

# Complex Workloads on HUBzero – Pegasus Workflow Management System

Karan Vahi

Science Automation Technologies Group USC Information Sciences Institute



Information Sciences Institute

## **HubZero**

### A valuable platform for scientific researchers

- For building analysis tools and sharing with researchers and educators.
- Made available to the community via a web browser

### Supports interfaces for

- Designing Analysis tools using the Rappture Toolkit
- Uploading and creating inputs
- Visualizing and plotting generated outputs
- Supports hundreds of analysis tools and thousands of users.





2

## **Hubzero - Scalability**

- Execution of the analysis tools for all users cannot be managed on the HubZero instance
- Need to decouple the analysis composition and user interaction layer from backend execution resources
- Scalability requires a need to support multiple types of execution backends
  - Local Campus Cluster
  - DiaGrid
  - Distributed Computational Grids such as Open Science Grid
  - Computational Clouds like Amazon EC2





## **Distributing Analysis - Challenges**

- Portability
  - Some Hubs are tied to local clusters. Others are connected to distributed computational grids. How do we get the analysis tool to run on local PBS cluster one day and OSG the next, or run across them.

### Data Management

- How do you ship in the small/large amounts data required by the analysis tool?
- You upload inputs via the web browser, but the analysis runs on a node in a cluster.
- Different protocols for different sites: Can I use SRM? How about GridFTP? HTTP and Squid proxies?

### Debug and Monitor Computations

- Users need automated tools to go through the log files
- Need to correlate data across lots of log files
- Need to know what host a job ran on and how it was invoked
- Restructure Analysis Steps for Improved Performance
  - Short running tasks or tightly coupled tasks
    - Run on local cluster a hub is connected to.
  - Data placement?





## HubZero – Separation of concerns

- Focus on user interface and provide users
  - means to design, launch analysis steps and inspect and visualize outputs
- Model analysis tools as scientific workflows
- Use a Workflow Management System to manage computation across varied execution resources.





## **Scientific Workflows**

School of Engineering

- Orchestrate complex, multi-stage scientific computations
- Often expressed as directed acyclic graphs (DAGs)
- Capture analysis pipelines for sharing and reuse
- Can execute in parallel on distributed resources



## Why Scientific Workflows?

- Automate complex processing pipelines
- Support parallel, distributed computations
- Use existing codes, no rewrites
- Relatively simple to construct
- Reusable, aid reproducibility
- Can be shared with others
- Capture provenance of data





## **Pegasus Workflow Management System (WMS)**

- Under development since 2001
- A collaboration between USC/ISI and the Condor Team at UW Madison
  - USC/ISI develops Pegasus
  - UW Madison develops DAGMan and Condor
- Maps abstract workflows to diverse computing infrastructure
  - Desktop, Condor Pool, HPC Cluster, Grid, Cloud
- Actively used by many applications in a variety of domains
  - Earth science, physics, astronomy, bioinformatics





## **Benefits of workflows in the Hub**

- 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 Submit 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





## **Pegasus Workflows are Directed Acyclic Graphs**

- Nodes are tasks
  - Typically, executables with arguments.
  - Each executable identified by a unique logical identifier e.g. fft , date, fast\_split
  - Nodes can also be other workflows
- File Aware
  - With each node you specify specify the input and output files referred to by logical identifiers.

### Edges are dependencies

- Represent data flow
- Can also be control dependencies
- Pegasus can infer edges from data use
- No loops, no branches

School of Engineering

- Recursion is possible
- Can generate workflows in a workflow
- Can conditionally skip tasks with wrapper
- Captures computational recipe, devoid of resource descriptions, devoid of data locations, that is portable and can be easily shared.





# **Abstract to Executable Workflow Mapping**

Pegasus compiles the Abstract Workflow to an Executable Workflow that can be executed on varied distributed execution environments **A** 



#### **Abstraction provides**

- Ease of Use (do not need to worry about low-level execution details)
- Portability (can use the same workflow description to run on a number of resources and/or across them)
- Gives opportunities for optimization and fault tolerance
  - automatically restructure the workflow
  - automatically provide fault recovery (retry, choose different resource)

### Pegasus Guarantee -

Wherever and whenever a job runs it's inputs will be in the directory where it is launched.





# Supported Data Staging Approaches - I

## Shared Filesystem setup (typical of XSEDE and HPC sites)

- Worker nodes and the head node have a shared filesystem, usually a parallel filesystem with great I/O characteristics
- Can leverage symlinking against existing datasets
- Staging site is the shared-fs.

School of Engineering

## Non-shared filesystem setup with staging site (typical of OSG and EC 2)

- Worker nodes don't share a filesystem.
- Data is pulled from / pushed to the existing storage element.
- A separate staging site such as S3.

HubZero uses Pegasus to run a single application worklow across sites, leveraging shared filesystem at local PBS cluster and non shared filesystem setup at OSG!







## Supported Data Staging Approaches - II

## Condor IO (Typical of large Condor Pools like CHTC)

- Worker nodes don't share a filesystem
- Symlink against datasets available locally
- Data is pulled from / pushed to the submit host via Condor file transfers
- Staging site is the submit host.



# Supported Transfer Protocols – for directory/file creation and removal, file transfers

- HTTP
- SCP
- GridFTP
- IRODS
- S3 / Google Cloud Storage
- Condor File IO
- File Copy
- OSG Stash

School of Engineering

Using Pegasus allows you to move from one deployment to another without changing the workflow description!

Jobs

Data

### **Pegasus Data Management Tools**

pegasus-transfer, pegasus-create-dir, pegasuscleanup support client discovery, parallel transfers, retries, and many other things to improve transfer performance and reliability



## **Workflow Reduction (Data Reuse)**



# Useful when you have done a part of computation and then realize the need to change the structure.





## File cleanup

Problem: Running out of disk space during workflow execution

### Why does it occur

- Workflows could bring in huge amounts of data
- Data is generated during workflow execution
- Users don't worry about cleaning up after they are done

### Solution

- Do cleanup after workflows finish
  - Add a leaf Cleanup Job
- Interleave cleanup automatically during workflow execution.
  - · Requires an analysis of the workflow to determine, when a file is no longer required
- Cluster the cleanup jobs by level for large workflows
- In 4.6 release, users should be able to specify maximum disk space that should not be exceeded. Pegasus will restructure the workflow accordingly.

# Real Life Example: Used by a UCLA genomics researcher to delete TB's of data automatically for long running workflows!!



## File cleanup (cont)





## **Workflow Restructuring to improve application performance**

- Cluster small running jobs together to achieve better performance
- Why?
  - Each job has scheduling overhead need to make this overhead worthwhile
  - Ideally users should run a job on the grid that takes at least 10/30/60/? minutes to execute
  - Clustered tasks can reuse common input data less data transfers



## **Workflow Monitoring - Stampede**

#### Leverage Stampede Monitoring framework with DB backend

- Populates data at runtime. A background daemon monitors the logs files and populates information about the workflow to a database
- Stores workflow structure, and runtime stats for each task.

#### Tools for querying the monitoring framework

pegasus-status

School of Engineering

- Status of the workflow
- pegasus-statistics
  - · Detailed statistics about your finished workflow
- Integrated into Hub infrastructure via the submit integration

Туре	Succeeded	Failed	Incomplete	Total	Retries	Total+Retries
Tasks	135002	0	0	135002	0	135002
Jobs	4529	0	0	4529	0	4529
Sub-Workflows	2	0	0	2	0	2

Workflow wall time: 13 hrs, 2 mins, (46973 secs)Workflow cumulative job wall time: 384 days, 5 hrs, (33195705 secs)Cumulative job walltime as seen from submit side: 384 days, 18 hrs, (33243709 secs)



## **Workflow Debugging Through Pegasus**

 After a workflow has completed, we can run pegasusanalyzer to analyze the workflow and provide a summary of the run

## pegasus-analyzer's output contains

- a brief summary section
  - showing how many jobs have succeeded
  - and how many have failed.
- For each failed job
  - showing its last known state
  - exitcode
  - working directory
  - the location of its submit, output, and error files.
  - any stdout and stderr from the job.

# Integrated with Submit . Alleviates the need for searching through lots of logs files and integrated with submit command.



## **Different Types of Analysis tools on HubZero**

### Parameter Sweep Analysis

 Execute independent analysis steps on a set of input data, and merge the data and inspect it

### Sequential Analysis

 Execute analysis steps sequentially one after the other, with each step taking in the input of the previous steps.

## DAG based analysis

Multistage analysis where each stage can be a single job or multiple stages

20





# **Hubzero Pegasus Integration**





## How to generate the workflow

 Tool description file tool.xml specifies a python run script (driver file) called by the Rappture interface when you hit the submit button

Two options on how to generate the DAX in python run script

- 1. Use the Pegasus Python DAX API to generate the DAX in the run script
- 2. Call out to an external application-specific DAX generator
- Compose the workflow programmatically and launch it using the Hub command line interfaces (Workspaces + submit).







### **USC**Viterbi



Acknowledgements: Steven Clark , Chris Thomson and Derrick Kearney, Purdue University







## Pegasus Tutorial tool now available in HUBZero

# https://hubzero.org/tools/pegtut







## Wrapper.py

- Python script that is the glue between Rappture and user
- Collects the data from the Rappture interface
- Generates the DAX
- Runs the workflow

Viterbi

 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")
```

## **Rappture (workflow definition)**

f.a

sayhi

f.b

inquire

f.c

Abstract Workflow

```
# Create a 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 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)
dax.addJob(inquire)
```

# Add control-flow dependencies
dax.addDependency(Dependency(parent=sayhi, child=inquire))

# User provides inputs to the workflow and clicks the "Submit" button







#### Workflow has completed. Outputs are available for browsing/downloading



### **USC**Viterbi

# **Submit Command**

- Used by Rappture interface to submit the workflow
- Submits the workflow through Pegasus to
  - OSG

Viterhi

- DIAGRID
- Local Cluster
- 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 pegasus tools.
  - Error reports in file pegasus.analysis
  - Job statistics and status pegasusjobstats.csv pegasusstatus.txt



# Examples of Pegasus Usage in Hubs





# **Pegasus on Hub Platforms**

## Using Workspaces

- Pegasus installed on all hubs
- Users can directly submit and compose workflows through Pegasus and execute on OSG, DiaGrid.

## Integrated into published tools on various Hubs

- Various tools in NanoHub, DiaGrid, NEESHub





# **Pegasus on Hub Platforms**

#### Using Workspaces

- Pegasus and submit installed on all hubs
- Users can directly submit and compose workflows through Pegasus and execute on OSG, Diagrid.

#### Nanohub

- Perform parameter sweep on variables (CNTFET Lab, nanoFET, NanoPlasiticity)
- Jobs are executed on Open Science Grid and DiaGrid
- **CNTFET Lab** Simulates using carbon nanotubes as field effect transistors
- nanoFET Simulates 2D mosfet devices.
- NanoPlasticity Investigates how nano-crystalline materials deform, includes uncertainty quantification.

#### DiaGrid

- **BLASTer** online tool to run BLAST (Basic Local Alignment Search Tool) on the DiaGrid Hub.
- Cryo-EM Electron cryo-microscopy (cryo-EM) for 3-D structure of large macromolecular machines
- **SubmitR** Run R scripts on high-performance computing resources.
- NEESHub
  - OpenSEES suite of simulation tools for submitting NEES scripts and for education and outreach.





# **Pegasus on Hub Platforms**

🛈 Input 🔸 🕑 Simulate 🤍 Addit his tool gestious	O Input * 🕑 Simulate 🔗 Abot	tris tool αστ	Questions?
Poput      Poput	€ Simulation Swring rimulation Swring rimulation Swring rimulation Swring rimulation Name of the selected of 1 yob instance. Wed Sep 9 12:26:12 2015 Runn 4774991 registered 1 yob instance. Wed Sep 9 12:26:32 2015 (10551555:0) BAB Subhattbrd at WF-Diskrid Wed Sep 9 12:25:40 2015 (10551555:0) Simulation Running at WF-Diskrid Wed Sep 9 12:29:24 2015 NanoPlasticity Tool in NanoPlasticity Tool in NanoPlasticity Tool in	Result Response (YoungsModulus vs PelerisEnergy)	A hart fit hold     (our dices)     (our
(111) d1 = 032 Note: The current version of this tool requires d1=d2=d3. Simulate >	NanoHub Stating (hput	7 0 000 0 0.60 0.70 0.60 1 result < input	8.95931e-09 Clear

## a) Specify and configure simulation in the tool UI

b) Workflow Executing through Pegasus

c) Inspect outputs in tool UI

### Rappture Based Tools

- Tools built using the Rappture GUI Builder
- Users configure the simulation, prepare inputs using Rappture tool UI
- Tool execution results in workflows launched through Submit/Pegasus
- Jobs execute on DiaGrid and Open Science Grid, and outputs staged back to Hub
- Users visualize the outputs in the tool UI
- Tools in NanoHub, DiaGrid, NEESHub

### **USC**Viterbi



## **Relevant Links**

## Pegasus: <u>http://pegasus.isi.edu</u>

- Tutorial and documentation: <u>http://pegasus.isi.edu/wms/docs/latest/</u>
- Support: <u>pegasus-users@isi.edu</u>
   <u>pegasus-support@isi.edu</u>

## Rappture & Submit:

- <u>https://hubzero.org/documentation/2.0.0/tooldevs/</u> <u>grid.rappture\_submit</u>
- Submit Command
  - <u>https://hubzero.org/documentation/2.0.0/tooldevs/grid.submitcmd</u>
- Pegasus Workflows on HubZero
  - https://hubzero.org/documentation/2.0.0/tooldevs/grid.pegasuswf
  - Tutorial <u>https://hubzero.org/tools/pegtut</u>



