

Science Automation with the Pegasus Workflow Management System

Ewa Deelman

USC Information Sciences Institute http://www.isi.edu/~deelman Funding from DOE, NSF, and NIH



Information Sciences Institute

The Problem

- Scientific data is being collected at an ever increasing rate
 - The "old days" -- big, focused experiments- LHC
 - Today also "cheap" DNA sequencers and an increasing number of them
- The complexity of the computational problems is ever increasing
- Local compute resources are often not enough
 - Too small, limited availability
 - Data sets are distributed
- The computing infrastructure keeps changing
 - Hardware, software, but also computational models





Our approach

- Provide a way to structure applications in such a way that enables them to be automatically managed
 - In a portable way: same description that works on different resources
 - In a way that scientists can interpret the results
- Develop a system that
 - Maps the application description onto the available resources
 - Manages its execution on heterogeneous resources
 - Sends results back to the user or archive
 - Provides good performance, reliability, scalability







- Scientific Workflows and Application Examples
- Managing scientific workflows
- Pegasus and its features
- Conclusions





Scientific Workflows

- Structure an overall computation
- Define the computation steps and their parameters
- Define the input/output data, parameters
- Invoke the overall computation

USC Viterbi

School of Engineering

- Reuse with other data/parameters/ algorithms and share
- Workflows can be hidden behind a nice user interface (e.g. portal)





Science-grade Mosaic of the Sky



Montage Workflow

Amazon M1 large with 2 cores

Size of mosaic in degrees Number of		Number of	Number of intermediate	Total data Cummulative		
	Square	input data mes	lasks	IIICS	ιοοιριπι	
	1	84	387	770	1.8 GB	11 mins
	2	300	1442	2880	6.4 GB	43 mins
					a program the	
	4	685	3738	7466	17 GB	1 hour, 56 mins
						3 hours, 42
	6	1461	7462	14904	35 GB	mins
		× 1			1.10	6 hours, 45
	8	2565	12757	25480	59 GB	mins

Some workflows are large-scale and data-intensive

Montage Galactic Plane Workflow



John Good (Caltech)

× 17

- Montage Galactic Plane Workflow
 - 18 million input images (~2.5 TB)
 - 900 output images (2.5 GB each, 2.4 TB total)
 - 10.5 million tasks (34,000 CPU hours)

Need to support hierarchical workflows and scale









- Scientific Workflows and Application Examples
- Managing scientific workflows
- Pegasus and its features
- Conclusions





Specification: Place Y = **F**(**x**) at **L**

- Find where x is--- {S1,S2, ...}
- Find where F can be computed---- {C1,C2, ...}
- Choose c and s subject to constraints (performance, space availability,....)
- Move x from s to c
 Error! x was not at s!
 Move F to c
- Compute F(x) at c
 Error! F(x) failed!
- Move Y from c to L
- Register Y in data registry
- Record provenance of Y, performance of F(x) at c

Error! *there is not enough* space at L!





Error! c crashed!



Workflows can be simple!







Sometimes you want to "hide" the workflow

The George E. Brown, Jr. Network	Thomas Hacker (tjhacker) 2651 New Messages Cogout My NEEShub A Dia You Dia For Earthquake Engineering Simulation
About NEES	Tools & Resources Learning & Outreach Project Warehouse Simulation Sites Collaborate Explore NEEShub Support 🕐
You are here: 🥎 Home 🤉	GROUPS » OpenSees Workflows on NeesHub - Pegasus » Wiki » Main Page
Pegasus	OpenSees Workflows on NeesHub - Pegasus ► Wiki Main Page Article Edit Comments History Delete Main Page
Group Member 👻	Introduction This page documents the effort to run <u>2 OpenSees</u> workflows through <u>2 NeesHub/Pegasus</u> on the OSG. The workflow setup is done using
Overview	Rappture interface on <u>PNeesHub</u> , and submitted via Pegasus on the OSG and other resources using the submit command. With HUBZERO
L Members	Rappture Interface
🛱 Wiki	The Rappture interface is being developed by Frank <u>McKenna</u> . The purpose is for the user to setup the workflow using the <u>OpenSees</u>
Resources	Same ecrossibility about appared preparties, record selections, column properties and floer properties are shown below.
🧙 Discussion	
Messages	🗷 Xnest 🔤 🔤
🗭 Blog	Application:
V Wish List	Obellogee
Data Sharing	O Constituir & O Constal Departition & O Colored Departition & O Colored Departition & O Constaling & O Constaling
🛱 Calendar	Graphic + Graeneral Propercies + Graecoro Selección + Grannin Propercies + Graon Propercies + Granniace
Discoverability: Visible Join Policy: Open	Earthquake Records Source: PEER NGA Steel Properties Credit: Frank McKenna UC Berkeley, NEES, HUBzero





Sometimes the environment is complex



School of Engineering



Sometimes the environment is just not exactly right

Single core workload





Cray XT System Environment / ALPS / aprun

• Designed for MPI codes





Sometime you want to change or combine resources



Workflow Management

- Assume a high-level workflow specification
- Assume the potential use of different resources within a workflow or over time
 - 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 provenance information
- Need to provide scalability, performance, reliability







- Scientific Workflows and Application Examples
- Managing scientific workflows
- Pegasus and its features
- Conclusions





Our Approach

Analysis Representation

- Support a declarative representation for the workflow (dataflow)
- Represent the workflow structure as a Directed Acyclic Graph (DAG)
- Tasks operate on files
- 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
 - 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





Submit locally, compute Globally



School of Engineering



Pegasus

Workflow Management System (est. 2001)

- A collaboration between USC and the Condor Team at UW Madison
- Maps a resource-independent "abstract" workflow onto resources and executes the "concrete" workflow
- 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, 1-10⁶ tasks)
- Infers data transfers, restructures workflows for performance
- Automatically captures provenance information
- Can run on resources distributed among institutions, laptop, campus cluster, Grid (OSG, XSEDE), Cloud (Amazon, FutureGrid)





Pegasus Workflow Management System

- A workflow "compiler"
 - Input: abstract workflow description, resource-independent
 - Auxiliary Info (catalogs): available resources, data, codes
 - Output: executable workflow with concrete resources
 - Automatically locates physical locations for both workflow tasks and data
 - Transforms the workflow for performance and reliability
- A workflow engine (DAGMan)
 - Executes the workflow on local or distributed resources (HPC, clouds)
 - Task executables are wrapped with pegasus-kickstart and managed by Condor schedd
- Provenance and execution traces are collected and stored
- Traces and DB can be mined for performance and overhead information







Generating executable workflows







Pegasus optimizations address issues of:

- Failures in the execution environment or application
- Data storage limitations on execution sites
- Performance
 - Small workflow tasks
- Heterogeneous execution architectures
 - Different file systems (shared/non-shared)
 - Different system architectures (Cray XT, Blue Gene, ...)





Sometimes fatal errors occur during workflow execution



Storage limitations create dir SI "Small" amount of space f.a f.a hello RM f.b f.a world RM f.c f.b ¥ SO f.c Reg RM f.c f.c

Automatically add tasks to "clean up" data no longer needed





LIGO and Montage



School of Engineering

Storage limitations

Variety of file system deployments: shared vs non-shared

School of Engineering

Workflow Restructuring to improve application performance

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

Level-based clustering Label-based clustering Time-based clustering

Southern California Earthquake Center

Workflows have different computational needs

a time

Develop an MPI-based workflow management engine to manage sub-workflows

time

USC Viterbi

School of Engineering

Pegasus-MPI-Cluster

- A master/worker task scheduler for running fine-grained workflows on batch systems
- Runs as an MPI job
 - Uses MPI to implement master/worker protocol
- Works on most HPC systems
 - Requires: MPI, a shared file system, and fork()
- Allows sub-graphs of a Pegasus workflow to be submitted as monolithic grid jobs to remote resources

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
 - Status of the workflow

- pegasus-statistics

Detailed statistics about your finished workflow

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)

Collaboration with Dan Gunter and Taghrid Samak, LBNL

Typical Deployment

Posix Access for Tasks in the workflow

- How do you ensure Posix access for the tasks?
 - Place it on a shared filesystem shared across nodes.
 - Place it directly on local filesystem of the worker node from the input site.
- Direct Transfers to local filesystem
 - Job starts and retrieves input data from input site.
 - Not efficient for large datasets that are shared across jobs.
- Shared Filesystem sounds appealing but problems for Big Data workflows
 - Shared storage at a computational site maybe limited. Cannot accommodate all files required for a large workflow.
 - In some cases, shared filesystem may have limited scalability---NFS
 - Harder to setup a shared filesystem in a dynamic environment like computational clouds
 - Some systems just don't support it

Object Storage for Workflows

- Clouds such as Amazon provide object stores--- S3
- Object Store: high level storage service with limited operations
 - Store, retrieve and delete data objects (files)
 - Doesn't provide byte level access
 - Cannot open a file in an object store, read and update it and then close it.
 - Instead a client needs to download the file, update it and then store as a new object.
- View traditional Grid services like GridFTP, SRM, IRODS as object stores
 - Store, retrieve and delete data (files)
 - Don't support random read or writes like object stores.
 - This generalization is important to lay out the different data management models.

General Workflow Execution Model, Cannot Assume Shared FS

- --- Task Flow

----- Data Flow

- Input Data Site, Compute Site and Output Data Sites can be co-located
 - Example: Input data is already present on the compute site.

Exclusive Use of Object Stores

Viterbi

School of Engineering

Advantages

- Can leverage scalable stores
- Distribute computations across resources, such as supporting spillover from local resources to cloud resources.
- Great bandwidth

Disadvantages

- Duplicate Transfers
- Latencies in transferring large number of files
- Added costs for duplicate transfers.

The Workflow System retrieves files from Object Store and makes it available to the workflow task on the local disk on a worker node.

Pegasus-kickstart

- Lightweight C based executable to launch jobs
- Captures job runtime provenance and logs it as a XML record
- Following information is captured about each job on all supported platforms
 - exit code with which the job it launched exited
 - start time and duration of the job
 - hostname and IP address of the host the job ran on
 - stdout and stderr of the job
 - arguments with which it launched the job
 - directory in which the job was launched
 - environment that was set for the job before it was launched
- Additional profiling
 - peak memory usage (resident set size, and vm size)
 - total I/O read and write,
 - Pid
 - all files accessed (total read and write per file)

Workflow Monitoring Dashboard – pegasus-dashboard

Workflow Statistics								
tistics								
	8 mins 53 se	cs						
	Workflow	Cumulative Job Wa	all Time			1 min 59 sec	5	
Cum	Cumulative Job Walltime as seen from Submit Side							
Workflow Retries								
Workflow Statistics								
Job Breakdown Statistics								
Show 50 + entries					5	Search:		
Transformation -	Count	Succeeded	Failed	≎ Min	◇ Max	Mean	Total	
dagman::post	32	32	0	5	6	5.063	162	
mAdd:3.3	1	1	0	1.203	1.203	1.203	1.203	
mBackground:3.3	32	32	0	0.054	0.197	0.130	4.174	
mBgModel:3.3	1	1	0	18.701	18.701	18.701	18.701	
mConcatFit:3.3	1	1	0	1.033	1.033	1.033	1.033	
mDiffFit:3.3	73	73	0	0.048	0.226	0.103	7.492	
mlmgtbl:3.3	1	1	0	0.107	0.107	0.107	0.107	
mJPEG:3.3	1	1	0	0.523	0.523	0.523	0.523	
mProjectPP:3.3	32	32	0	0.915	0.978	0.926	29.633	
mShrink:3.3	1	1	0	0.485	0.485	0.485	0.485	
pegasus::cleanup	1	1	0	5	5	5	5	
pegasus::dirmanager	1	1	0	10	10	10	10	
pegasus::rc-client	4	4	0	0.706	0.868	0.783	3.134	
pegasus::transfer	14	14	0	0	5.229	2.724	38.135	
Showing 1 to 14 of 14 entries					F	irst Previous	1 Next Last	

Status, statistics, timeline of jobs

Helps pinpoint errors

School of Engineering

USC Viterbi

Tools to calculate job statistics

Task Type	Count	Runtime(s)	IO Read (MB)	IO Write (MB)	Memory Peak(MB)	CPU Utilization(%)
mProjectPP	2102	1.73	2.05	8.09	11.81	86.96
mDiffFit	6172	0.66	16.56	0.64	5.76	28.39
mConcatFit	1	143.26	1.95	1.22	8.13	53.17
mBgModel	1	384.49	1.56	0.10	13.64	99.89
mBackground	2102	1.72	8.36	8.09	16.19	8.46
mImgtbl	17	2.78	1.55	0.12	8.06	3.48
mAdd	17	282.37	1102	775.45	16.04	8.48
mShrink	16	66.10	412	0.49	4.62	2.30
mJPEG	1	0.64	25.33	0.39	3.96	77.14

 Table 1. Execution profile of the Montage workflow, averages calculated

dv/dt – Accelerating the rate of progress towards extreme scale collaborative science, (9/12-8/15, DOE)

Objectives

Miron Livny UWMadison, Bill Allcock ANL, Ewa Deelman USC, Doug Thain UND, Frank Wuerthwein UCSD)

- Design a computational framework that enables computational experimentation at scale while supporting the model of "submit locally, compute globally"
- Focus on Estimating application resource needs, Finding the appropriate computing resources, Acquiring those resources, Deploying the applications and data on the resources, Managing applications and resources during run
- Task resource profiling and resource estimation

Task Characterization/Execution

- Collect and archive data from existing infrastructure deployments
- Understand the resource needs of a task (memory, disk, CPU)
- Establish expected values and limits for task resource consumption
- Launch tasks on the correct resources
- Monitor task execution and resource consumption, interrupt tasks that reach resource limits
- Possibly re-launch tasks on different resources
- Task characterization needs to be an online process

Predictive Modeling and Diagnostic Monitoring of Extreme Science Workflow (9/14-8/17, DOE)

Ewa Deelman USC, Chris Carothers RPI, Anirban Mandal RENCI Brian Tierney LBNL, Jeff Vetter ORNL

Objective: Understand complex scientific workflow applications and infrastructure behaviors and to translate this understanding into flexible, end-to-end analytical models that can effectively predict the behavior of extreme scale workflows on current and future infrastructures

Approach:

- Engage DOE science teams from simulation (e.g., Earth Systems Modeling (ESM) and instrument facilities (e.g., Spallation Neutron Source(SNS) to create example workflow scenarios
- Develop a general analytical modeling methodology that captures the endto-end performance of these workflow scenarios using a structured modeling approach
- Validate the analytical models using empirical measurement and simulation
- Employ the analytical performance models to facilitate prototype capabilities that include anomaly detection and diagnosis, resource management and adaptation, and infrastructure design and planning

Benefits of Pegasus

- Provides Support for Complex Computations
 - Can be hidden behind a portal
- Portability / Reuse
 - User created workflows can easily be run in different environments without alteration (XSEDE, OSG, FutureGrid, Amazon)
- 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

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.
- Analysis tools

If you are interested in Pegasus

- Pegasus: <u>http://pegasus.isi.edu</u>
- Tutorial and documentation: <u>http://pegasus.isi.edu/wms/docs/latest/</u>
- Virtual Machine with all software and examples <u>http://pegasus.isi.edu/downloads</u>
- Take look at some Pegasus applications: <u>http://pegasus.isi.edu/applications</u>
- Support: <u>pegasus-users@isi.edu</u> <u>pegasus-support@isi.edu</u>

