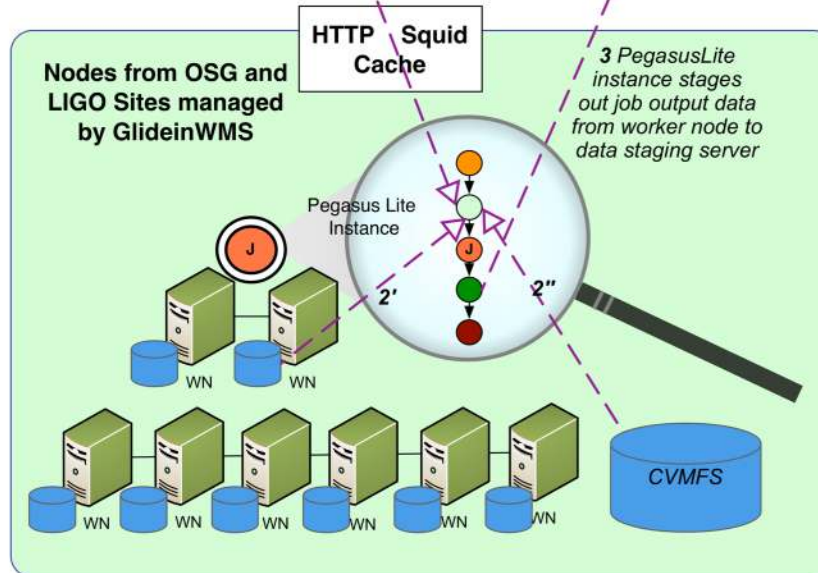
1 Workflow Stagein Job stages in the input data for workflow from user server



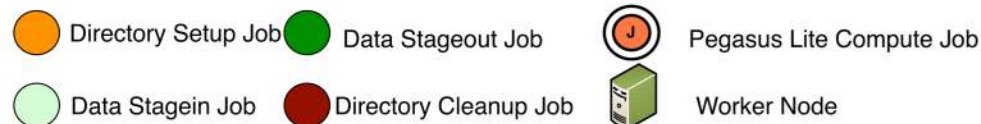
2 PegasusLite instance looks up input data on the compute node/ CVMFS
If not present, stage-in data from remote data staging server



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



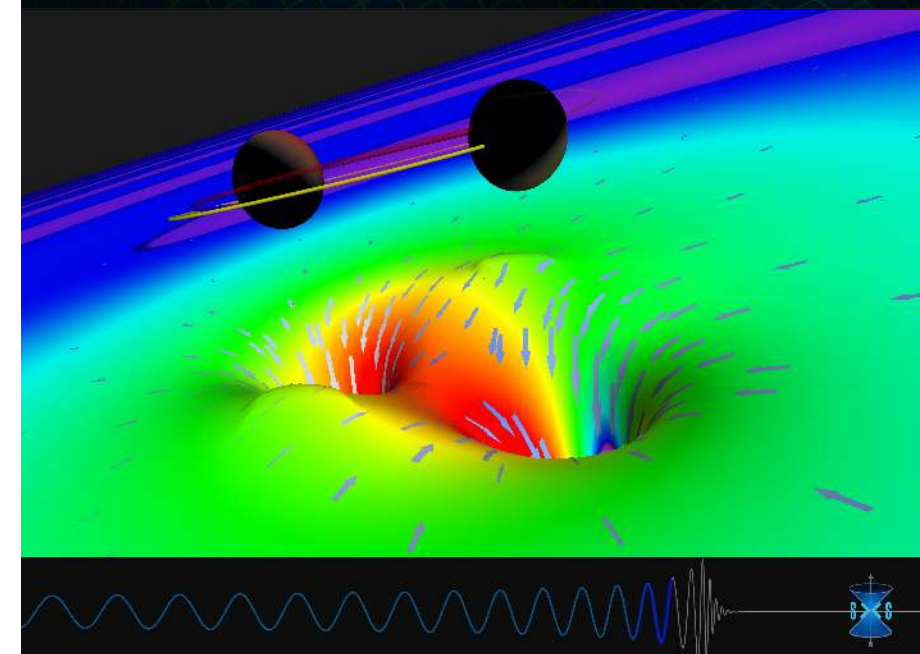
LEGEND



Advanced LIGO – Laser Interferometer Gravitational Wave Observatory

60,000 compute tasks
Input Data: 5000 files (10GB total)
Output Data: 60,000 files (60GB total)

executed on LIGO Data Grid,
Open Science Grid and XSEDE



Advanced LIGO PyCBC Workflow

One of the main pipelines to measure the statistical significance of data needed for discovery

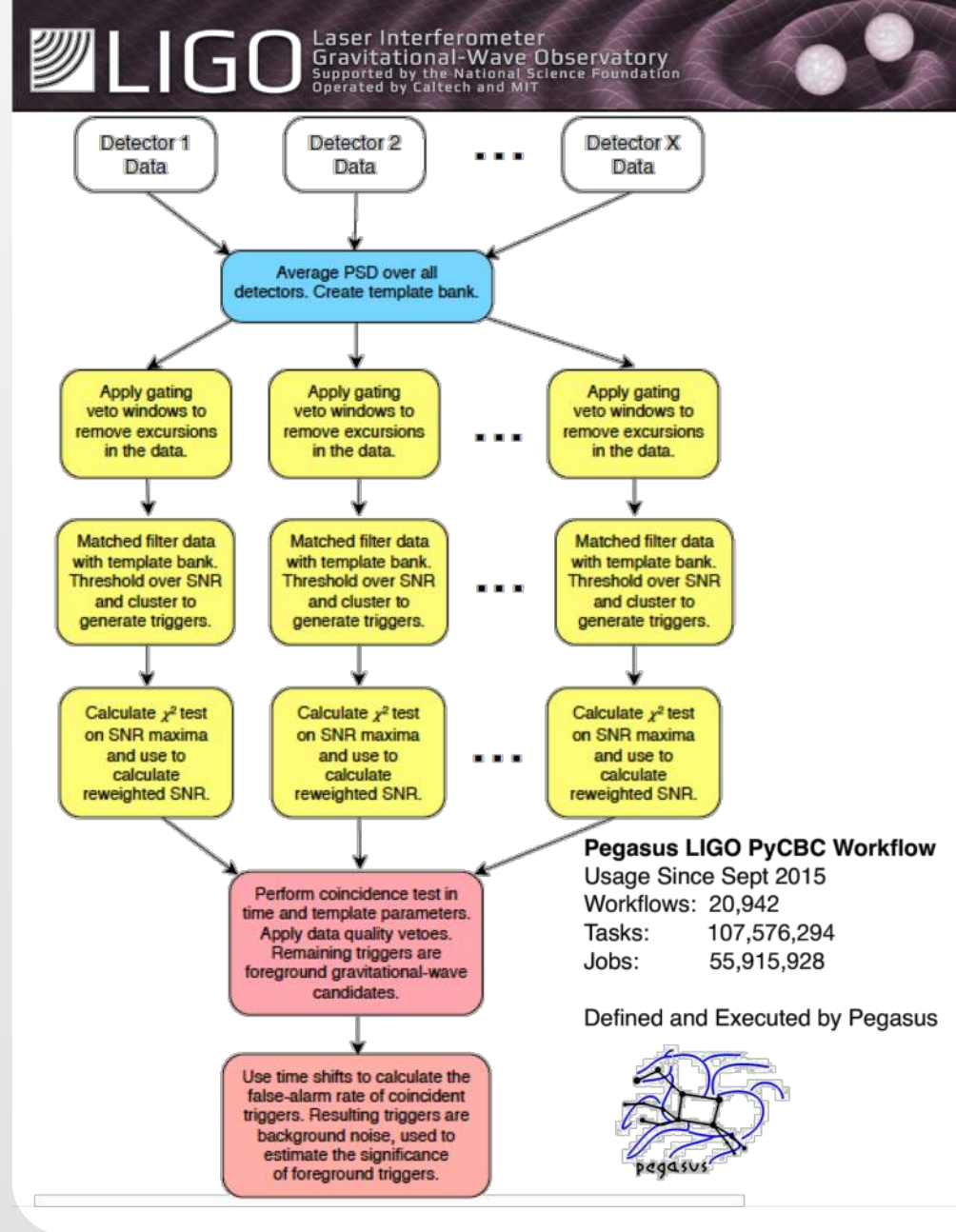
Contains **100's of thousands of jobs** and accesses on order of **terabytes of data**

Uses data from multiple detectors

For the detection, the pipeline was executed on Syracuse and Albert Einstein Institute Hannover

A single run of the binary black hole + binary neutron star search through the O1 data (about 3 calendar months of data with 50% duty cycle) requires a **workflow** with **194,364 jobs**

Generating the final O1 results with all the review required for the first discovery took about **20 million core hours**



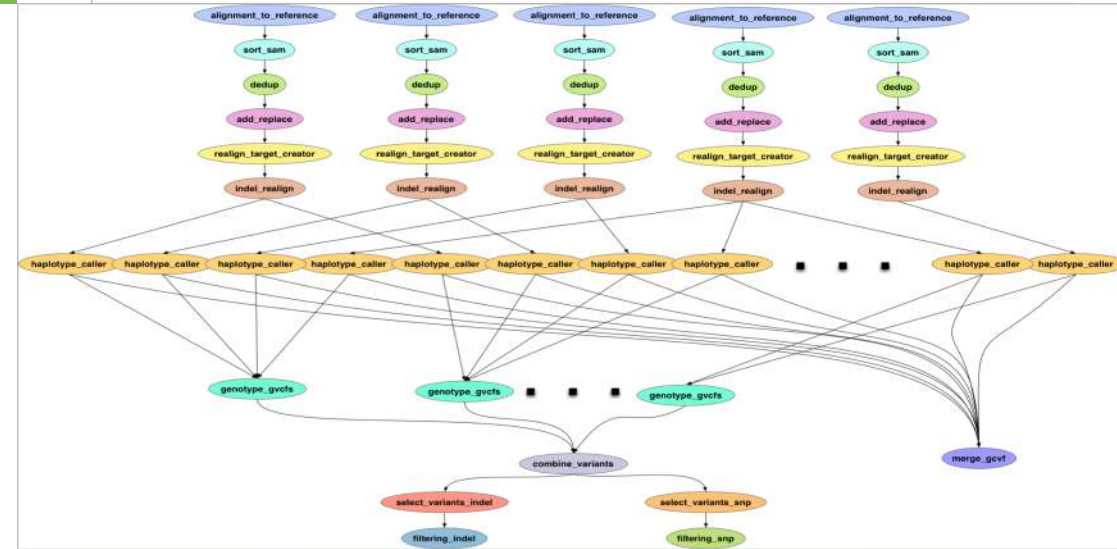
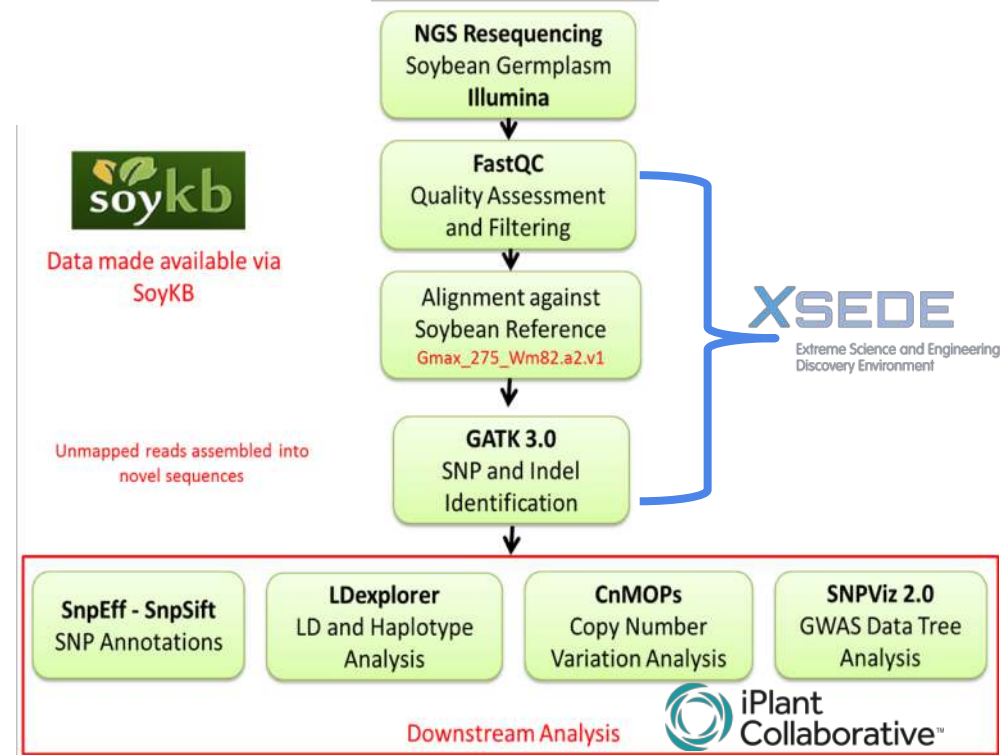
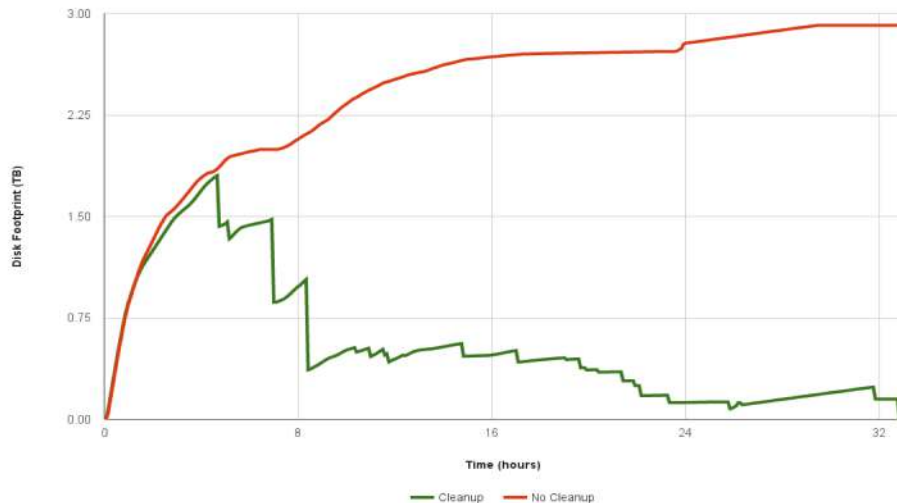
Soybean Workflow

TACC Wrangler as Execution Environment

Flash Based Shared Storage

Switched to glideins (pilot jobs) - Brings in remote compute nodes and joins them to the HTCondor pool on the submit host - Workflow runs at a finer granularity

Works well on Wrangler due to more cores and memory per node (48 cores, 128 GB RAM)





Pegasus

est. 2001

Automate, recover, and debug scientific computations.

Get Started

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

Support

pegasus-support@isi.edu

Pegasus Website

<http://pegasus.isi.edu>

Users Mailing List

pegasus-users@isi.edu



Pegasus

Workflow Management System

<https://pegasus.isi.edu/jobs>

 **USC** University of
Southern California
Information Sciences Institute

We're Hiring!