



Managing Workflows Within HUBzero: How to Use Pegasus to Execute Computational Pipelines

Ewa Deelman

USC Information Sciences Institute

Acknowledgement:

Steven Clark, Derrick Kearney, Michael McLennan (HUBzero)

Frank McKenna (OpenSees)

Gideon Juve, Gaurang Mehta, Mats Rynge, Karan Vahi (Pegasus)



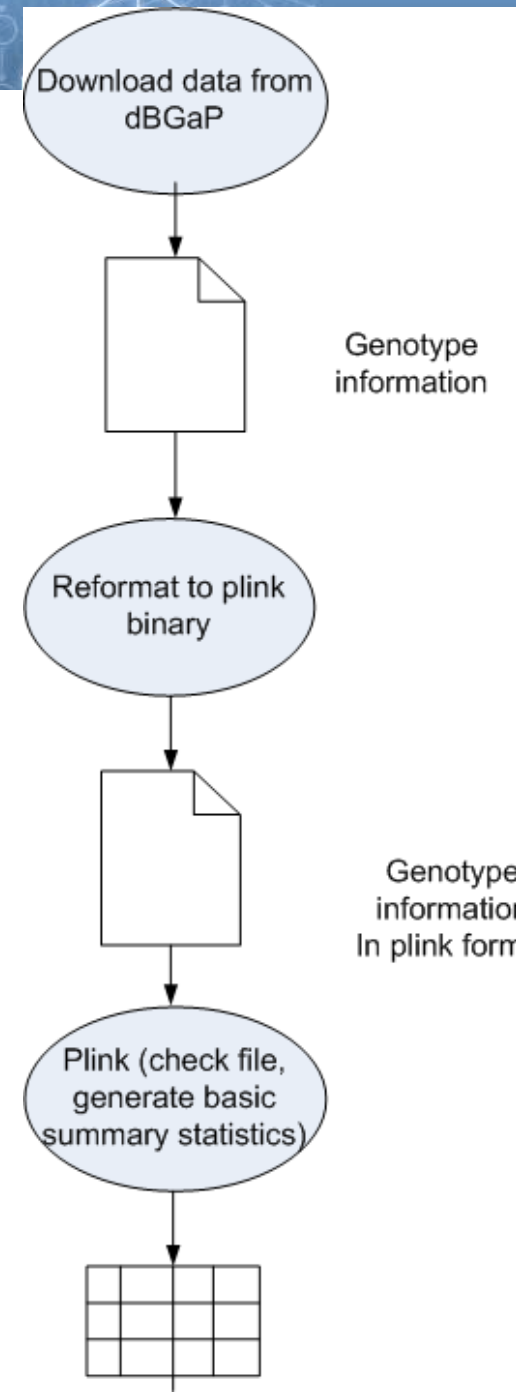
Outline

- Introduction to Pegasus and workflows
- HUB Integration
 - Rappture and Pegasus
 - Submit command and Pegasus
- Example: OpenSEES / NEESHub
- Future directions



Computational workflows

- Help express multi-step computations in a declarative way
- Can support automation, minimize human involvement
 - Makes analyses easier to run
- Can be high-level and portable across execution platforms
- Keep track of provenance to support reproducibility
- Foster collaboration—code and data sharing





Workflow Management

- You may want to use different resources within a workflow or over time
 - Need a high-level workflow specification
 - Need a **planning** capability to map from high-level to executable workflow
 - Need to **manage the task dependencies**
 - Need to **manage the execution of tasks** on the remote resources
- Need to provide scalability, performance, reliability



Our Approach

● Analysis Representation

- Support a declarative representation for the workflow (dataflow)
- Represent the workflow structure as a Directed Acyclic Graph (DAG) in a resource-independent way
- Use recursion to achieve scalability

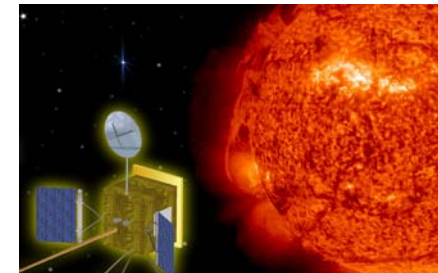
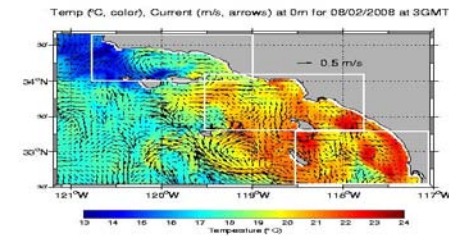
● System (Plan for the resources, Execute the Plan, Manage tasks)

- Layered architecture, each layer is responsible for a particular function (Pegasus Planner, DAGMan, Condor schedd)
- Mask errors at different levels of the system
- Modular, composed of well-defined components, where different components can be swapped in
- Use and adapt existing graph and other relevant algorithms
- Can be embedded into



Pegasus Workflow Management System (est. 2001)

- A collaboration with University of Wisconsin Madison
- Used by a number of applications in a **variety of domains**
- **Provides reliability**—can retry computations from the point of failure
- **Provides scalability**—can handle large data and many computations (kbytes-TB of data, $1-10^6$ tasks)
- **Optimizes workflows for performance**
- Automatically captures **provenance** information
- Runs workflows on distributed resources: laptop, campus cluster, Grids (DiaGrid, OSG, XSEDE), Clouds (FutureGrid, EC2, etc..)





Planning Process

- Assume data may be distributed in the Environment
- Assume you may want to use local and/or remote resources
- Pegasus needs information about the environment
 - data, executables, execution and data storage sites
- Pegasus generates an executable workflow
- Data transfer protocols
 - Gridftp, Condor I/O, HTTP, scp, S3, iRods, SRM, FDT (partial)
- Scheduling to interfaces
 - Local, Gram, Condor, Condor-C (for remote Condor pools), via Condor Glideins – PBS, LSF, SGE

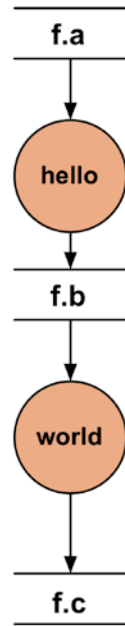


Generating executable workflows

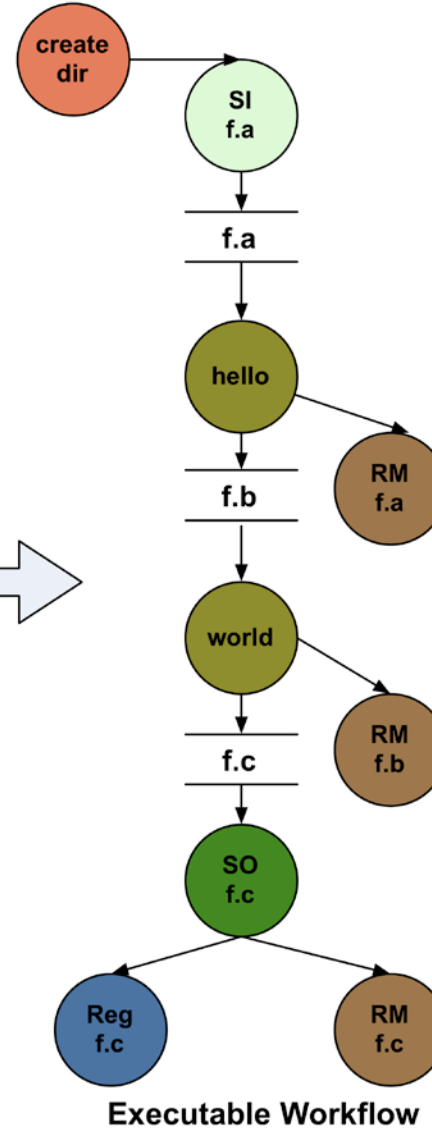
APIs for
workflow
specification
(DAX---
DAG in XML)

Java, Perl, Python

(DAX)



Abstract Workflow



Executable Workflow

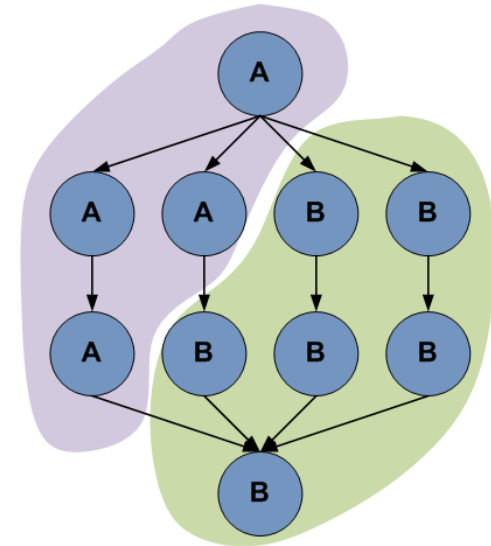
LEGEND

- Unmapped Job
- Compute Job mapped to a site
- Stage-in Job
- Stage-Out Job
- Registration Job
- Make Dir Job
- Cleanup Job



Advanced features

- Performs data reuse
- Registers data in data catalogs
- Manages storage—deletes data no longer needed
- Can cluster tasks together for performance
- Can manage complex data architectures (shared and non-shared filesystem, distributed data sources)
- Different execution modes which leverage different computing architectures (Condor pools, HPC resources, etc..)





HUBzero Integration

Pegasus with



Pegasus Tutorial

Launch Tool

Version 1.0 - published on 22 Sep 2012

Open source: [license](#) | [download](#)

[View All Supporting Documents](#)

- 0 Citation(s)
- 0 review(s) ([Review this](#))
- 0 questions ([Ask a question](#))
- user(s), [detailed usage](#)
- 0 wish(es) ([Add a new wish](#))
- [Add to your favorites!](#)

Share: [f](#) [t](#) [g+](#) ...

By [Steven Clark](#)¹, [Derrick Kearney](#), [Mats Ryngé](#)², [Karan Vahi](#)

1. [Purdue University](#) 2. [USC Information Sciences Institute](#)

Demonstration of simple Pegasus examples combined with Rapture interface.

About Citations Reviews Questions Usage Wish List Supporting Docs

SEE ALSO

No results found.

Category [Tools](#)

Abstract **Goal:** The goal of this tutorial is to demonstrate how to use the [Pegasus Workflow Management System](#) together with a Rapture interface to present the users an easy way to run complex computations in the Hub.

The workflow in this tutorial is a simple Hello World example. Attached to this tool is a PDF which shows how the tool was created.

To see the code for the Hub integration, please use a Workspace and check out a copy:



Benefits of Pegasus for HUB Users

- **Provides Support for Complex Computations**
 - Can connect the existing HUB models into larger computations
- **Portability / Reuse**
 - User created workflows can easily be run in different environments without alteration (today DiaGrid, OSG)
- **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.



Benefits of Pegasus for HUB Users

- **Provenance**

- Performance and provenance data is collected in a database, and the data can be summaries with tools such as **pegasus-statistics**, **pegasus-plots**, or directly with SQL queries.

- **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 failures.

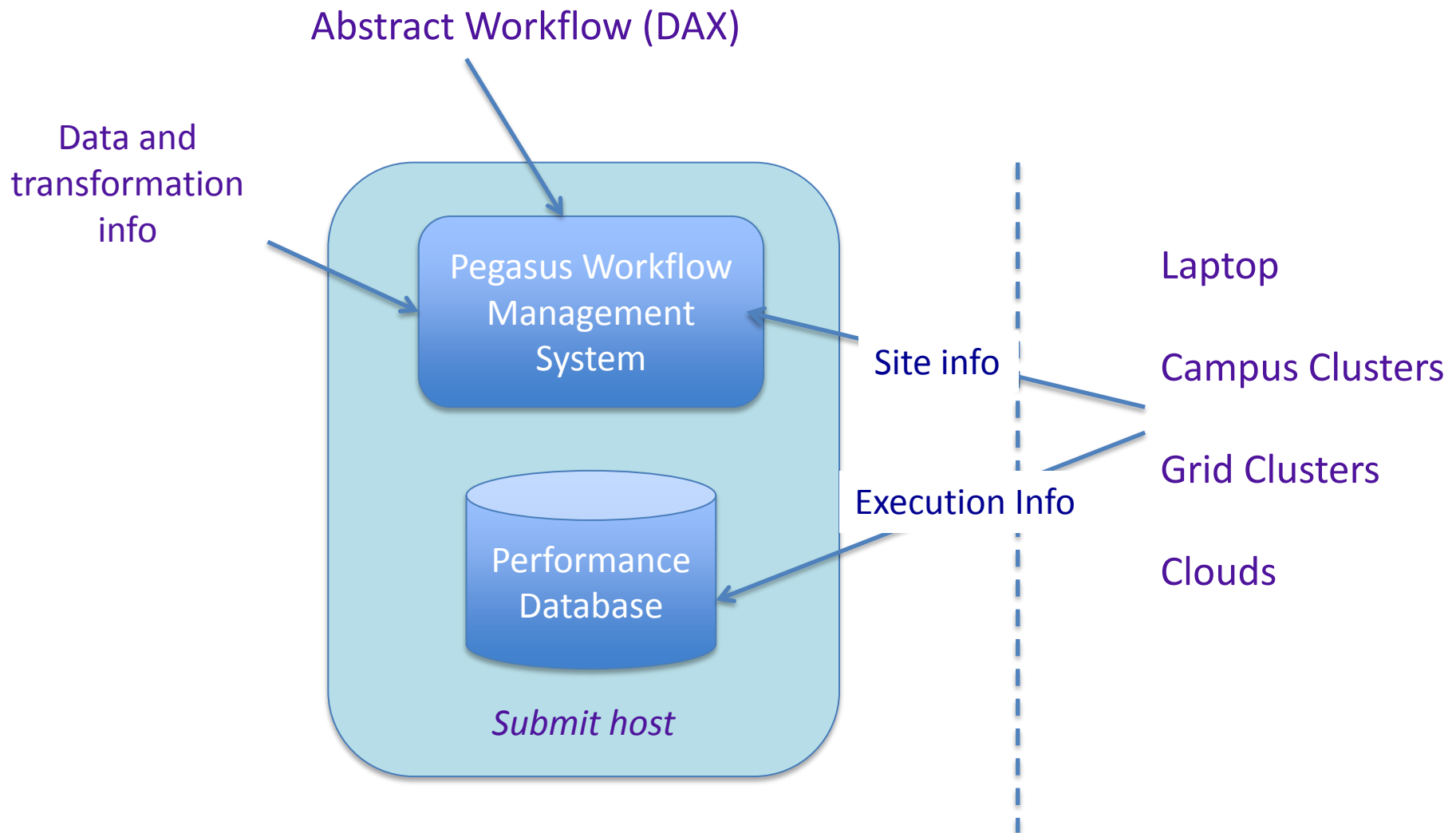


Pegasus in HUBzero

- Pegasus as a backend to the *submit* command
- Pegasus workflows composed in *Rappture*
 - *Build workflow within Rappture*
 - *Have Rappture collect inputs, call a workflow generator, and collect outputs*
- Pegasus Tutorial tool now available in HUBzero

<http://hubzero.org/tools/pegtut>

- Session that includes Pegasus on Tuesday 1:30 – 5:30
Room 206 #2 **Creating and Deploying Scientific Tools (part 2)**
“... Scientific Workflows with Pegasus” by George Howlett & Derrick Kearney,
Purdue University



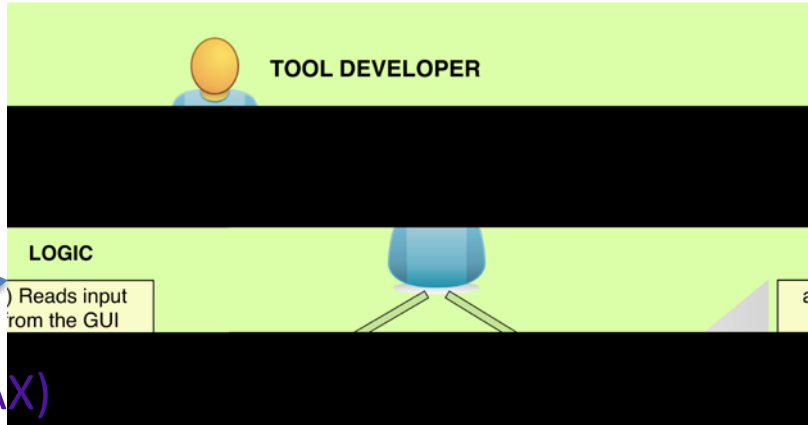


Use of Pegasus with Submit Command

- Used by Rappture interface to submit the workflow
- Submits the workflow through Pegasus to
 - OSG
 - DIAGRID
- Prepares the site catalog and other configuration files for Pegasus
- Uses pegasus-status to track the workflow
- Generates statistics and report about job failures using

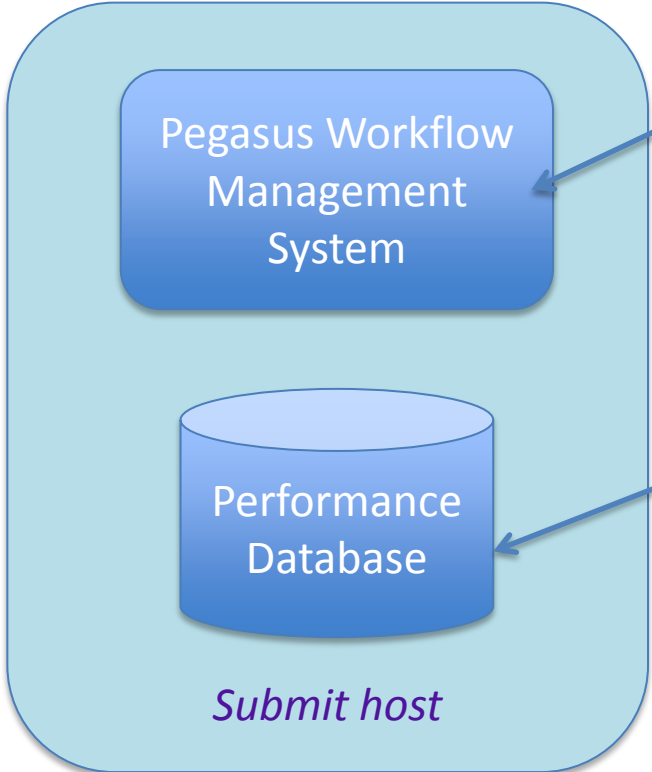


Data and transformation info



Hub

Abstract Workflow (DAX)



Site info

- Laptop
- Campus Clusters
- Grid Clusters
- Clouds

Execution Info



Pegasus Workflows in the HUB



Tool description

Inputs

Outputs

```
?xml version="1.0"?  
<run>  
  <tool>  
    <title>Hello World Workflow</title>  
    <about>Simple interface for Pegasus Hello World example</about>  
    <command>python @tool/wrapper.py @driver</command>  
  </tool>  
  <input>  
    <note>  
      <contents>file://workflow.html</contents>  
    </note>  
    <string id="textinput">  
      <about>  
        <label>Your Name</label>  
        <description>Enter your name</description>  
      </about>  
      <default>Pete</default>  
    </string>  
  </input>  
  <output>  
    <string id="greeting">  
      <about>  
        <label>Greeting</label>  
        <description>Greeting generated by the two jobs submitted to the grid through</description>  
      </about>  
    </string>  
    <string id="fa">  
      <about>  
        <label>f,a</label>  
        <description>Input to Hello job</description>  
      </about>  
    </string>  
    <string id="fb">  
      <about>  
        <label>f,b</label>  
        <description>Output from sayhi job</description>  
      </about>  
    </string>  
  </output>  
</run>
```

Rappture (data definitions)

Calls an external DAX generator

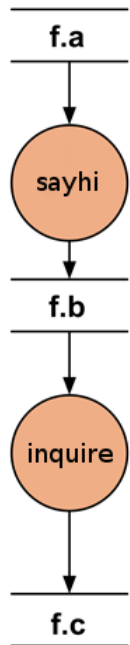
wrapper.py

- Python script
- Collects the data from the Rappture interface
- Generates the DAX
- Runs the workflow
- Presents the outputs to Rappture

```
#####  
# MAIN PROGRAM - generated by the Rappture Builder  
#####  
import sys  
import os  
from math import *  
  
import Rappture  
from Rappture.tools import getCommandOutput as RapptureExec  
from Pegasus.DAX3 import *  
  
# open the XML file containing the run parameters  
io = Rappture.library(sys.argv[1])  
  
# setup paths to our executables  
scriptpath = os.path.realpath(__file__)  
scriptdir = os.path.dirname(scriptpath)  
tooldir = os.path.dirname(scriptdir)  
sayhipath = os.path.join(tooldir, 'bin', 'sayhi.sh')  
inquirepath = os.path.join(tooldir, 'bin', 'inquire.sh')  
  
#####  
# Get input values from Rappture  
#####  
  
# get input value for input.string(textinput)  
textinput = io.get('input.string(textinput),current')  
if not textinput:  
    sys.stderr.write("Input data is missing\n")  
    sys.exit(1)  
  
#####  
# Add your code here for the main body of your program  
#####  
  
fp = open('f.a', 'w')  
if fp:  
    fp.write(textinput + '\n')  
    fp.close()  
else:  
    sys.stderr.write("Could not create datafile\n")
```



Workflow generation



Abstract Workflow

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

# Add input file to the DAX-level replica catalog
a = File("f.a")
a.addPFN(PFN("file://" + os.getcwd() + "/f.a", "local"))
dax.addFile(a)

# Add executables to the DAX-level replica catalog
e_sayhi = Executable(namespace="sayhi_inquire", name="sayhi", version="1.0", \
                    os="linux", arch="x86_64", installed=False)
e_sayhi.addPFN(PFN("file://" + sayhipath, "condorpool"))
dax.addExecutable(e_sayhi)

e_inquire = Executable(namespace="sayhi_inquire", name="inquire", version="1.0", \
                    os="linux", arch="x86_64", installed=False)
e_inquire.addPFN(PFN("file://" + inquirepath, "condorpool"))
dax.addExecutable(e_inquire)

# Add the sayhi job
sayhi = Job(namespace="sayhi_inquire", name="sayhi", version="1.0")
sayhi.addArguments('f.a')
b = File("f.b")
sayhi.uses(a, link=Link.INPUT)
sayhi.uses(b, link=Link.OUTPUT)
dax.addJob(sayhi)

# Add the inquire job (depends on the sayhi job)
inquire = Job(namespace="sayhi_inquire", name="inquire", version="1.0")
inquire.addArguments('f.b')
c = File("f.c")
inquire.uses(b, link=Link.INPUT)
inquire.uses(c, link=Link.OUTPUT)
```



User provides inputs to the workflow and clicks the “Submit” button

Hello World with Pegasus

With Pegasus, you can hook together multiple executables and pass data between them. In this example, two executables are used to build a greeting.

To start, enter your name below. After pressing the Simulate button, the information is saved in the file named f.a, and passed to the **sayhi** program. **sayhi** prints a specialized hello message for you. The greeting is saved in the file named f.b and sent to the **inquire** program. **inquire** finishes the greeting by asking how you are doing.

This job is executed on the grid using Pegasus and Submit!

Learn about software that powers this tool:

- [Pegasus Workflow Management System](#)
- [Rapture Toolkit](#)
- [Submit](#)

Give it a try:

Enter Your Name Here **Then** Press Simulate Up There

Your Name:

Simulate new input parameters

About this tool Questions?

Hello World Workflow

Simple interface for Pegasus Hello World example



Workflow has completed. Outputs are available for browsing/downloading

Hello World with Pegasus

With Pegasus, you can hook together multiple executables and pass data between them. In this example, two executables are used to build a greeting.

To start, enter your name below. After pressing the Simulate button, the information is saved in the file named f.a, and passed to the **sayhi** program. **sayhi** prints a specialized hello message for you. The greeting is saved in the file named f.b and sent to the **inquire** program. **inquire** finishes the greeting by asking how you are doing.

This job is executed on the grid using Pegasus and Submit!

Learn about software that powers this tool:

- [Pegasus Workflow Management System](#)
- [Rapture Toolkit](#)
- [Submit](#)

Give it a try:

Enter Your Name Here **Then** Press Simulate Up There

Your Name:

? About this tool Questions?

Result: f.c

Hello

- Greeting
- f.a
- f.b
- f.c**
- Analysis
- Run metrics
-
- Download

Find:

1 result



OpenSEES / NEEShub



The OpenSeesLab tool:

<http://nees.org/resources/tools/openseeslab>

OpenSeesLab NEEShub

A Collection of Tools for Structural/Geotechnical Engineers that use the **Open** System for **E**arthquake **E**ngineering **S**imulation

Response SOOF to Earthquake

$$m(\dot{u}) = c\dot{u} + k(u - m\ddot{u}_g)$$
$$\dot{u} = \dot{u}_g - \dot{u}_g$$
$$T = \frac{2\beta - 2\beta_1}{\beta - \beta_1} \frac{c}{\omega} \frac{1 + \beta_1 \omega^2}{1 + \beta \omega^2}$$

Damping Ratio: $\zeta = \frac{c}{2\beta\omega}$

Analysis Type: Single Analysis

Source Ground Motion: PEER NGA

Earthquake Name: []

Min Mag of EQ to search: 0.0

Max Mag of EQ to search: 6

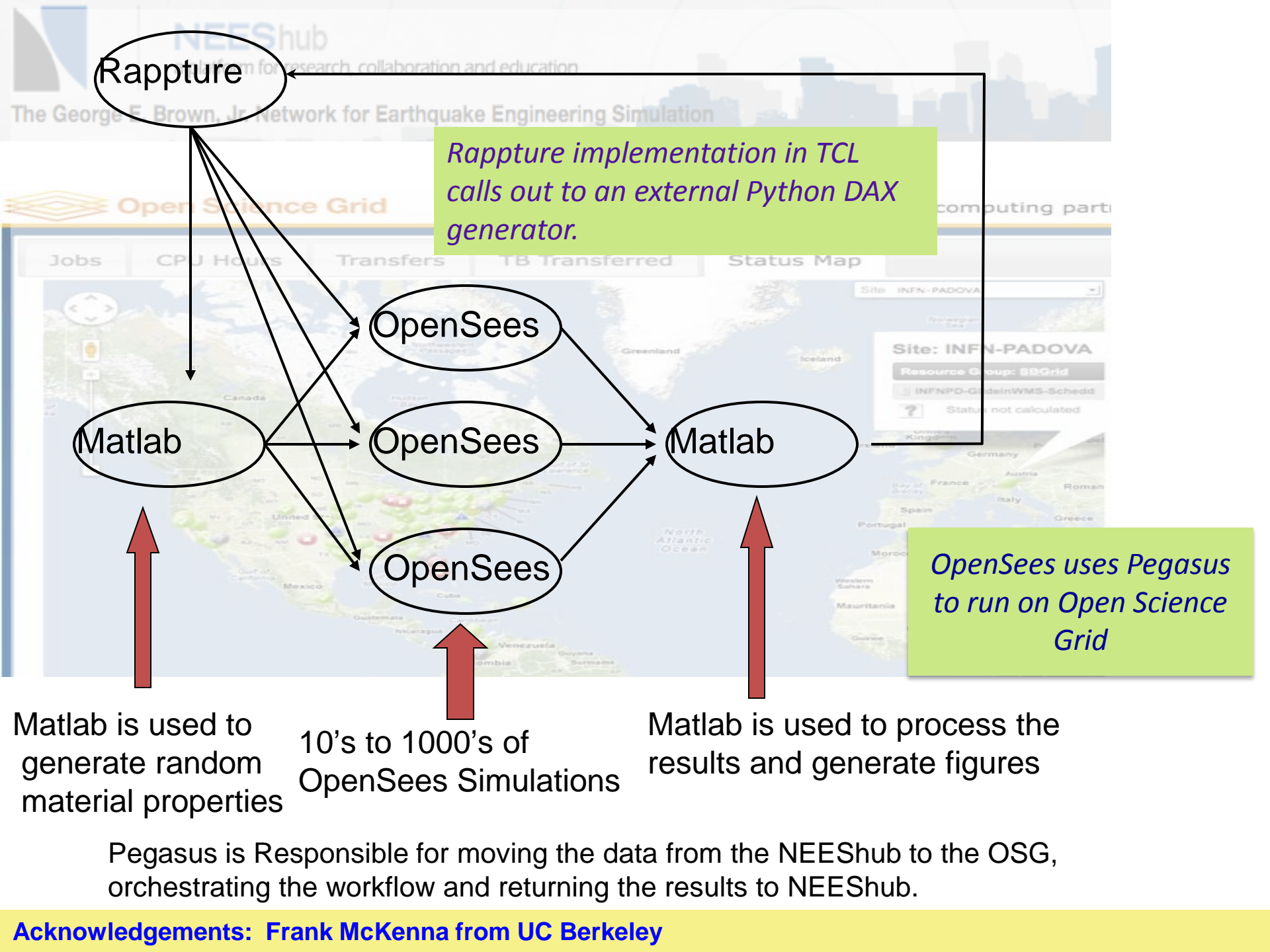
EQ Type: []

Min PGA of records to search: 0.12

Max PGA of record to search: 1.0

Is a suite of Simulation Tools powered by OpenSees for:

1. Submitting OpenSees scripts to NEEShub resources
2. Educating students and practicing engineers

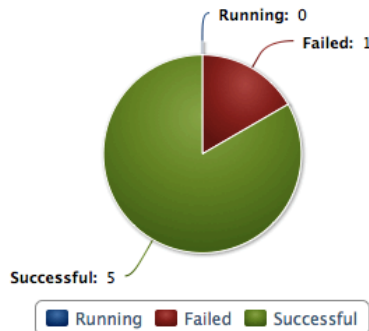




Future Directions

- Submit to manage parameter sweep computations (now only on HUBzer0)
- Web-based monitoring

Workflow Listing



Label	Type	Progress	Submit Host	User	Submit Directory
hello_world	root-wf	Failed	mayani.isi.edu	mayani	/ifs1/software/pegasus/trunk/dist/pegasus-4.2.0cvs/share/pegasus/examples/hello-world/work/mayani/pegasus/hello_world/20120919T180724-0700

Wall Time	Cumulative Wall Time	Successful Jobs	Failed Jobs	Total Jobs
2 mins 47 secs	5 secs	7	1	11

Charts

Sub Workflows	Failed	Running	Successful
---------------	--------	---------	------------

Show 10 entries

Workflow Label	Submit Directory
hello_world	/ifs1/software/pegasus/trunk/dist/pegasus-4.2.0cvs/share/pegasus/examples/hello-world/work/mayani/pegasus/hello_world/20120919T180724-0700
hello_world	/ifs1/software/pegasus/trunk/dist/pegasus-4.2.0cvs/share/pegasus/examples/hello-world/work/mayani/pegasus/hello_world/20120919T180434-0700

Show 10 entries Search:

Job Name	Exit Code	Standard Out	Standard Error
world_ID0000002	256	stdout	stderr

Showing 1 to 1 of 1 entries First Previous 1 Next Last



Benefits of workflows in the HUB

- Support for complex applications/ builds on existing domain tools
- Clean separations for users/developers/operator
 - User: Nice high level interface via Rappture
 - Tool developer: Only has to build/provide a description of the workflow (DAX)
 - Hub operator: Ties the Hub to an existing distributed computing infrastructure (DiaGrid, OSG, ...)
- The Hub and Pegasus handle low level details
 - Job scheduling to various execution environments
 - Data staging in a distributed environment
 - Job retries
 - Workflow analysis
 - Support for large workflows



Benefits of the HUB to Pegasus

- Provides a nice, easy to use interface to Pegasus workflows
- Broadens the user base
- Improves the software based on user's feedback
- Drives innovation—new deployment scenarios, use cases
- I look forward to a continued collaboration



Further Information

- Session that includes Pegasus on Tuesday 1:30 – 5:30
 - Room 206 #2 **Creating and Deploying Scientific Tools (part 2)**
 - “... Scientific Workflows with Pegasus” by George Howlett & Derrick Kearney, Purdue University
- Pegasus Tutorial on the HUB
 - <https://hubzero.org/tools/pegtut>
- General Pegasus Information <http://pegasus.isi.edu>
- Pegasus in a VM—allows you to develop DAXes
 - <http://pegasus.isi.edu/downloads>
 - **We are happy to help!**
- Support mailing lists pegasus-support@isi.edu pegasus-users@isi.edu, pegasus-announce@isi.edu
- Contact me deelman@isi.edu

Big Thank You to the HUBzero and OpenSees teams!