

Pegasus

Workflow Management System

Karan Vahi
Ewa Deelman

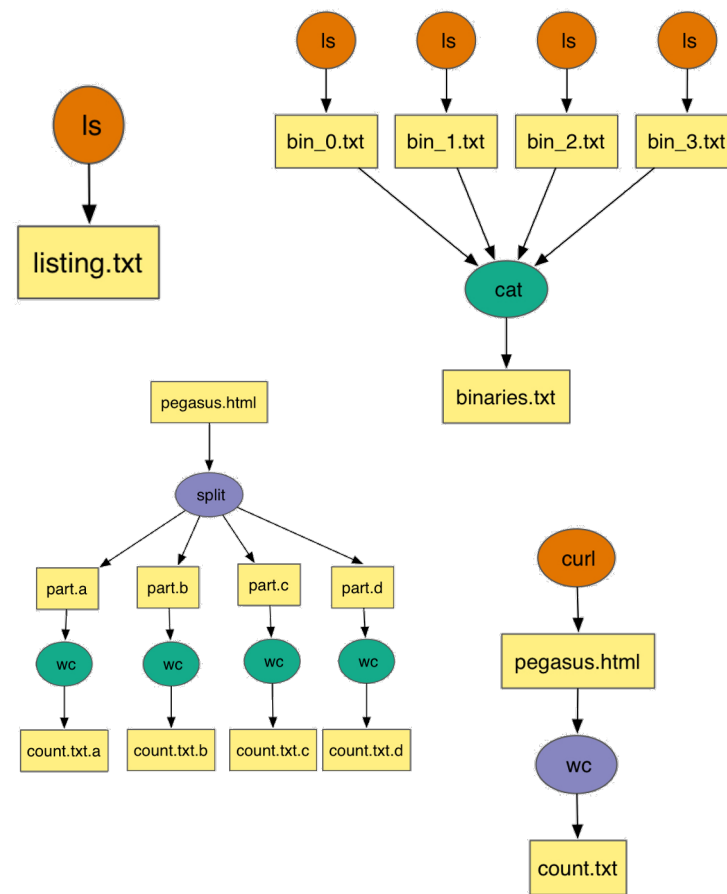
University of Southern California, School of Engineering
Information Sciences Institute
vahi@isi.edu, deelman@isi.edu



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.
 - 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.

Workflow Building Blocks



Slide Content Courtesy of David Okaya, SCEC, USC



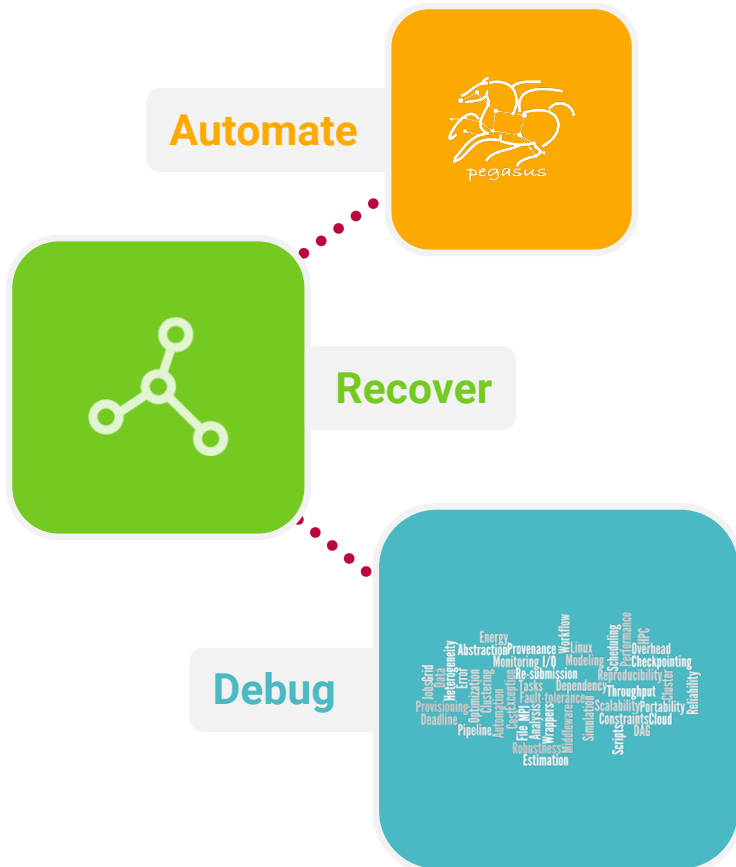
Workflow Challenges Across Domains

- 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
- Ease of use. Ability to debug and monitor large workflows

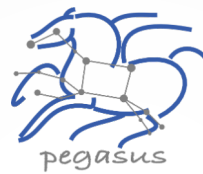
Our Focus

- 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.

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



Some of The Success Stories...

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



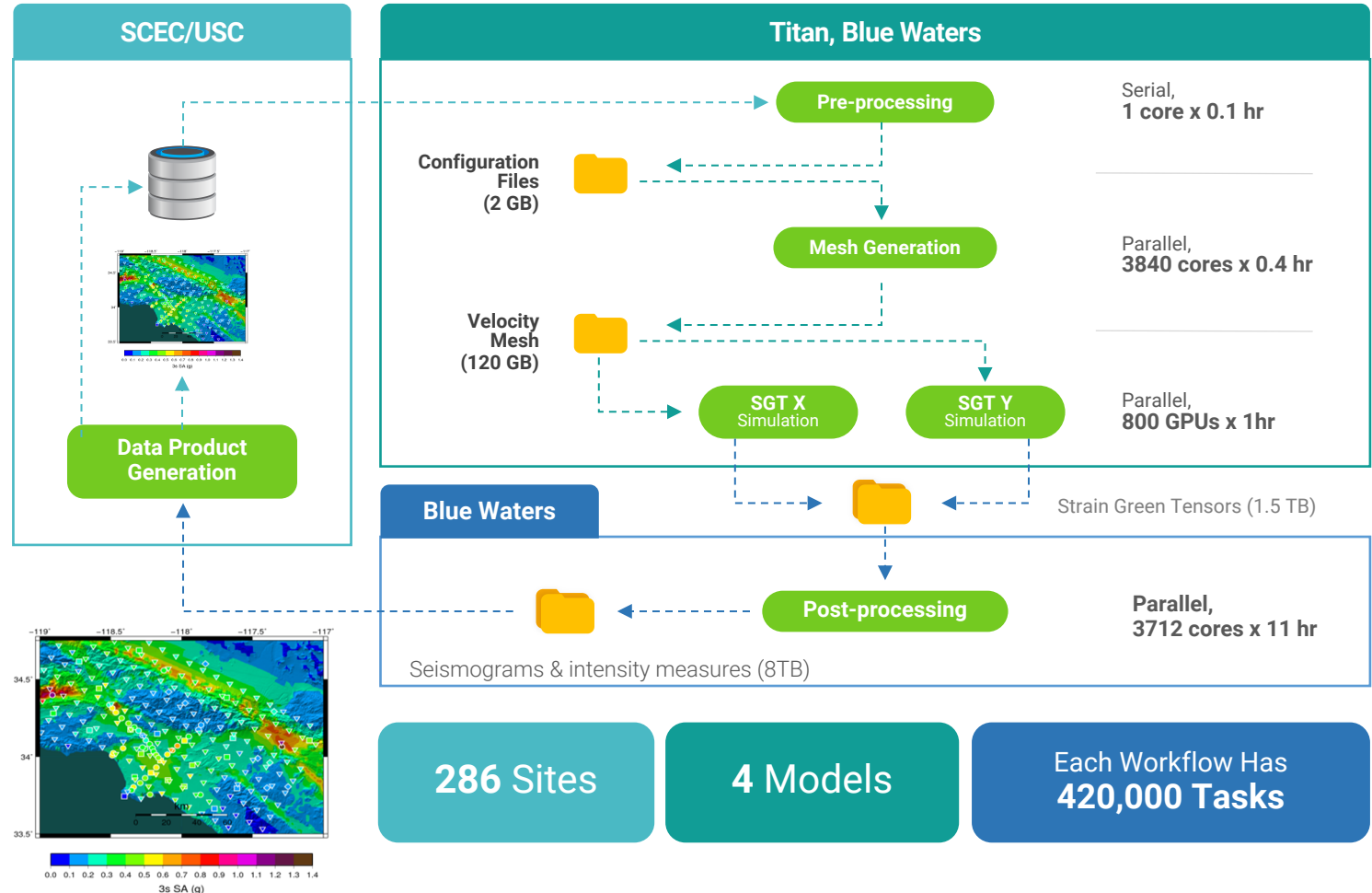
GPU jobs:
439,000 node-hours
AWP-ODC finite-difference code
5 billion points per volume, 23,000 timesteps
200 GPUs for 1 hour



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



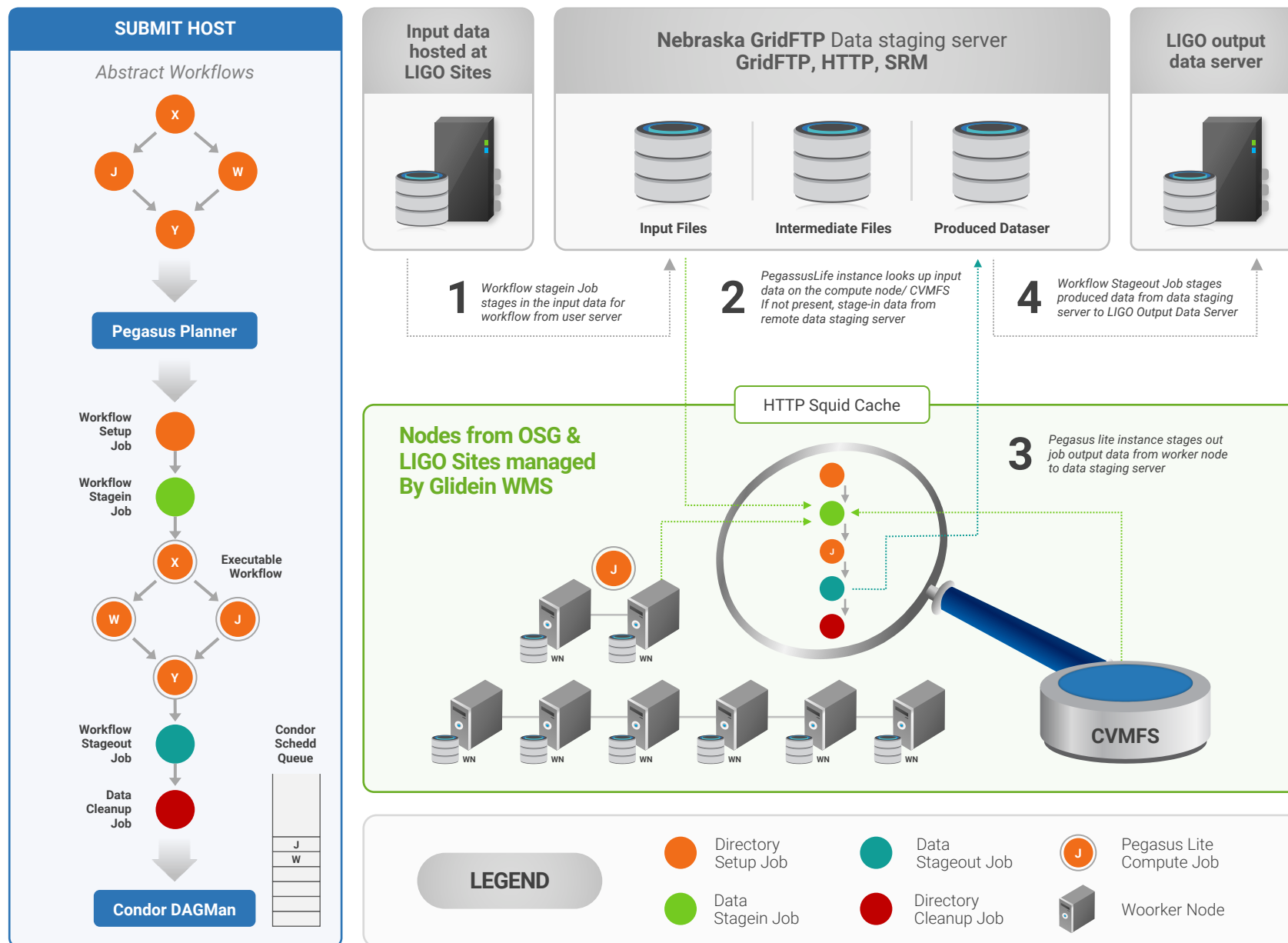
Data Flow for LIGO Pegasus Workflows in OSG

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



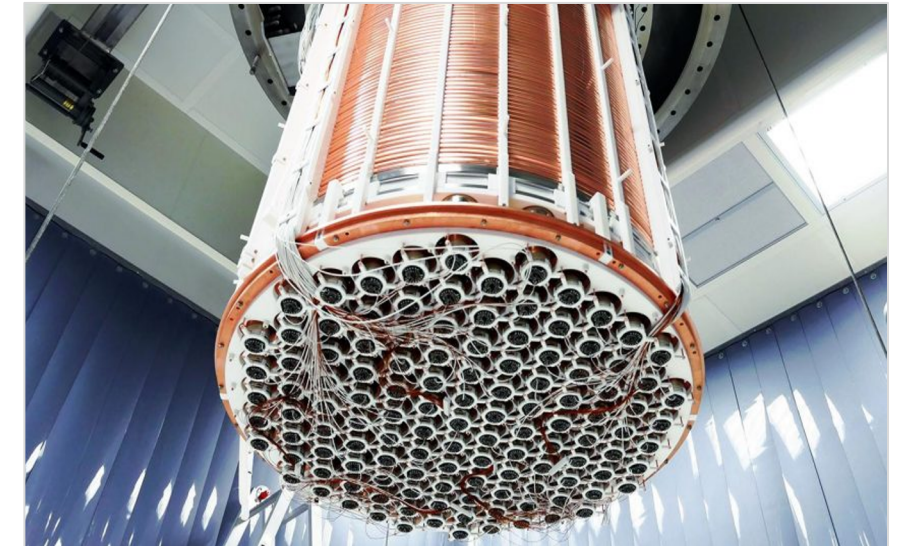
XENONnT - Dark Matter Search



Two Workflows

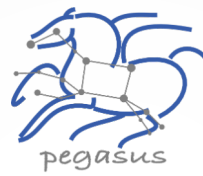
Monte Carlo simulations and the main processing pipeline.

- 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



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

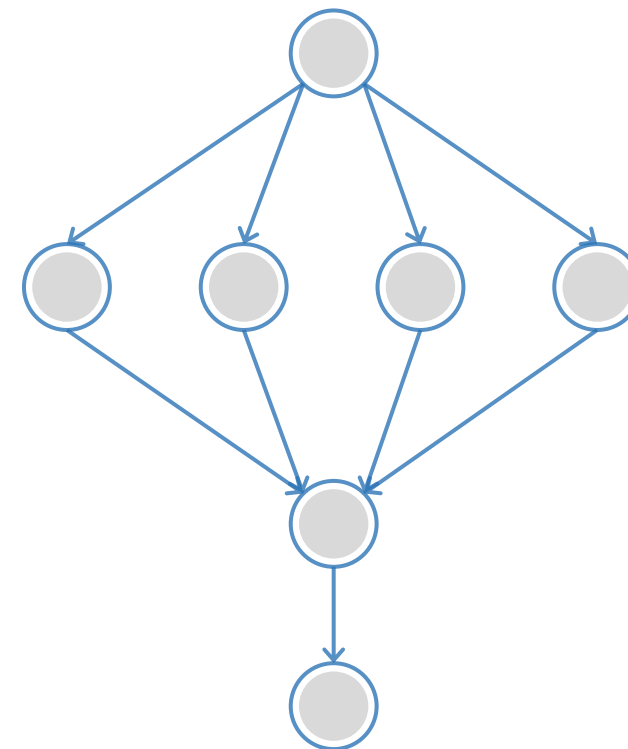
▲ Workflows are DAGs

- 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



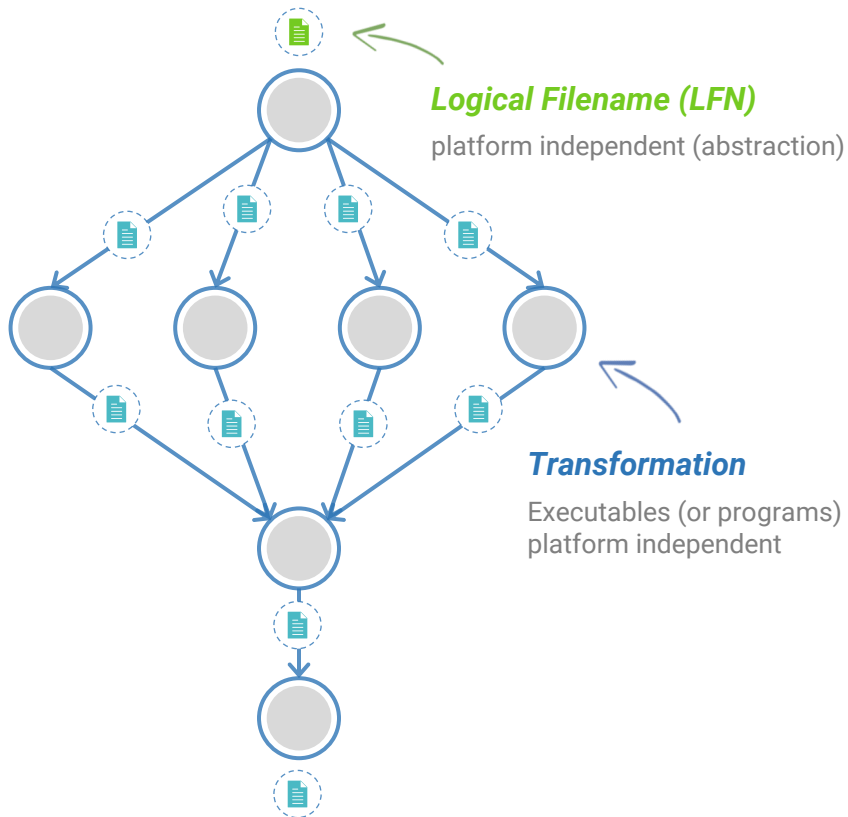


Input Workflow Specification **YAML formatted**

Portable Description

Users do not worry about low level execution details

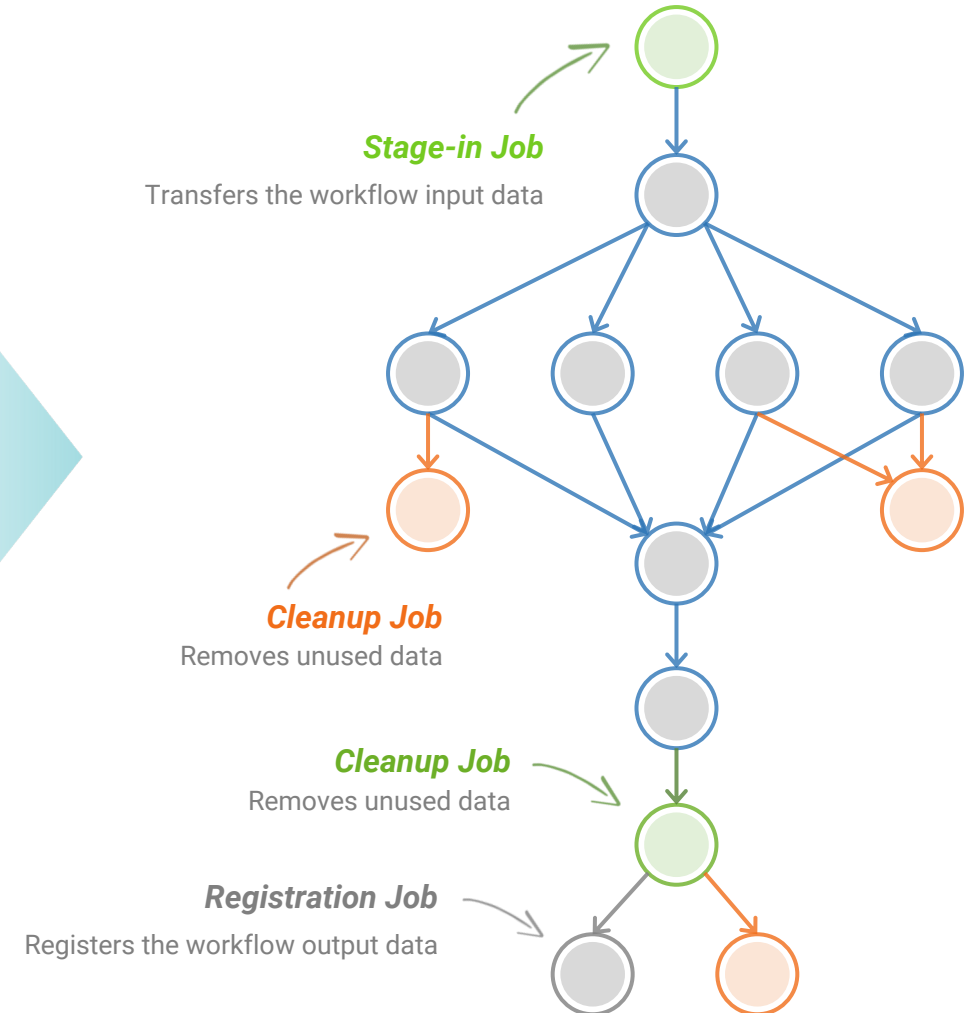
ABSTRACT WORKFLOW



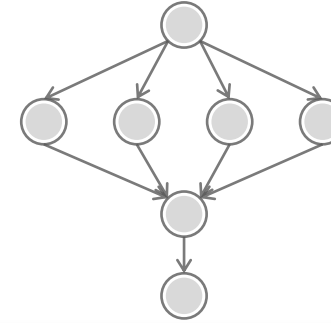
directed-acyclic graphs

Output Workflow

EXECUTABLE WORKFLOW



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()
```



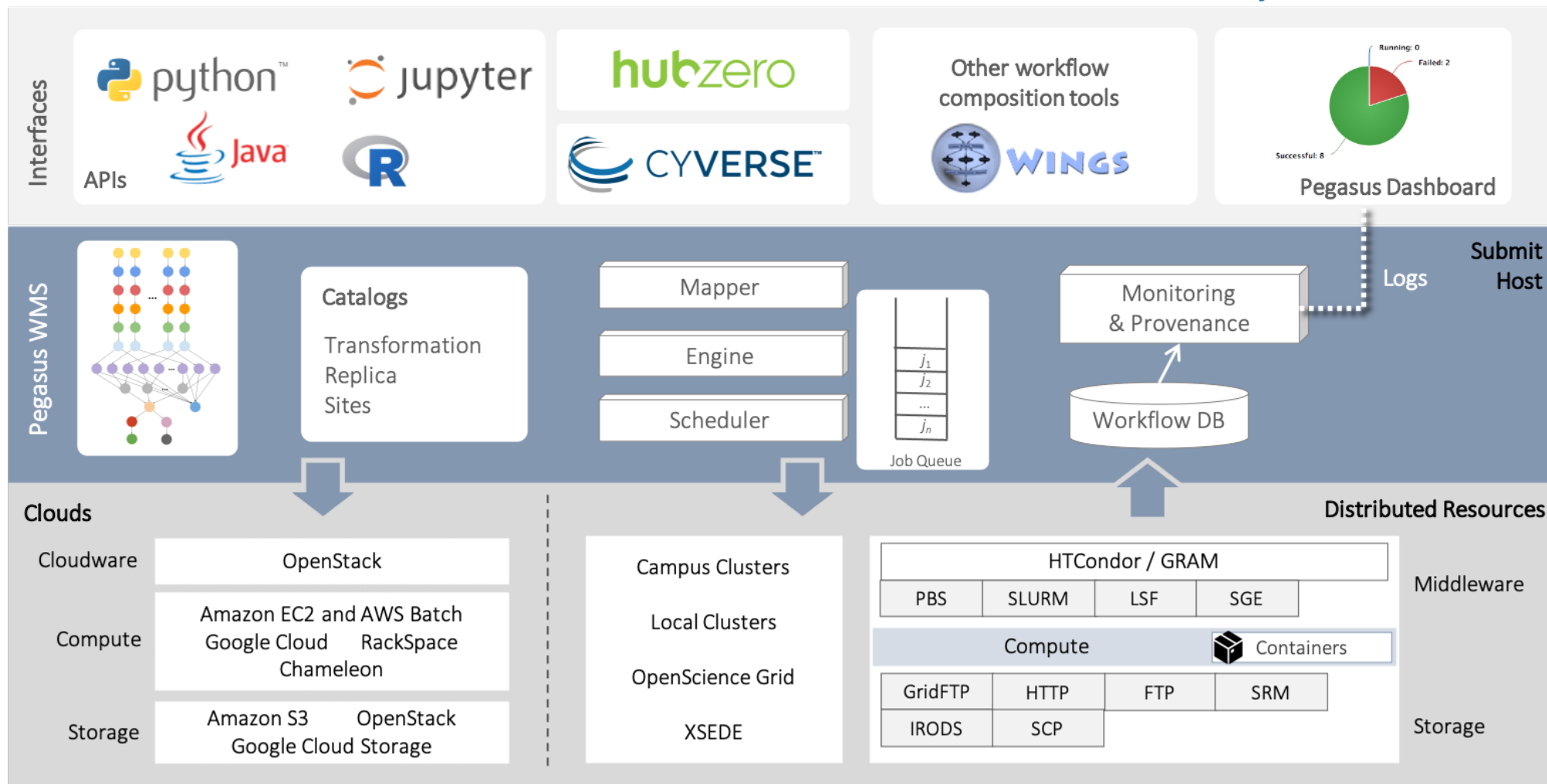
```
x-pegasus:
  apiLang: python
  createdBy: vahi
  createdOn: 11-19-20T14:57:58Z
  pegasus: '5.0'
  name: pipeline
  jobs:
    - type: job
      name: curl
      id: ID0000001
      arguments:
        - -o
        - pegasus.html
        - http://pegasus.isi.edu
      uses:
        - lfn: pegasus.html
          type: output
          stageOut: false
          registerReplica: false
    - type: job
      name: wc
      id: ID0000002
      stdout: count.txt
      arguments:
        - -l
        - pegasus.html
      uses:
        - lfn: count.txt
          type: output
          stageOut: true
          registerReplica: true
        - lfn: pegasus.html
          type: input
      jobDependencies:
        - id: ID0000001
          children:
            - ID0000002
```

YAML Formatted

Abstract Workflow

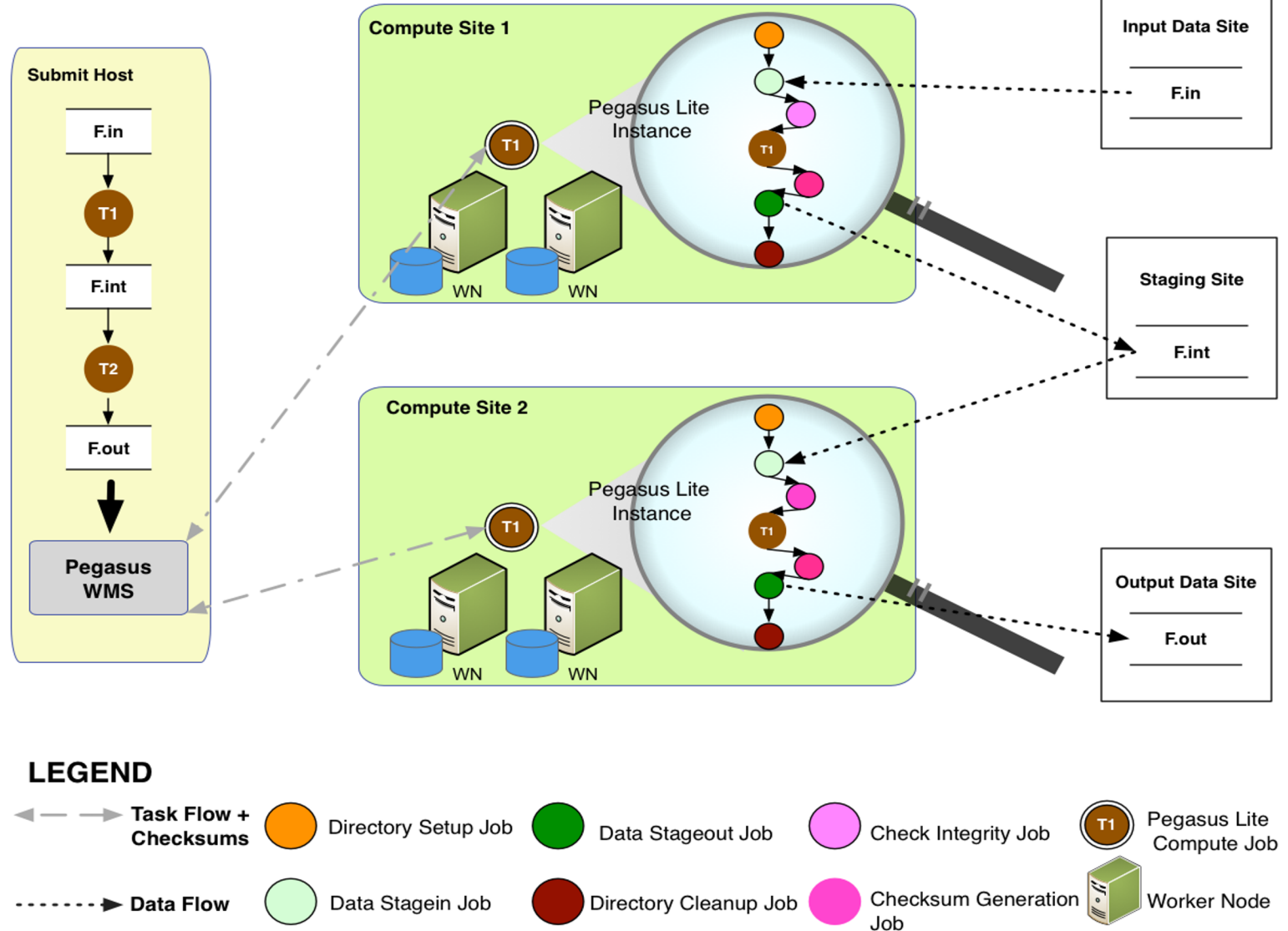


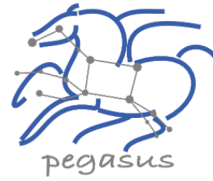
System Architecture



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 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						
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

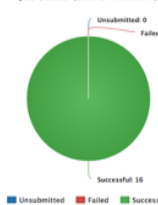
Real-time **monitoring** of workflow executions. It shows the **status** of the workflows and jobs, job **characteristics, statistics** and **performance** metrics.

Provenance data is stored into a relational database.

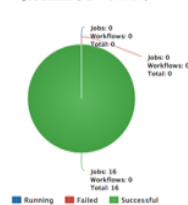
Workflow Details

Label	split
Type	root-ef
Progress	Successful
Submit Host	workflow.isi.edu
User	pegtrain01
Submit Directory	/nfs/ccg3/cog/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...

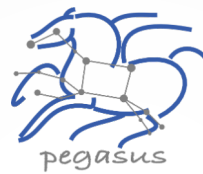
*****Summary*****

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**



Understanding Pegasus Features...

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

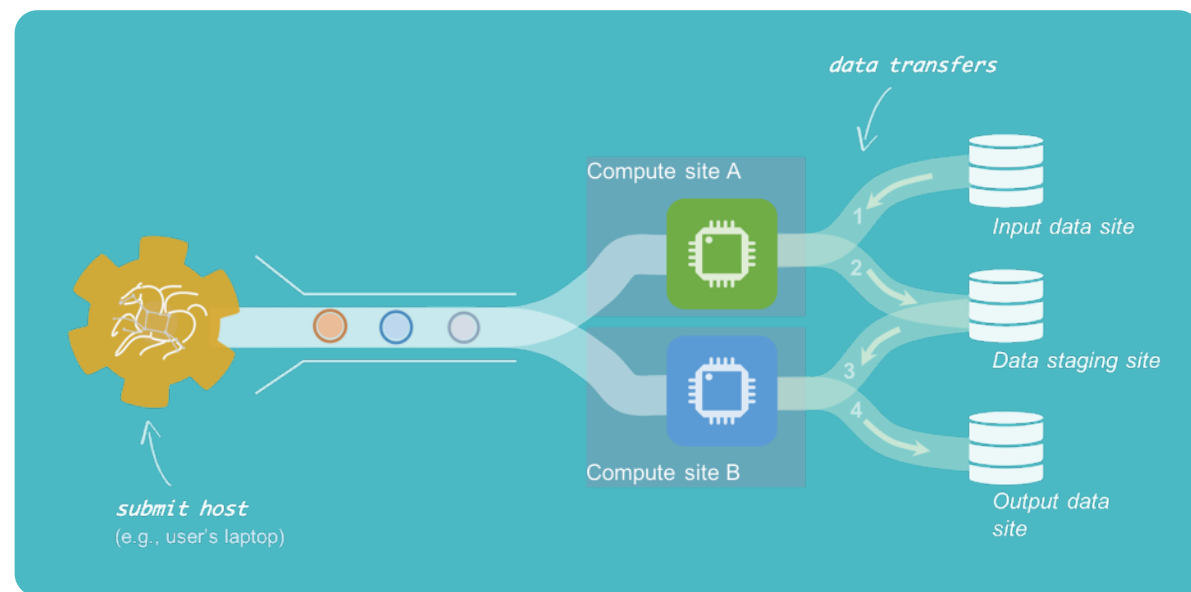
Non-shared File System (clouds, OSG, ...)

- 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



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

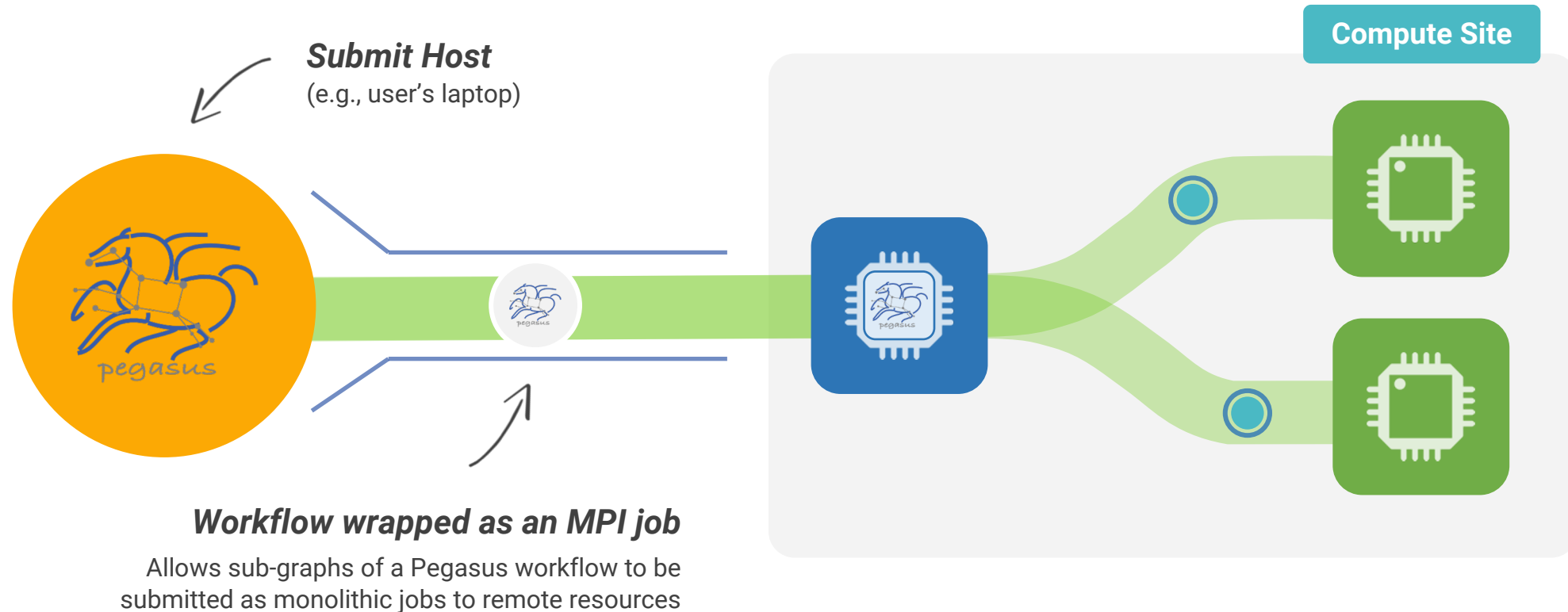


Credential management

- 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

Running fine-grained workflows on HPC systems...





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.

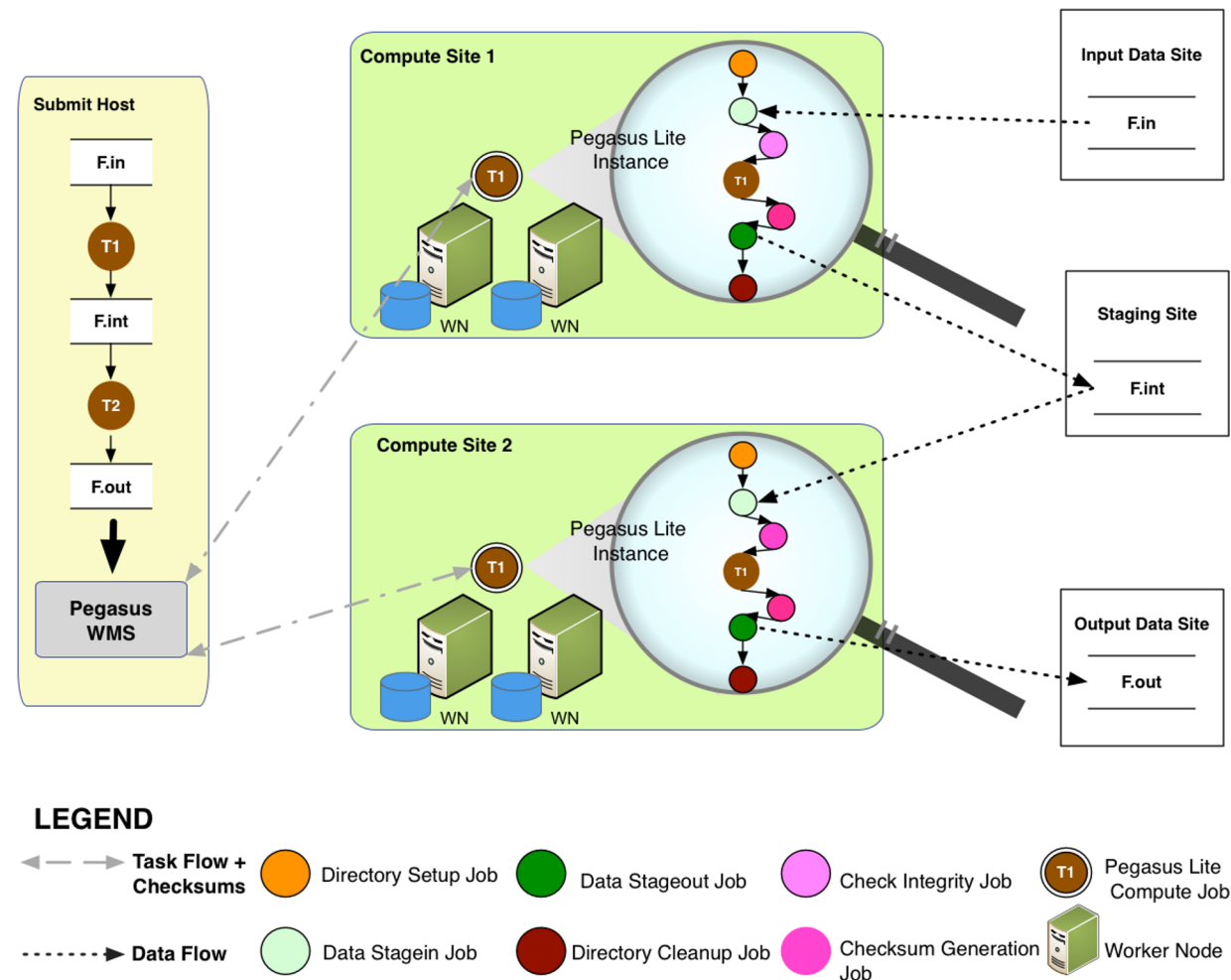
Plus there is the threat of intentional changes: malicious attackers, insider threats, etc.

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 Container Support



Containers Execution Model



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



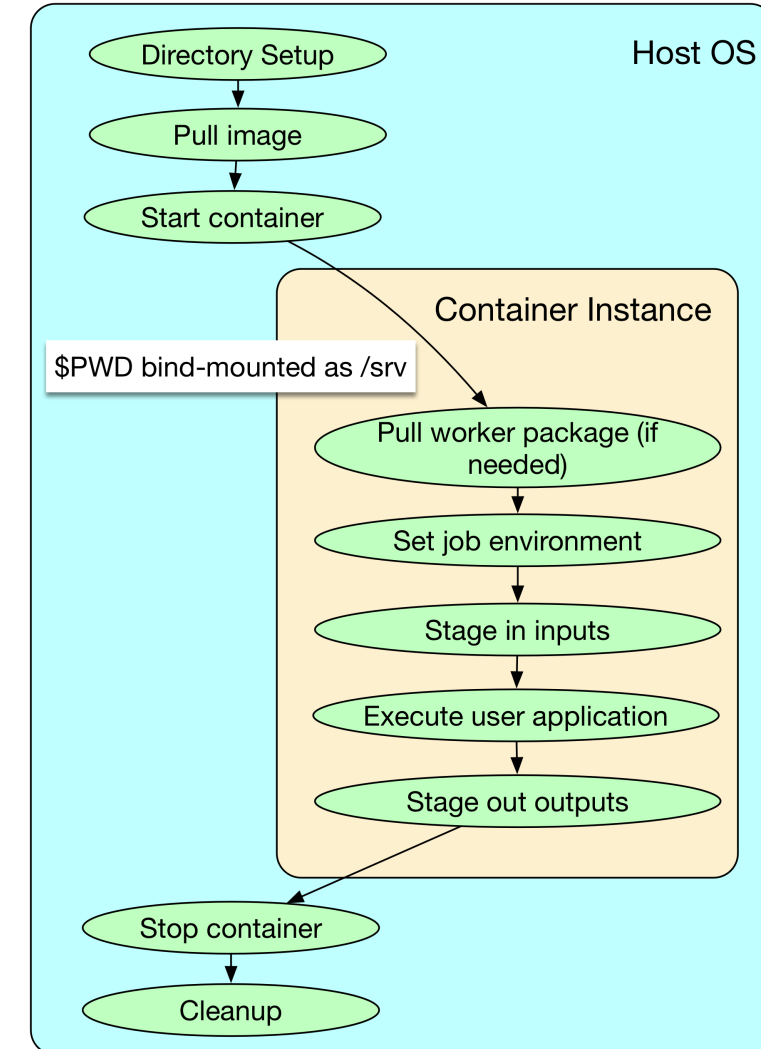
Users can **refer** to a **container** they want to **use** – **Pegasus stages** their executables and containers to the node

- 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

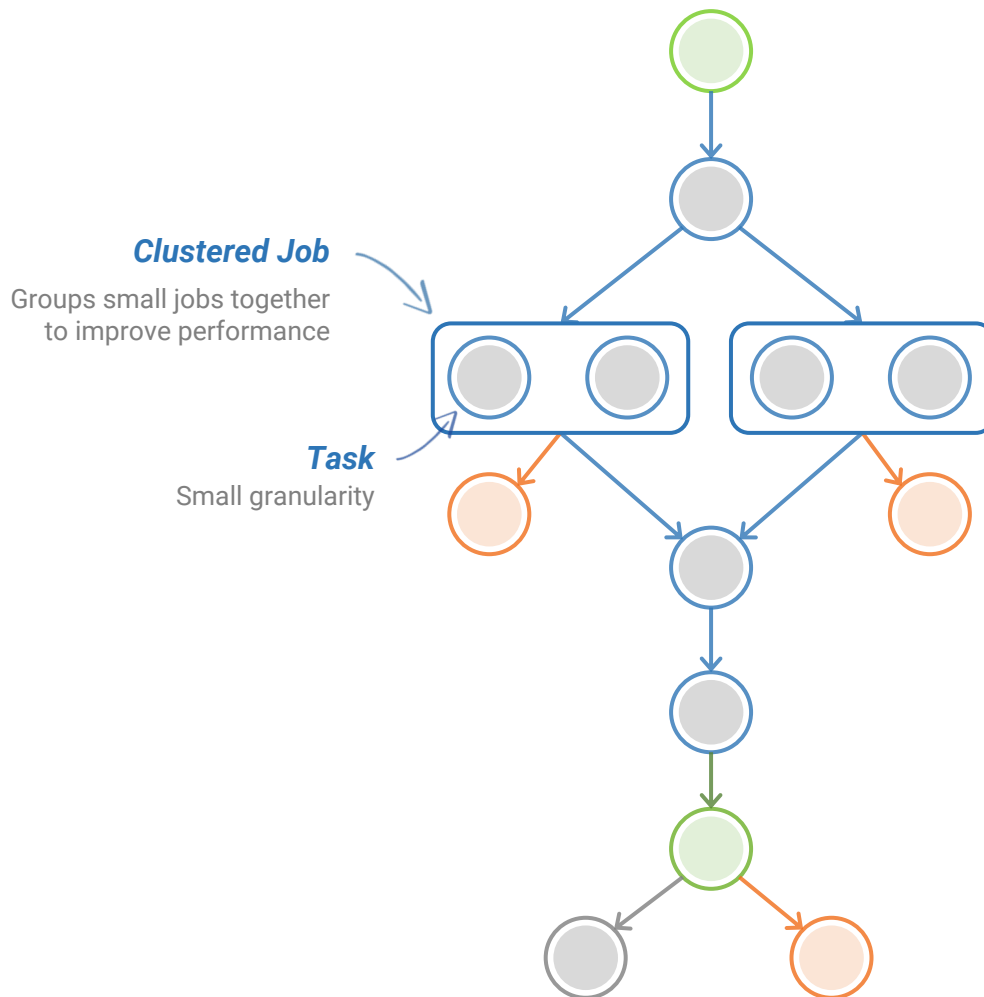
- The image is pulled down as a tar file as part of data stage-in jobs in the workflow
- The exported tar file is then shipped with the workflow and made available to the jobs
- Pricing considerations. You are now charged if you exceed a certain rate of pulls from Hubs

Other Optimizations

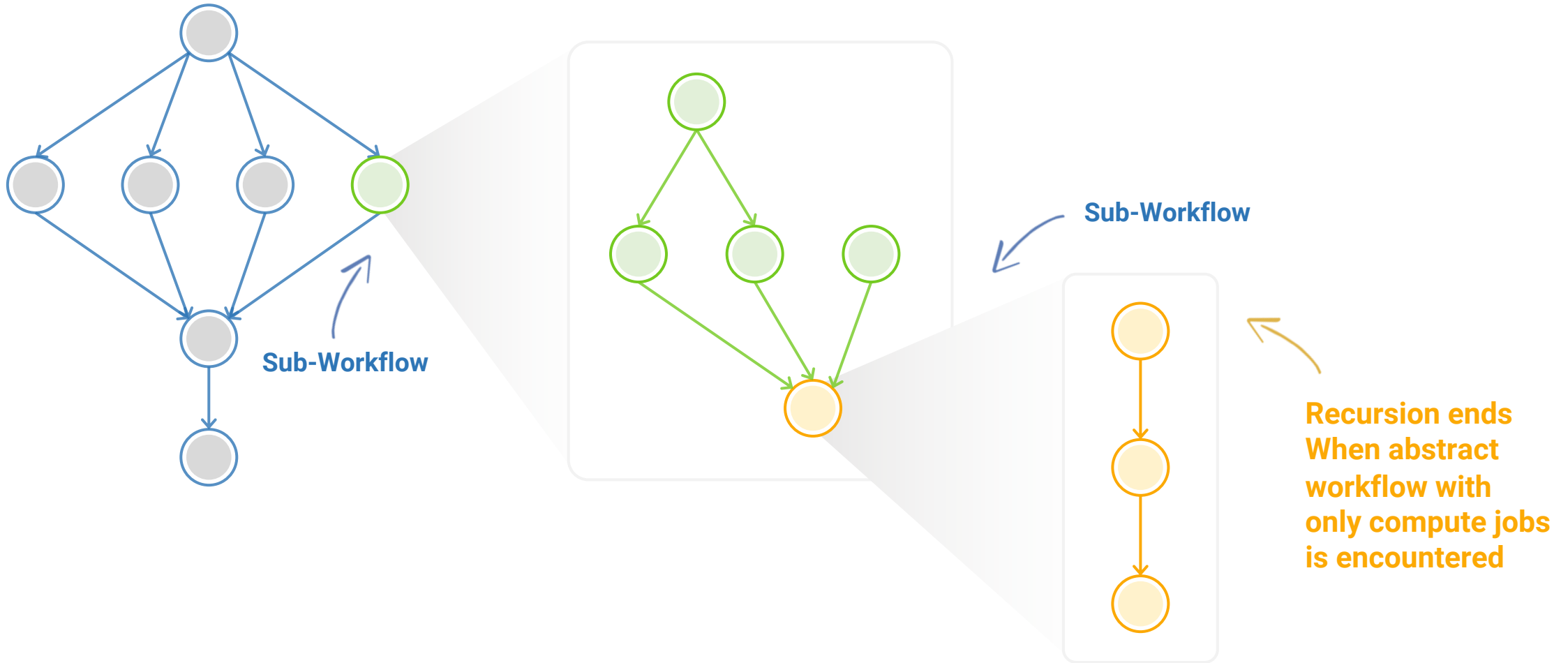
- Symlink against **existing images** on shared filesystem such as **CVMFS**
- The exported tar file is then shipped with the workflow and made available to the jobs



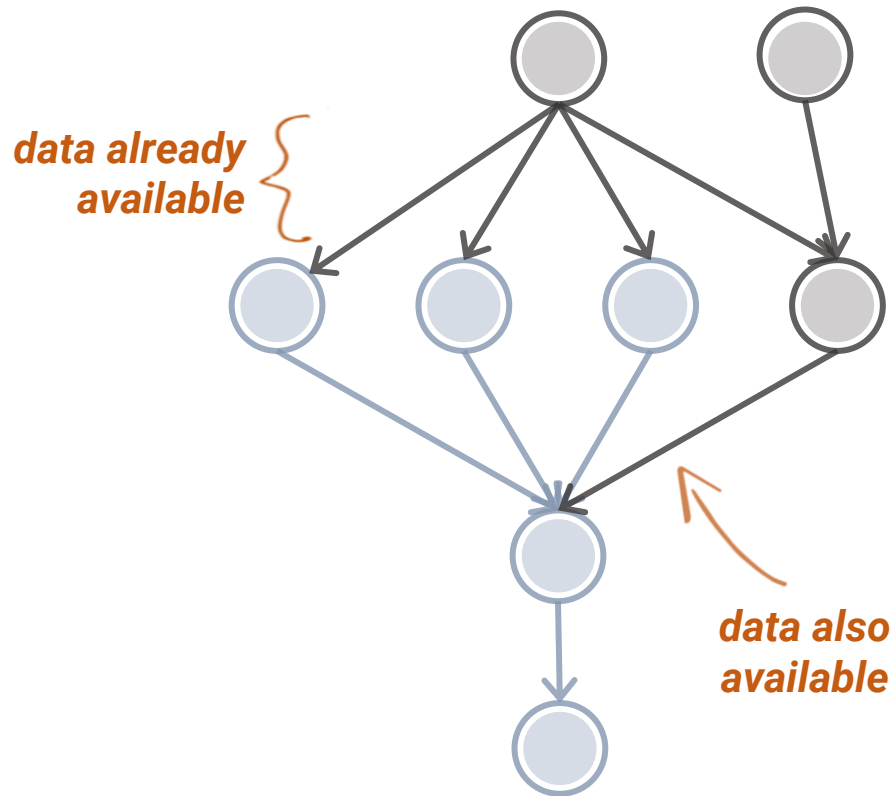
Performance. Why not improve it?



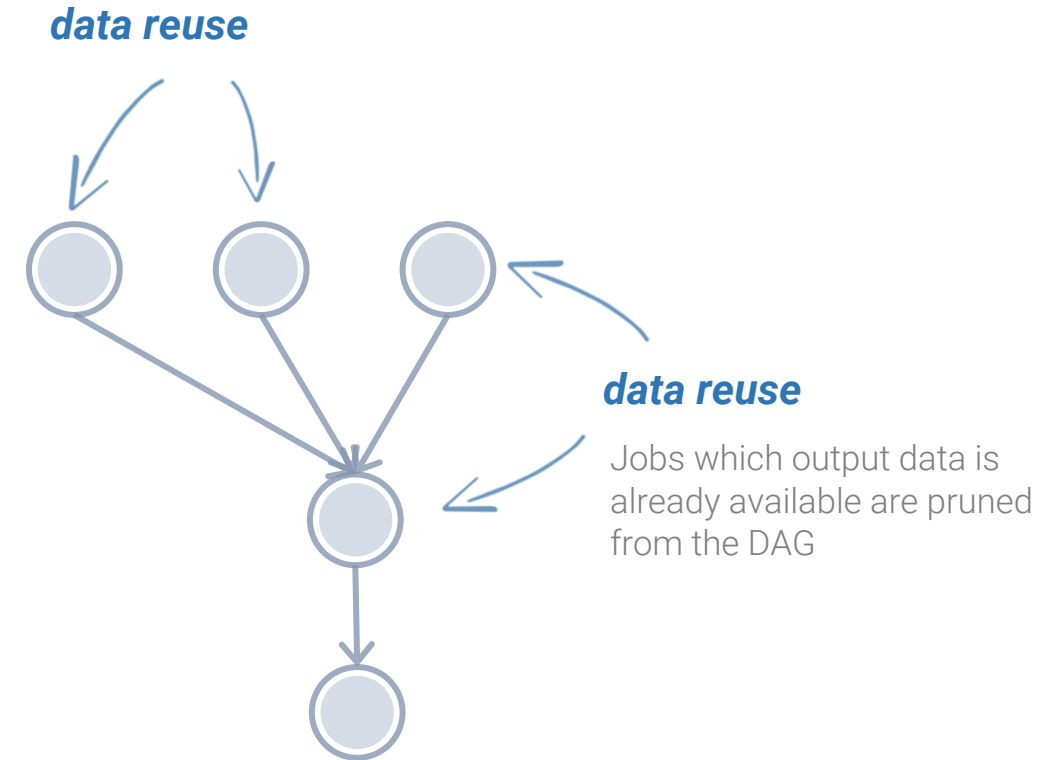
Pegasus also handles large-scale workflows



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

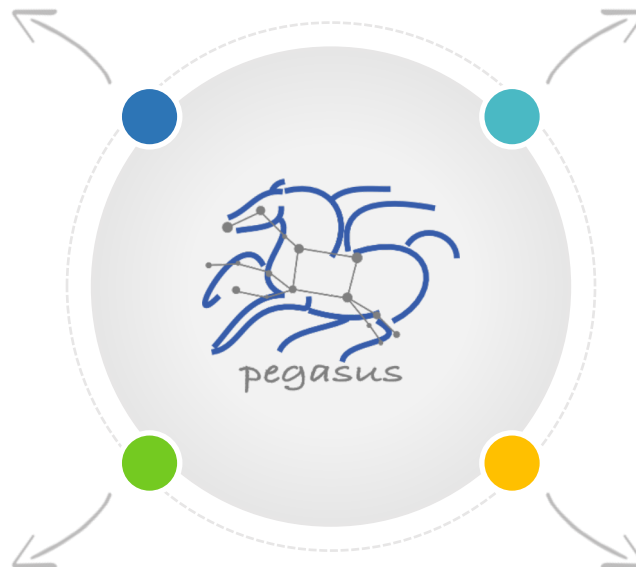


helps with transient failures
set number of retries per
job and run

Rescue DAGs



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



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

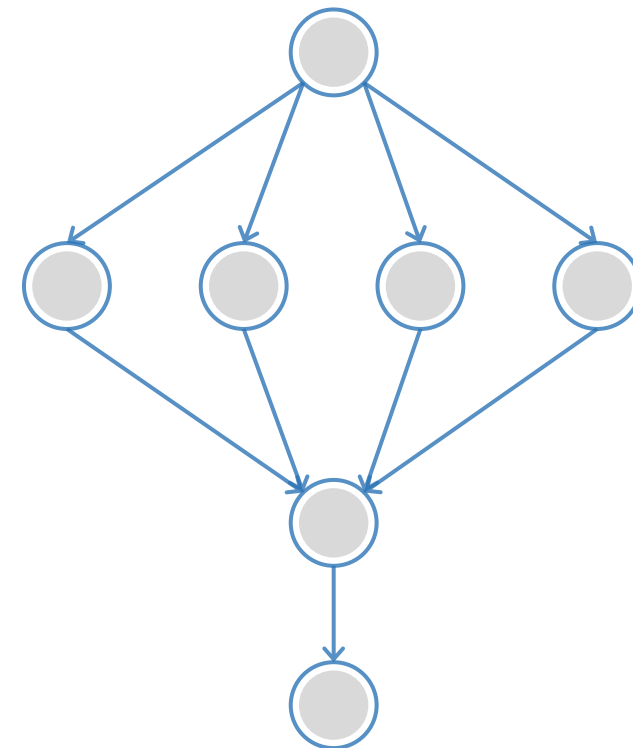
- Users can specify credentials in a **generic credentials file** on submit host
- Associate credentials with sites in site catalog

Approach

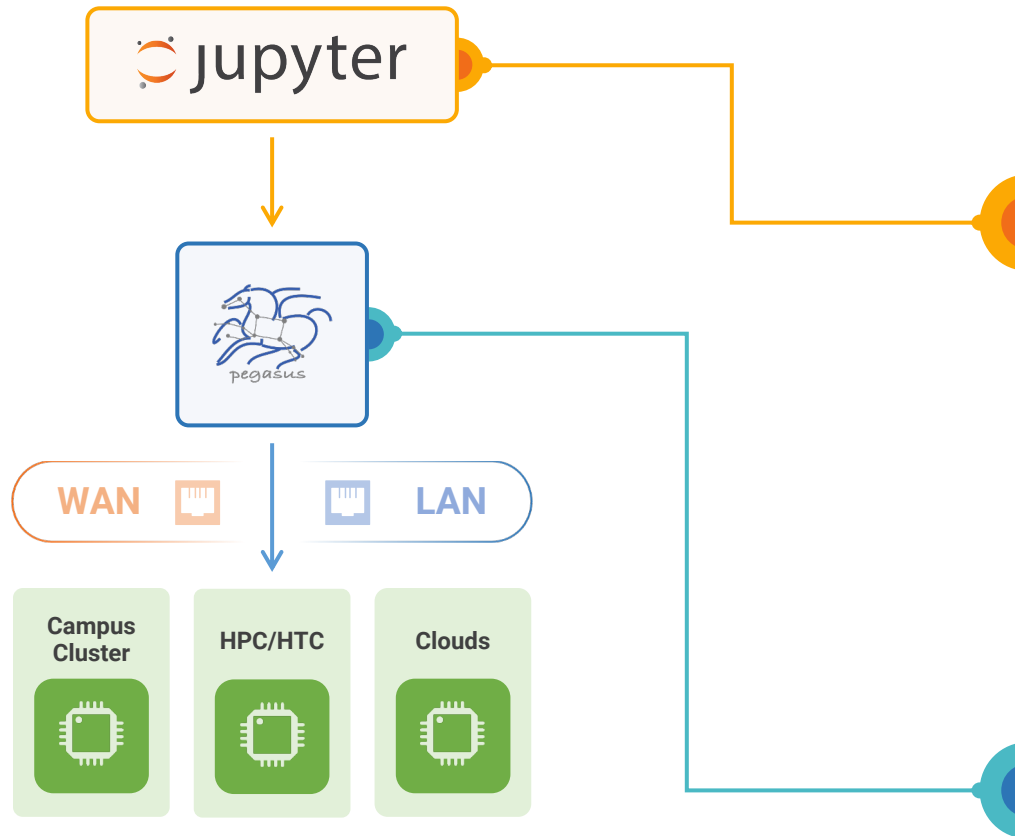
- Planner will **automatically** associate the **required credentials** with each job
- The credentials are **transferred** along with the job
- Usually available **only for the duration** of the job **execution**

▲ Supported Credentials

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



Running Pegasus workflows with Jupyter



The screenshot shows a Jupyter notebook titled "Pegasus-Tutorial-Split" with a last checkpoint of "03/15/2017 (autosaved)". The notebook displays a Directed Acyclic Graph (DAG) representing a workflow. The DAG consists of several nodes: four orange nodes at the top (wc_ID000003, wc_ID000002, wc_ID000005, wc_ID000004), several green nodes in the middle (clean_up_local_level_4_0, stage_out_local_level_1_1, stage_out_local_level_1_0, clean_up_local_level_4_1, clean_up_local_level_3_0), and several yellow nodes at the bottom (clean_up_local_level_5_0, register_local_1_1, clean_up_local_level_5_1, register_local_1_0). Arrows indicate the flow of the workflow from top to bottom, converging on a final node labeled "cleanup_split_0_local".

After the workflow has been submitted you can monitor it using the `status()` method. This method takes two arguments:

- `loop`: whether the status command should be invoked once or continuously until the workflow is completed or a failure is detected.
- `delay`: The delay (in seconds) the status will be refreshed. Default value is 10s.

```
In [6]: instance.status(loop=True, delay=5)
```

Progress: 100.0% (Success) (Completed: 17, Queued: 0, Running: 0, Failed: 0)

File for submitting this DAG to Condor: split-0.dag.condor.sub
Log of DAGMan debugging messages: split-0.dag.dagman.out
Log of Condor library output: split-0.dag.lib.out
Log of Condor library error messages: split-0.dag.lib.err
Log of the life of condor_dagman itself: split-0.dag.dagman.log

Your database is compatible with Pegasus version: 4.7.0
Submitting to condor split-0.dag.condor.sub
Submitting job(s).
1 job(s) submitted to cluster 1068.

Your workflow has been started and is running in the base directory: /Users/silva/Downloads/split-submit-host-2017-03-27T10:17:45/submit/silva/pegasus/split/run0002

*** To monitor the workflow you can run ***

```
pegasus-status -l /Users/silva/Downloads/split-submit-host-2017-03-27T10:17:45/submit/silva/pegasus/split/run0002
```




Released Nov, 2020

- 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

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

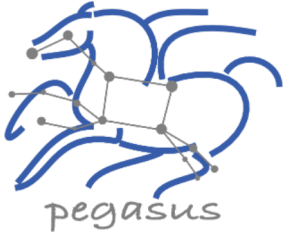
from Pegasus.api import *

# 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

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

▶ 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