

Pegasus short tutorial



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

Pegasus has a number of features that contribute to its useability and effectiveness

- **Portability / Reuse**

User created workflows can easily be run in different environments without alteration. Pegasus currently runs workflows on top of Condor, Grid infrastructures such as Open Science Grid and TeraGrid, Amazon EC2, Nimbus, and many campus clusters. The same workflow can run on a single system or across a heterogeneous set of resources.

- **Performance**

The Pegasus mapper can reorder, group, and prioritize tasks in order to increase the overall workflow performance.

- **Scalability**

Pegasus can easily scale both the size of the workflow, and the resources that the workflow is distributed over. Pegasus runs workflows ranging from just a few computational tasks up to 1 million. The number of resources involved in executing a workflow can scale as needed without any impediments to performance.

- **Provenance**

By default, all jobs in Pegasus are launched via the kickstart process that captures runtime provenance of the job and helps in debugging. The provenance data is collected in a database, and the data can be summarised with tools such as `pegasus-statistics`, `pegasus-plots`, or directly with SQL queries.

- **Data Management**

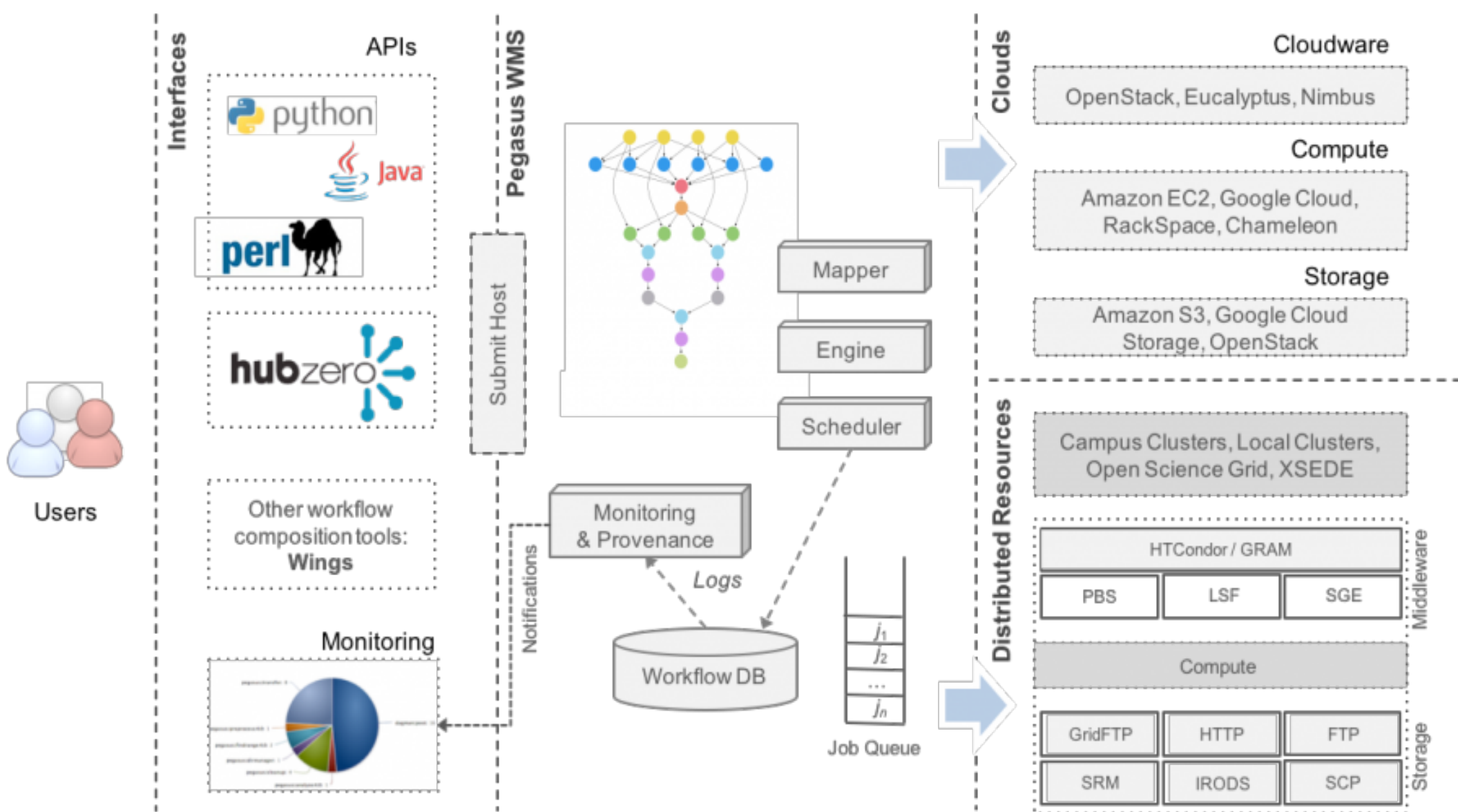
Pegasus handles replica selection, data transfers and output registrations in data catalogs. These tasks are added to a workflow as auxiliary jobs by the Pegasus planner.

- **Reliability**

Jobs and data transfers are automatically retried in case of failures. Debugging tools such as pegasus-analyzer helps the user to debug the workflow in case of non-recoverable failure

- **Error Recovery**

When errors occur, Pegasus tries to recover when possible by retrying tasks, by retrying the entire workflow, by providing workflow-level checkpointing, by re-mapping portions of the workflow, by trying alternative data sources for staging data, and, when all else fails, by providing a rescue workflow containing a description of only the work that remains to be done. It cleans up storage as the workflow is executed so that data-intensive workflows have enough space to execute on storage-constrained resource. Pegasus keeps track of what has been done (provenance) including the locations of data used and produced, and which software was used with which parameters.



- Gallery of workflows:
https://pegasus.isi.edu/workflow_gallery/
- Virtual box image (tutorial):
demo...

```

tutorial@localhost [tutorial] 208 Jan
[tutorial@localhost split]$ history
 1 ll
 2 mkdir test1
 3 cd test1
 4 pegasus-init split
 5 ll
 6 ll ./split
 7 cd split/
 8 ./generate_dax.sh split.dax
 9 ls -lt
10 ./plan_dax.sh split.dax
11 condor_status
12 condor_status
13 hostname

```

Workflow Listing

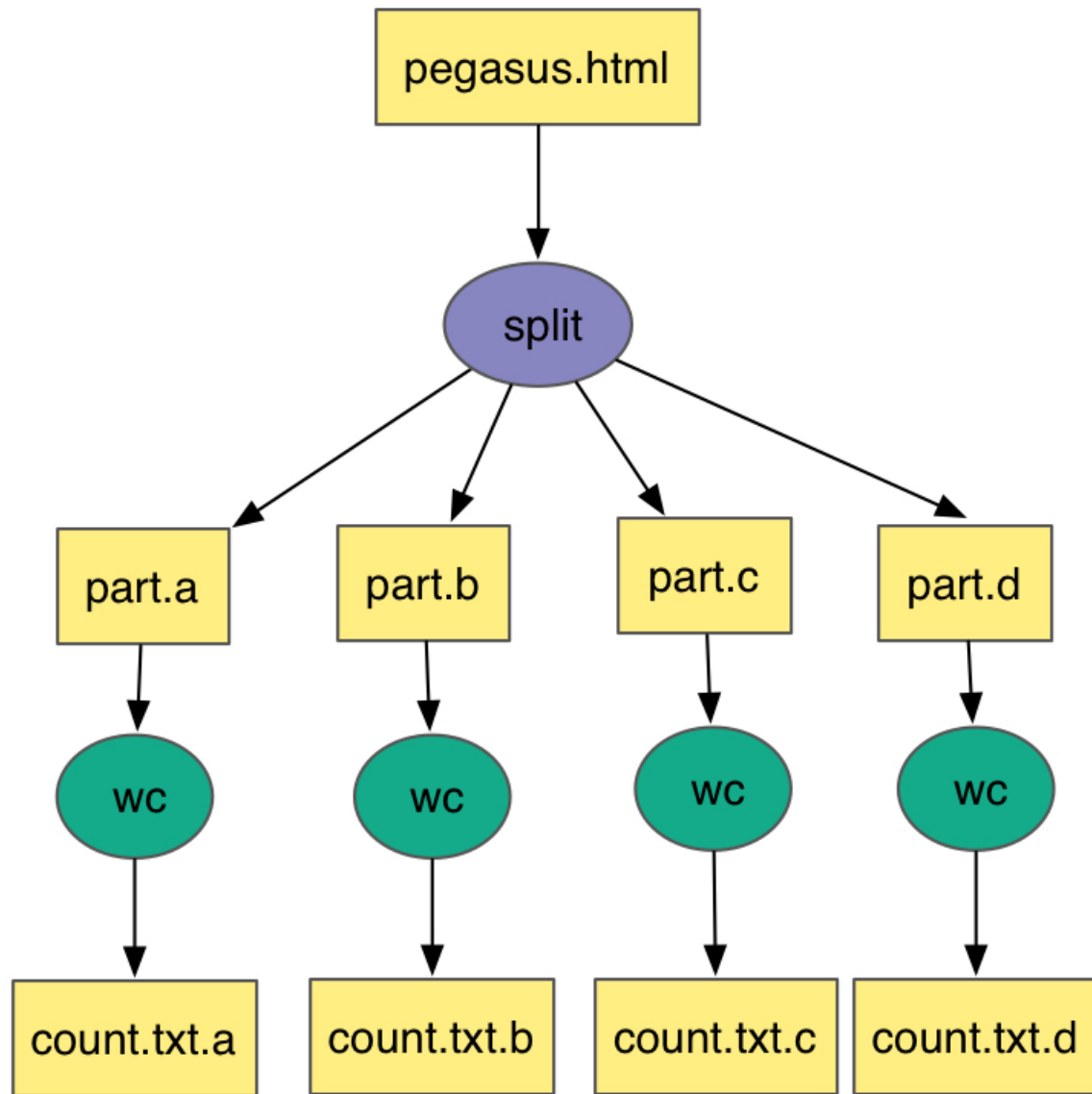
Running: 0
Failed: 0
Successful: 1

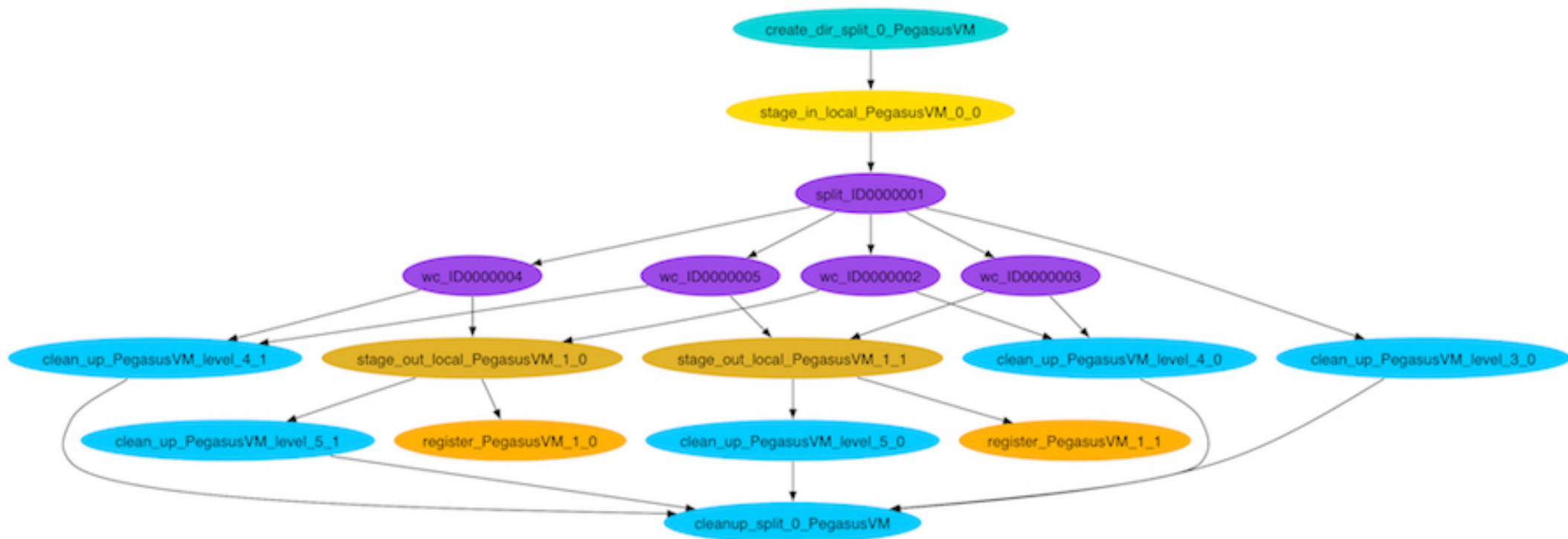
Show results for

Workflow Label	Submit Host	Submit Directory	State	Submitted On
split	localhost	/home/tutorial/test1/split/submit/tutorial/pegasus/split/run0001	Successful	Sat, 11 Jan 2020 09:57:42

Showing 1 to 1 of 1 entries

STAMPEDE INFORMATION SCIENCES INSTITUTE USC







Pegasus

Workflow Management
System