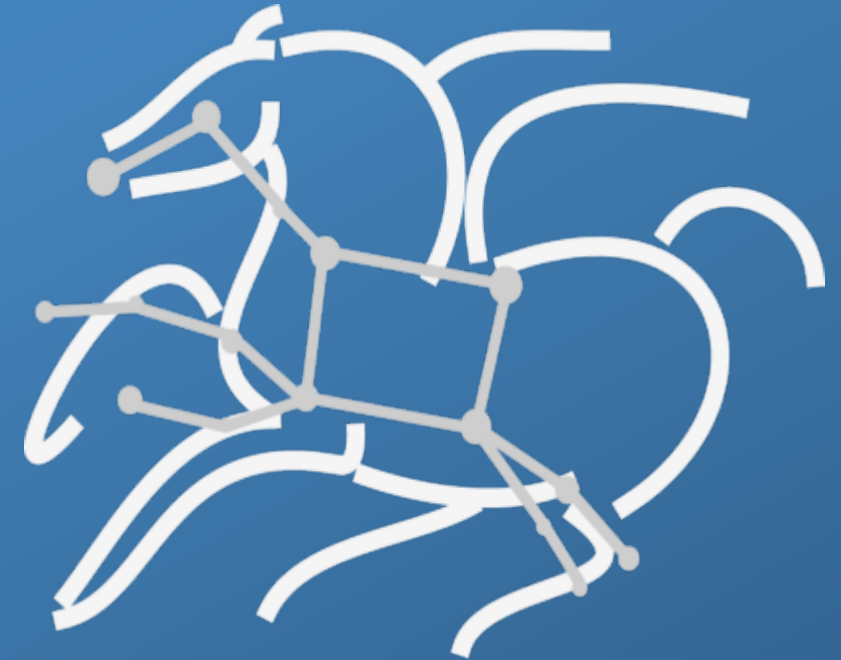# Pegasus

Automate, recover, and debug scientific computations.

**Mats Rynge**
rynge@isi.edu

USC Viterbi
School of Engineering
Information Sciences Institute

https://pegasus.isi.edu

# *Why* Pegasus *?*

Automates complex, multi-stage processing pipelines

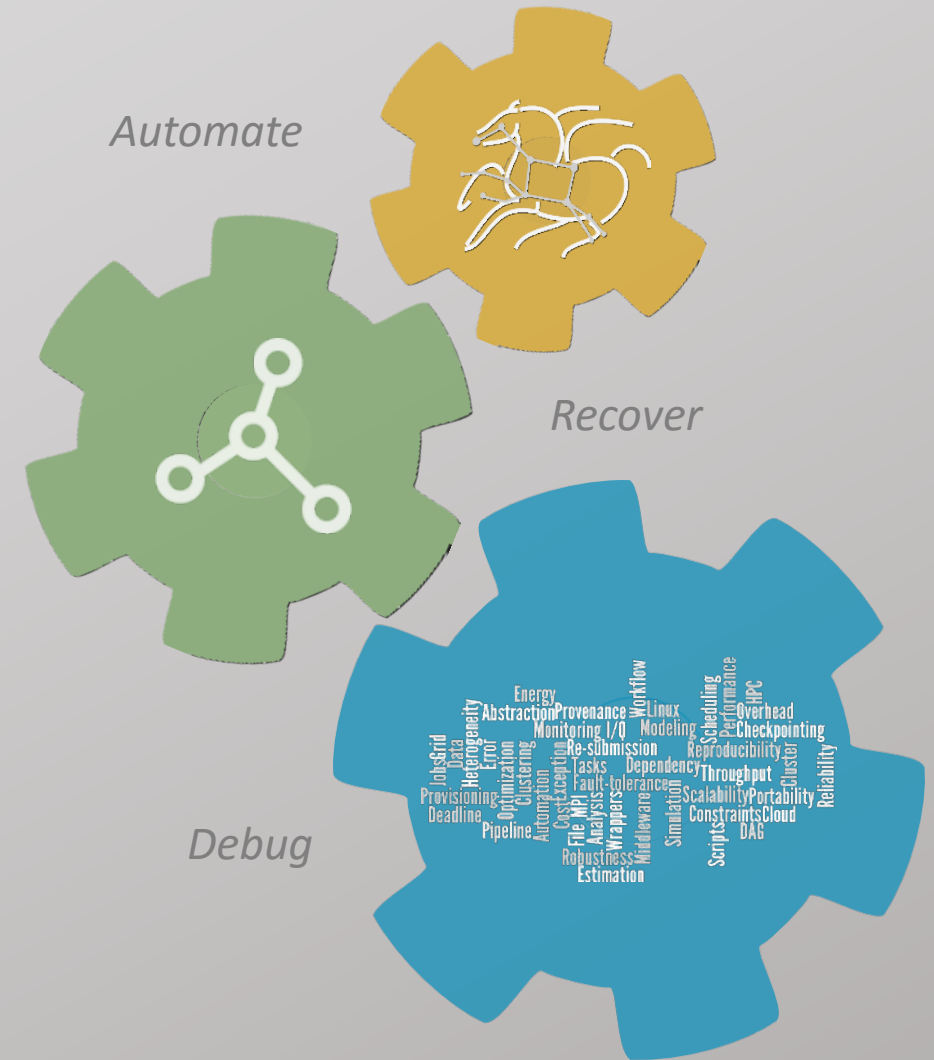Enables parallel, distributed computations

Automatically executes data transfers

Reusable, aids reproducibility

Records how data was produced (provenance)

Handles failures with to provide reliability
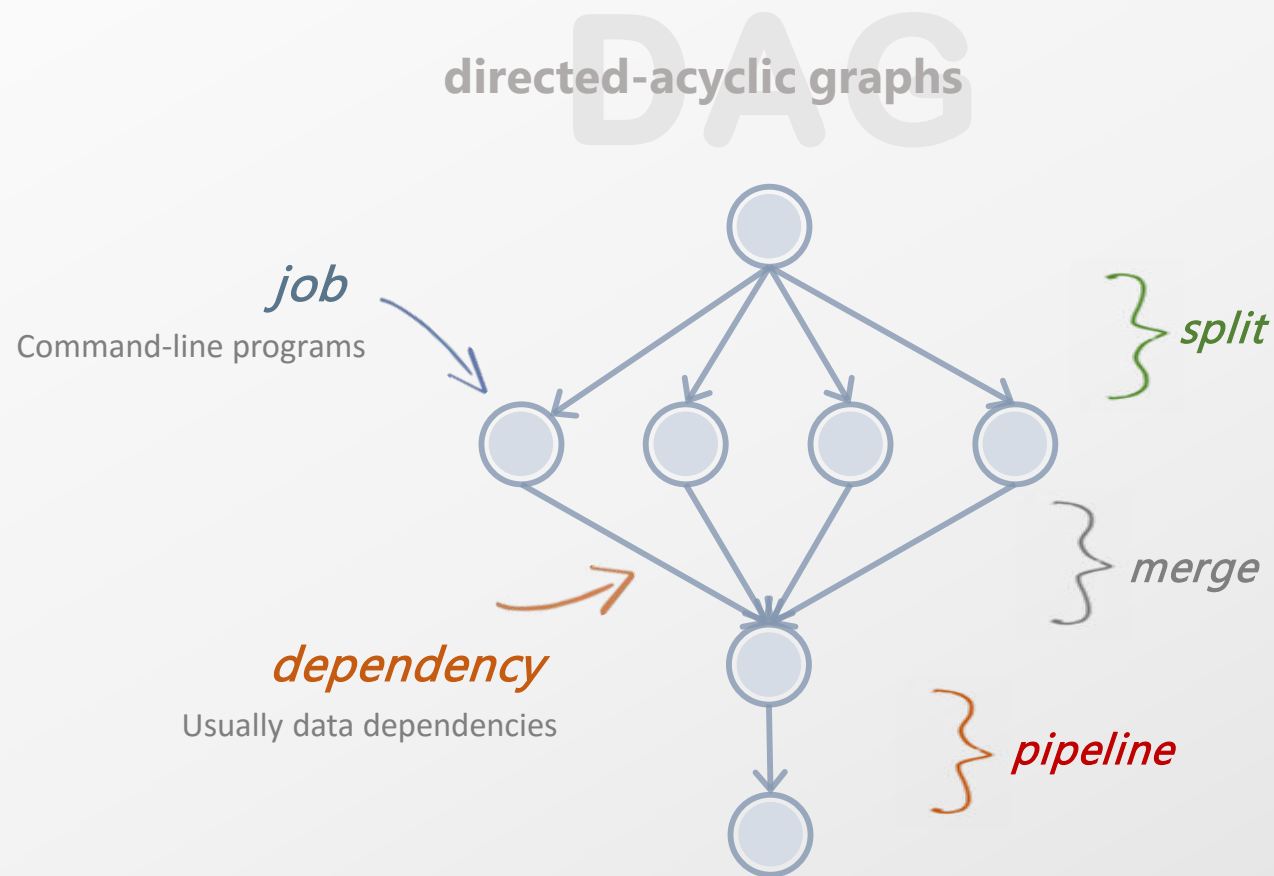
Keeps track of data and files

*Automate*

*Recover*

*Debug*

**HTC**ondor
**High Throughput Computing**

Pegasus

# Taking a closer look into a workflow...

**directed-acyclic graphs**

DAG

*job*

Command-line programs

*split*

*merge*

*dependency*

Usually data dependencies

*pipeline*

**DAG in XML**

DAX

# 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*

# 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()
```
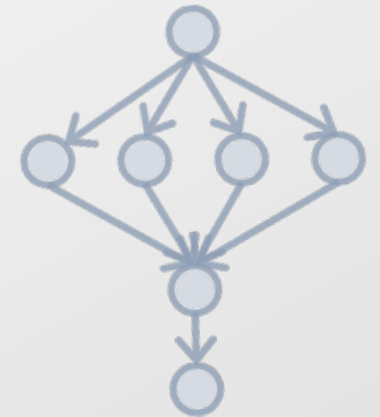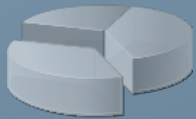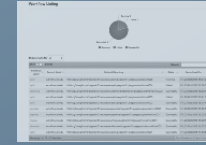
DAG in XML

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.

### Statistics

Workflow execution and job performance metrics

### Debug

Set of debugging tools to unveil issues

### RESTful API

Monitoring and reporting information on your own application interface

**Past executions?**

### Command-line tools

Tools for monitor and debug workflows

**Pegasus**

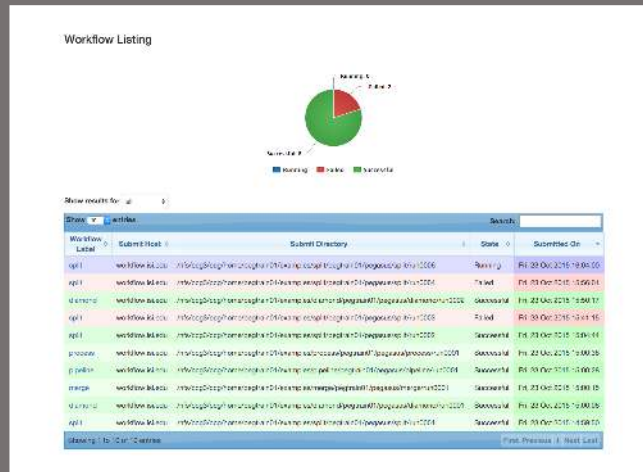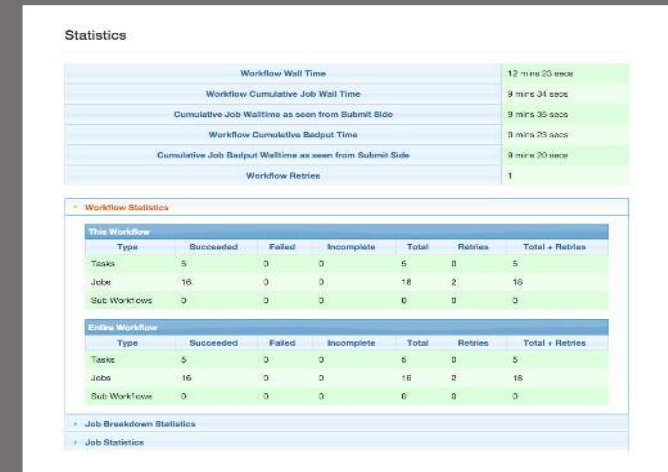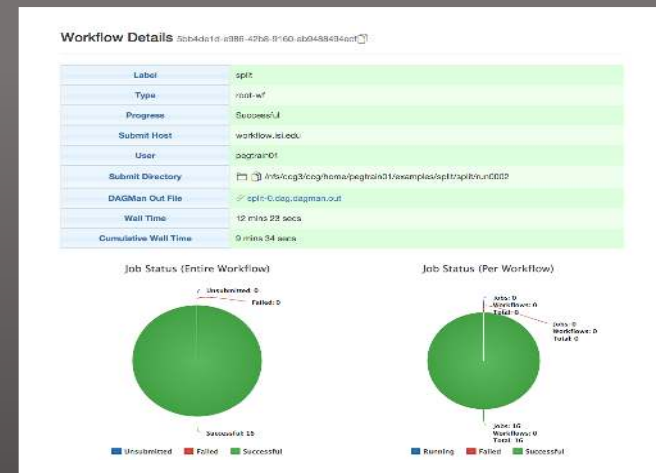*http://pegasus.isi.edu*

**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**

# 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?

Shared filesystem

Let Pegasus manage it for you
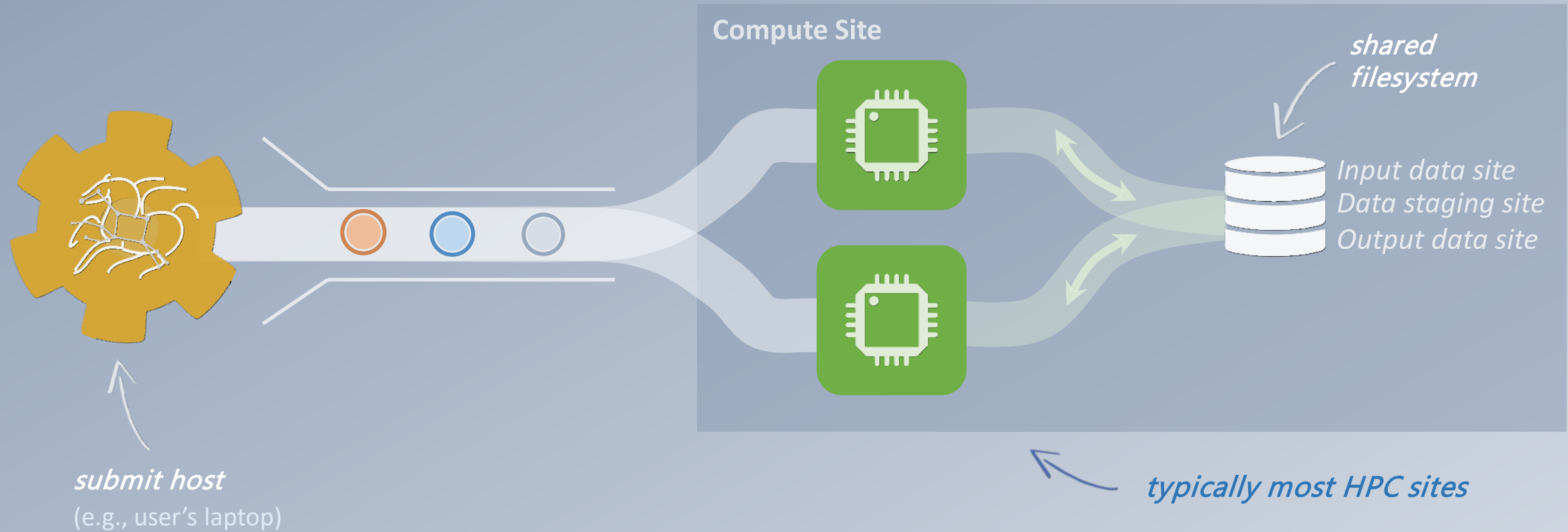
ftp

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

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

# 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**

# Pegasus can also manage data over the submit host...



Compute Site

Typical OSG sites
Open Science Grid

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

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

Compute site B

Data staging site

Input data site
Output data site

Pegasus

# 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**

*http://pegasus.isi.edu*

# 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

# A few more features...

# Performance, why not improve it?

*clustered job*

Groups small jobs together
to improve performance

*task*

small granularity

*http://pegasus.isi.edu*

**Pegasus**

22

# 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 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)

*workflow wrapped as an MPI job*

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

**HPC System**

*Master (rank 0)*

*worker*

*rank 1*

*rank n-1*

**Pegasus**

*http://pegasus.isi.edu*

# 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

**Remote Workflow
Engine**
Site Catalog
Transformation Catalog

**Pegasus**

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

18 million input images (~2.5TB)

1,100 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

## 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

**Pegasus**

# http://soykb.org

**XSEDE Allocation**
**PI: Dong Xu**
**Trupti Joshi, Saad Kahn, Yang Liu, Juexin Wang, Badu Valliyodan, Jiaojiao Wang**

https://github.com/pegasus-isi/Soybean-Workflow

| Task | Base Code | Cores (Threads) | Memory (GB) |
|------|-----------|-----------------|-------------|
| Alignment_to_reference | BWA | 7 | 8 |
| Sort_sam | Picard | 1 | 21 |
| Dedup | Picard | 1 | 21 |
| Add_replace | Picard | 1 | 21 |
| Realign_target_creator | GATK | 15 | 10 |
| Indel_realign | GATK | 1 | 10 |
| Haplotype_caller | GATK | 1 | 3 |
| Genotype_gvcfs | GATK | 1 | 10 |
| Merge_gvcf | GATK | 10 | 20 |
| Combine_variants | GATK | 1 | 10 |
| Select_variants | GATK | 14 | 10 |
| Filtering | GATK | 1 | 10 |

## 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 in the submit host - Workflow runs at a finer granularity

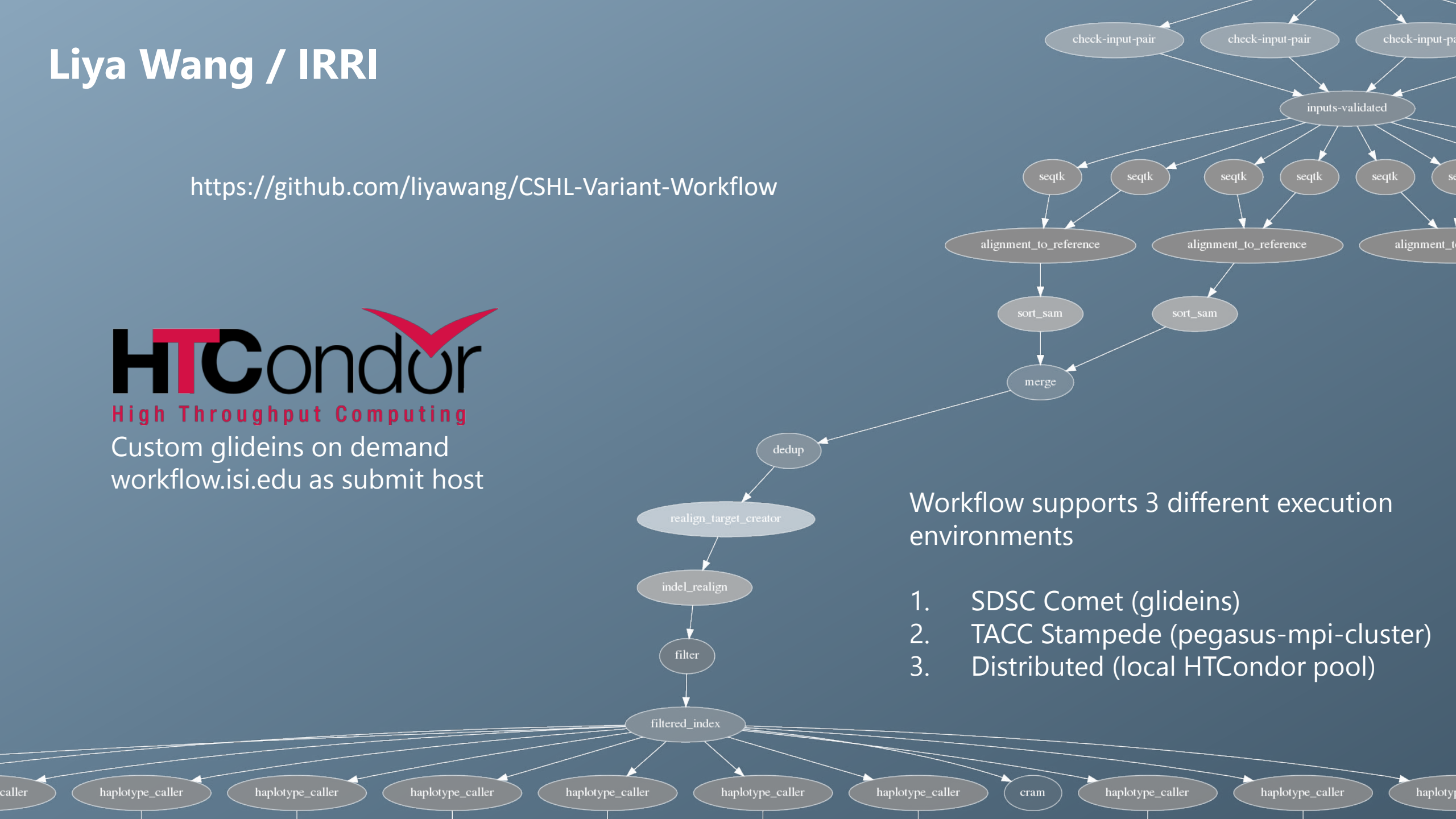Works well on Wrangler due to more cores and memory per node (48 cores, 128 GB RAM)

**Pegasus**

# Liya Wang / IRRI

https://github.com/liyawang/CSHL-Variant-Workflow
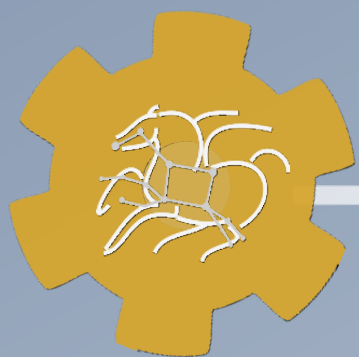
**HTCondor**
**High Throughput Computing**

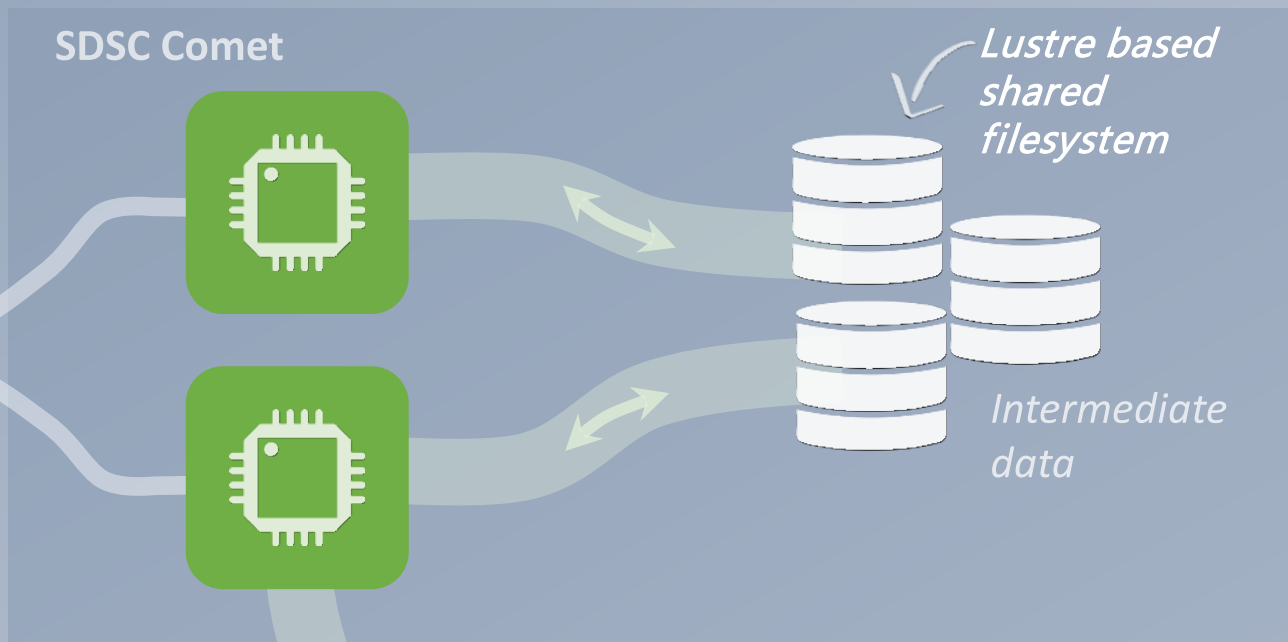Custom glideins on demand
workflow.isi.edu as submit host

Workflow supports 3 different execution environments

1. SDSC Comet (glideins)
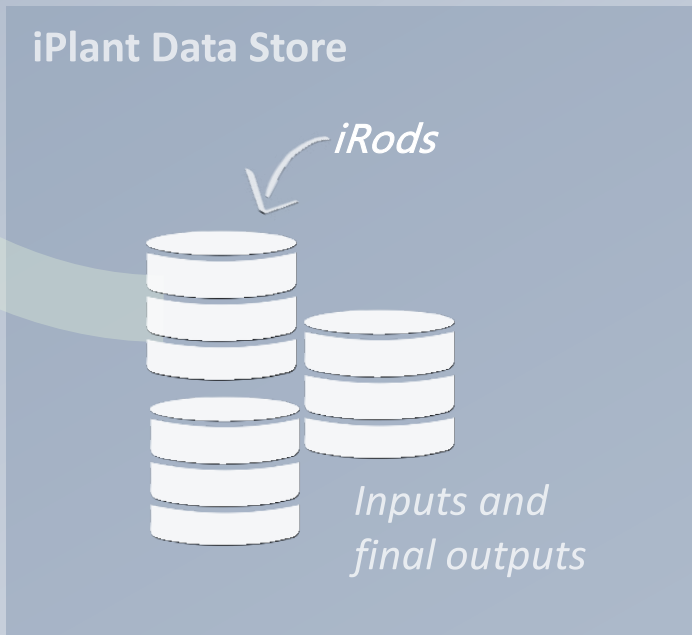2. TACC Stampede (pegasus-mpi-cluster)
3. Distributed (local HTCondor pool)

check-input-pair    check-input-pair    check-input-pa

inputs-validated

seqtk    seqtk    seqtk    seqtk    seqtk    se

alignment_to_reference    alignment_to_reference    alignment_to

sort_sam    sort_sam

merge

dedup

realign_target_creator

indel_realign

filter

filtered_index

caller    haplotype_caller    haplotype_caller    haplotype_caller    haplotype_caller    haplotype_caller    cram    haplotype_caller    haplotype_caller    haploty

**SDSC Comet**

Lustre based shared filesystem

Intermediate data

submit host (workflow.isi.edu)

**iPlant Data Store**

iRods

Inputs and final outputs

**Pegasus**

http://pegasus.isi.edu

31

# Pegasus est. 2001

Automate, recover, and debug scientific computations.

# Get Started

**Pegasus Website**

http://pegasus.isi.edu

**Users Mailing List**

pegasus-users@isi.edu

**Support**

pegasus-support@isi.edu

**HipChat**