# Pegasus

Automate, recover, and debug scientific computations.

**Rafael Ferreira da Silva**

USCViterbi
School of Engineering
Information Sciences Institute

http://pegasus.isi.edu

# Experiment Timeline

**Scientific Problem**

Earth Science, Astronomy, Neuroinformatics, Bioinformatics, etc.

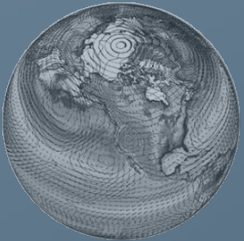**Computational Scripts**

Shell scripts, Python, Matlab, etc.

**Distributed Computing**

Clusters, HPC, Cloud, Grid, etc.

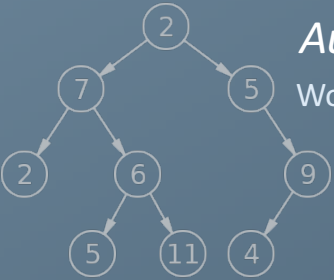**Scientific Result**

Models, Quality Control, Image Analysis, etc.

**Analytical Solution**

**Automation**

Workflows, MapReduce, etc.

**Monitoring and Debug**

Fault-tolerance, Provenance, etc.



Pegasus

# *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

*Automate*

*Recover*

*Debug*

# Taking a closer look into a workflow...

**DAG**

directed-acyclic graphs

*job*

Command-line programs

*split*

*merge*

*dependency*

Usually data dependencies

*pipeline*

**DAX**

DAG in XML

**Pegasus**

# From the abstraction to execution!

*stage-in job*

Transfers the workflow input data

*stage-out job*

Transfers the workflow output data

*registration job*

Registers the workflow output data

**Pegasus**

*http://pegasus.isi.edu*

# Optimizing storage usage...

*cleanup job*

Removes unused data

**Pegasus**

*http://pegasus.isi.edu*

7

# In a nutshell…

# …and all automatically!

executable
workflow

abstract
workflow

storage
constraints

**Pegasus**

# Pegasus also provides tools to generate the abstract workflow

```
dax = ADAG("test_dax")
firstJob = Job(name="first_job")
firstInputFile = File("input.txt")
firstOutputFile = File("tmp.txt")
firstJob.addArgument("input=input.txt", "output=tmp.txt")
firstJob.uses(firstInputFile, link=Link.INPUT)
firstJob.uses(firstOutputFile, link=Link.OUTPUT)
dax.addJob(firstJob)
for i in range(0, 5):
    simulJob = Job(id="%s" % (i+1), name="simul_job")
    simulInputFile = File("tmp.txt")
    simulOutputFile = File("output.%d.dat" % i)
    simulJob.addArgument("parameter=%d" % i, "input=tmp.txt",
        output=%s" % simulOutputFile.getName())
    simulJob.uses(simulInputFile, link=Link.INPUT)
    simulJob.uses(simulOutputFile, line=Link.OUTPUT)
dax.addJob(simulJob)
dax.depends(parent=firstJob, child=simulJob)
fp = open("test.dax", "w")
dax.writeXML(fp)
fp.close()
```

DAG in XML

DAX

# While you wait…

## …or the execution is finished.

Does everything
executed successfully?

How my workflow
behaves?

*Web-based interface*

Real-time monitoring, graphs,
provenance, etc.

*Debug*

Set of debugging tools to
unveil issues

*Statistics*

Workflow execution and
job performance metrics

*RESTful API*

Monitoring and reporting information
on your own application interface

Past executions?

*Command-line tools*

Tools for monitor and debug workflows

**Pegasus**

*http://pegasus.isi.edu*

# 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

# 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

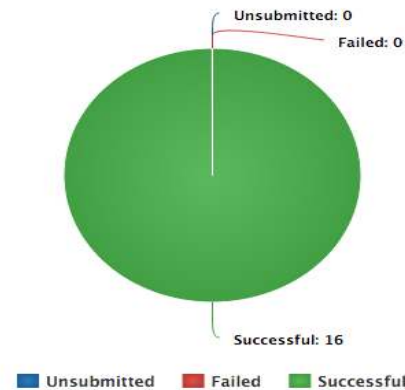| | |
|---|---|
| Label | split |
| Type | root-wf |
| Progress | Successful |
| Submit Host | workflow.isi.edu |
| User | pegtrain01 |
| Submit Directory | /nfs/ccg3/ccg/home/pegtrain01/examples/split/split/run0002 |
| DAGMan Out File | split-0.dag.dagman.out |
| Wall Time | 12 mins 23 secs |
| Cumulative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

**Job Status (Entire Workflow)**

Unsubmitted: 0
Failed: 0
Successful: 16

Unsubmitted   Failed   Successful

**Job Status (Per Workflow)**

Jobs: 0
Workflows: 0
Total: 0

Jobs: 0
Workflows: 0
Total: 0

Jobs: 16
Workflows: 0
Total: 16

Running   Failed   Successful

## Pegasus

# >_  But, if you prefer the 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 :
```

...Pegasus provides a set of concise and powerful tools

## Pegasus

# And if a job fails?

## Job Failure Detection

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

## Job Retry

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

## Checkpoint Files

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

## Rescue DAGs

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

**Pegasus**

SRM

http

Local disk

Amazon S3

GridFTP

**Worried about**
# data?
Let Pegasus manage it for you

ftp

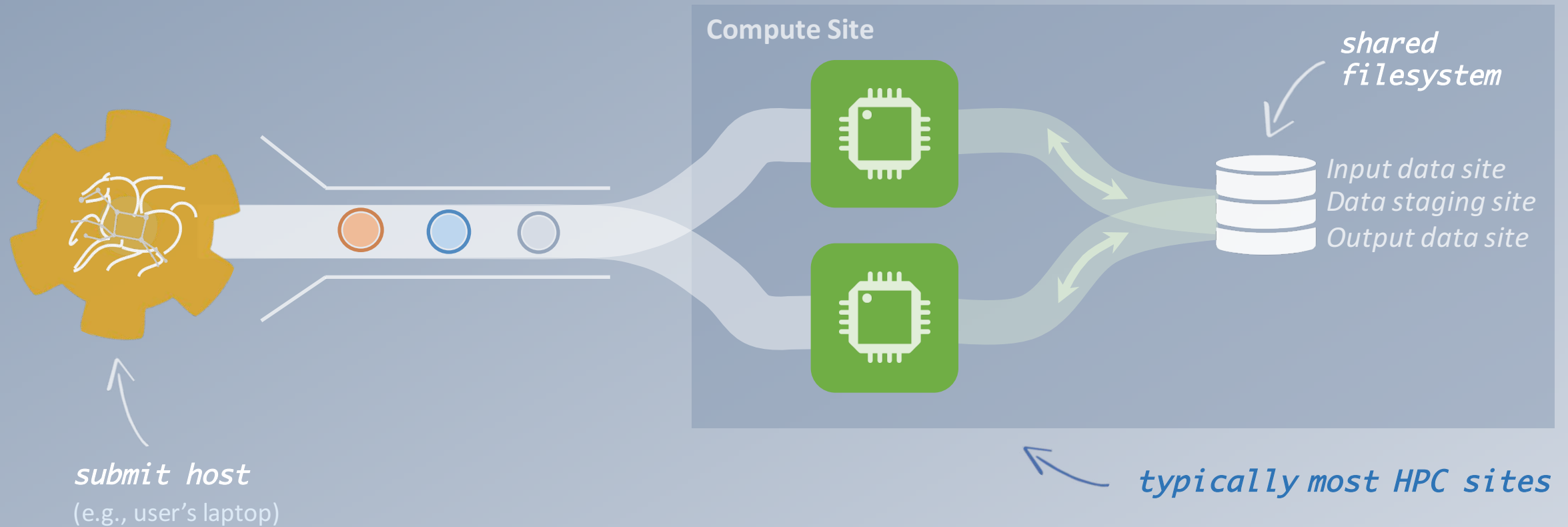Shared filesystem

scp

Google Storage

StashCache

iRODS

**Pegasus**

# How we handle it:

*data transfers*

Compute site A

Compute site B

Input data site

Data staging site

Output data site

1

2

3

4

*submit host*
(e.g., user's laptop)

Pegasus

# However, there are several possible configurations for data sites…



**Compute Site**

*shared filesystem*

Input data site
Data staging site
Output data site

*submit host*
(e.g., user's laptop)

*typically most HPC sites*

![Pegasus logo] Pegasus

# Pegasus also handles high-scalable object storages

Compute Site

object storage

Input data site
Data staging site
Output data site

Staging Site

submit host
(e.g., user's laptop)

Typical cloud computing deployment (Amazon S3, Google Storage)

# Pegasus can also manage data over the submit host…



Typical OSG sites
Open Science Grid

Compute Site

submit host
(e.g., user's laptop)

Pegasus

# And yes… you can mix everything!



Output data site

shared filesystem

Compute site A

Input data site
Data staging site

Data staging site

Compute site B

submit host
(e.g., user's laptop)

Input data site
Output data site

object storage

Pegasus

http://pegasus.isi.edu

20

# So, what information does Pegasus need?

## Transformation Catalog

describes all of the executables (called "transformations") used by the workflow

## Site Catalog

describes the sites where the workflow jobs are to be executed

## Replica Catalog

describes all of the input data stored on external servers

**Pegasus**

# How does Pegasus decide where to execute?

*site description*

describes the compute resources

*scratch*

tells where temporary data is stored

*storage*

tells where output data is stored

*profiles*

key-pair values associated per job level

```
...
  <!-- The local site contains information about the submit host -->
    <!-- The arch and os keywords are used to match binaries in the transformation
catalog -->
    <site handle="local" arch="x86_64" os="LINUX">

      <!-- These are the paths on the submit host were Pegasus stores data -->
      <!-- Scratch is where temporary files go -->
      <directory type="shared-scratch" path="/home/tutorial/run">
        <file-server operation="all" url="file:///home/tutorial/run"/>
      </directory>

      <!-- Storage is where pegasus stores output files -->
      <directory type="local-storage" path="/home/tutorial/outputs">
        <file-server operation="all" url="file:///home/tutorial/outputs"/>
      </directory>

      <!-- This profile tells Pegasus where to find the user's private key for SCP
transfers -->
      <profile namespace="env" key="SSH_PRIVATE_KEY">/home/tutorial/.ssh/id_rsa</profile>

    </site>
...
```

## Pegasus

# How does it know where the executables are or which ones to use?

**executables description**

list of executables locations per site

**physical executables**

mapped from logical transformations

**transformation type**

whether it is installed or available to stage

```
...
# This is the transformation catalog. It lists information about each of the
# executables that are used by the workflow.

tr ls {
    site PegasusVM {
        pfn "/bin/ls"
        arch "x86_64"
        os "linux"
        type "INSTALLED"
    }
}
...
```

## Pegasus

# What if data is not local to the submit host?

```
# This is the replica catalog. It lists information about each of the
# input files used by the workflow. You can use this to specify locations to input files
present on external servers.

# The format is:
# LFN PFN site="SITE"

f.a    file:///home/tutorial/examples/diamond/input/f.a    site="local"
```

*logical filename*

abstract data name

*physical filename*

data physical location on site
different transfer protocols
can be used (e.g., scp, http,
ftp, gridFTP, etc.)

*site name*

in which site the file is available

Pegasus

*http://pegasus.isi.edu*

# A few more features…

# Performance, why not improve it?

*clustered job*

Groups small jobs together
to improve performance

*task*

small granularity

**Pegasus**

*http://pegasus.isi.edu*

# What about **data reuse**?

*data already available*

*data also available*

*workflow reduction*

*data reuse*

*data reuse*

Jobs which output data is already available are pruned from the DAG

**Pegasus**

# Pegasus also handles
## large-scale workflows

*sub-workflow*

*sub-workflow*

*recursion ends
when DAX with
only compute jobs
is encountered*

# Running **fine-grained** workflows on HPC systems…

*submit host*
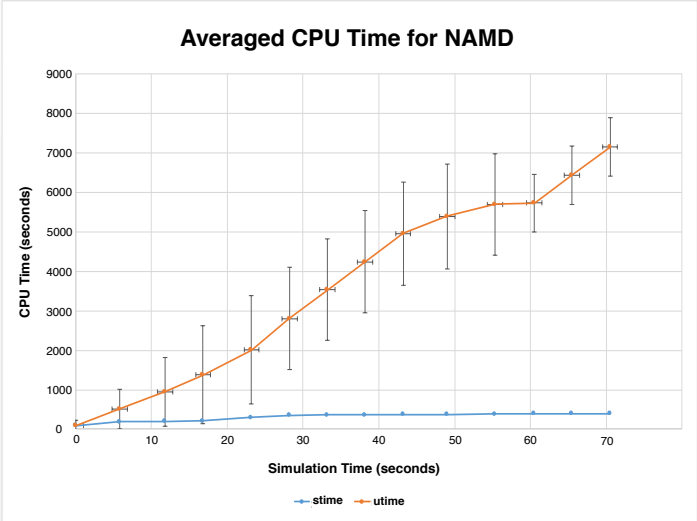(e.g., user's laptop)

**HPC System**

*Master*
*(rank 0)*

*worker*

*rank 1*

*rank n-1*

*workflow wrapped as an MPI job*

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

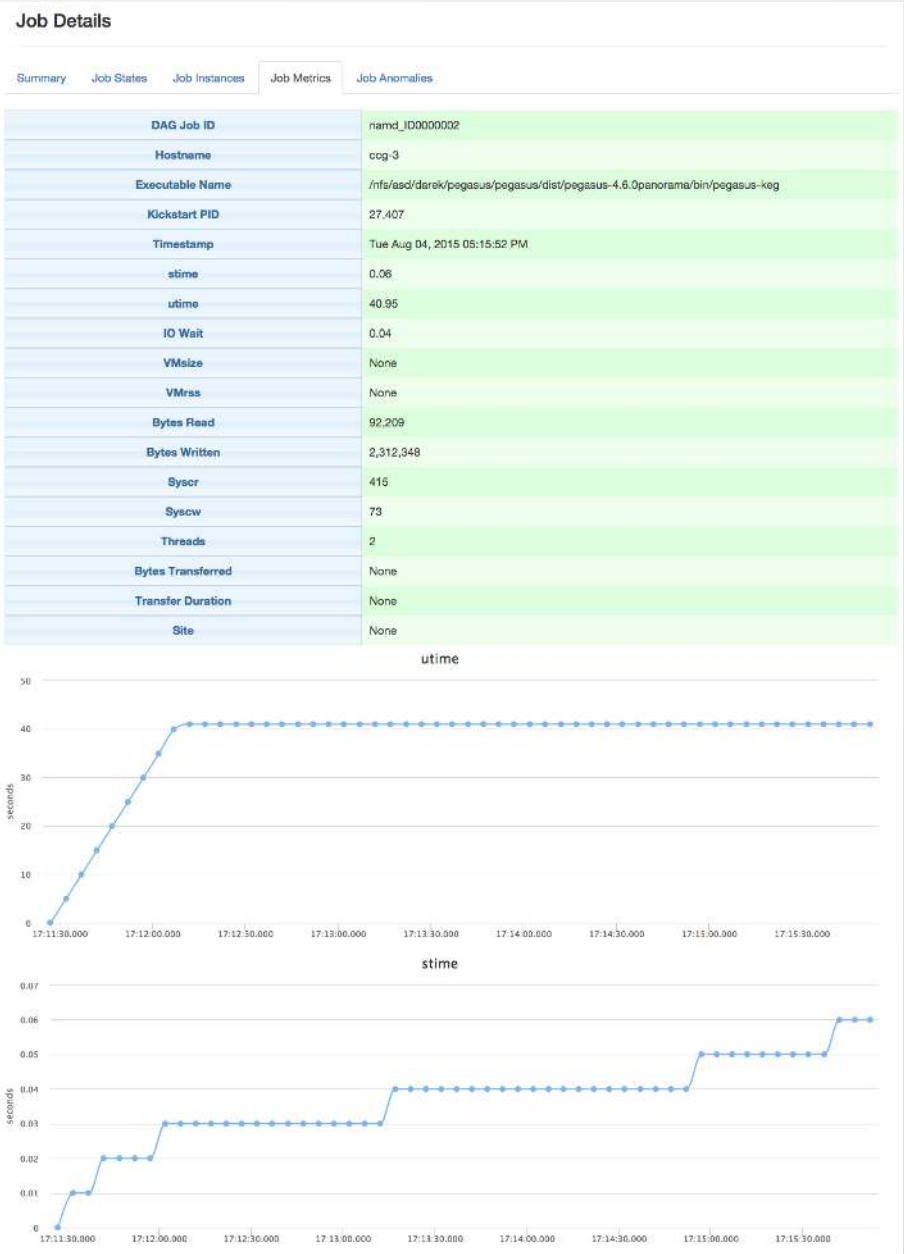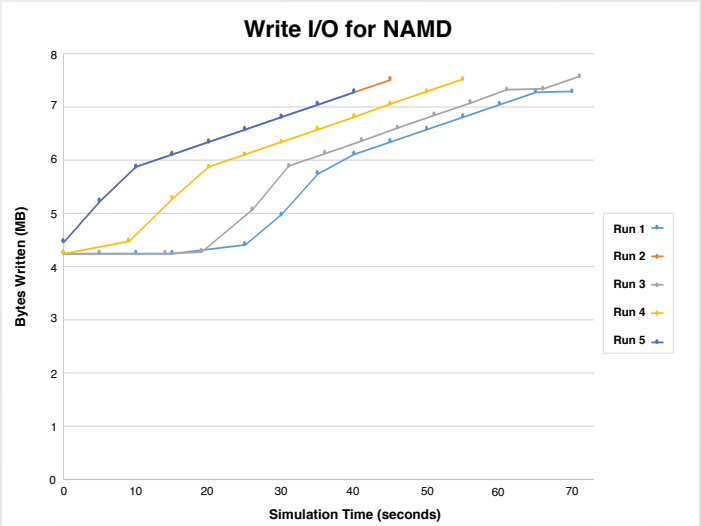**Pegasus**

*http://pegasus.isi.edu*

# Real-time collection of **time-series** of workflow performance metrics
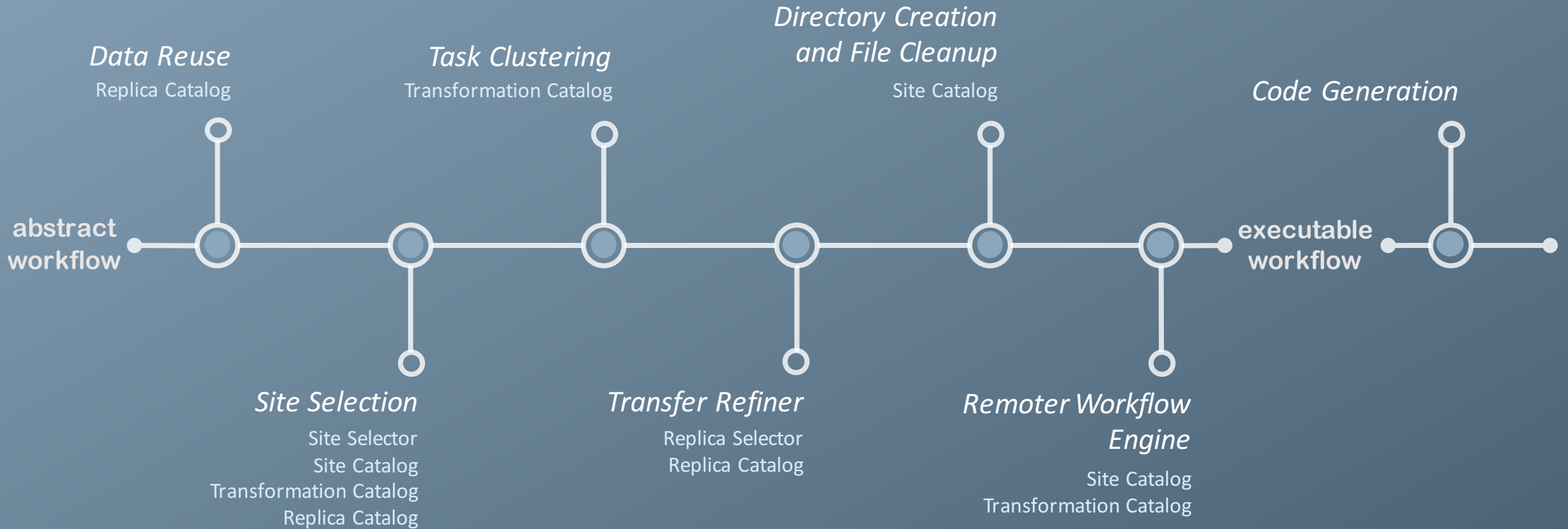


*time-series data in real-time*

integrated with Pegasus dashboard

*time-series data*

I/O (read, write), memory, CPU

**Pegasus**

# Pegasus' flow at a glance



*Data Reuse*
Replica Catalog

*Task Clustering*
Transformation Catalog

*Directory Creation
and File Cleanup*
Site Catalog

*Code Generation*

**abstract
workflow**

**executable
workflow**

*Site Selection*
Site Selector
Site Catalog
Transformation Catalog
Replica Catalog

*Transfer Refiner*
Replica Selector
Replica Catalog

*Remoter Workflow
Engine*
Site Catalog
Transformation Catalog

**Pegasus**

*http://pegasus.isi.edu*

**Science-grade Mosaic of the Sky
(Galatic Plane - Montage)**

18 million input images (~2.5TB)
900 output images (2.5GB each, 2.4TB total)
17 workflows, each of which contains
900 sub-workflows (hierarchical workflows)
10.5 million tasks (34,000 CPU hours)

executed on the cloud (Amazon EC2)

**SCEC CyberShake**

286 sites, 4 models
each workflow has 420,000 tasks
described as 21 jobs using PMC

executed on BlueWaters (NCSA)
and Stampede (TACC)

## How Pegasus has been used?

**Periodogram**

1.1M tasks grouped into 180 jobs
1.1M input, 12M output files
~101,000 CPU hours
16 TB output data

executed at SDSC

**ORNL Spallation Neutron Source
(SNS)**

5 jobs that consumes about
900 cores for more than 12 hours

executed on Hopper (NERSC)

🐎 **Pegasus**