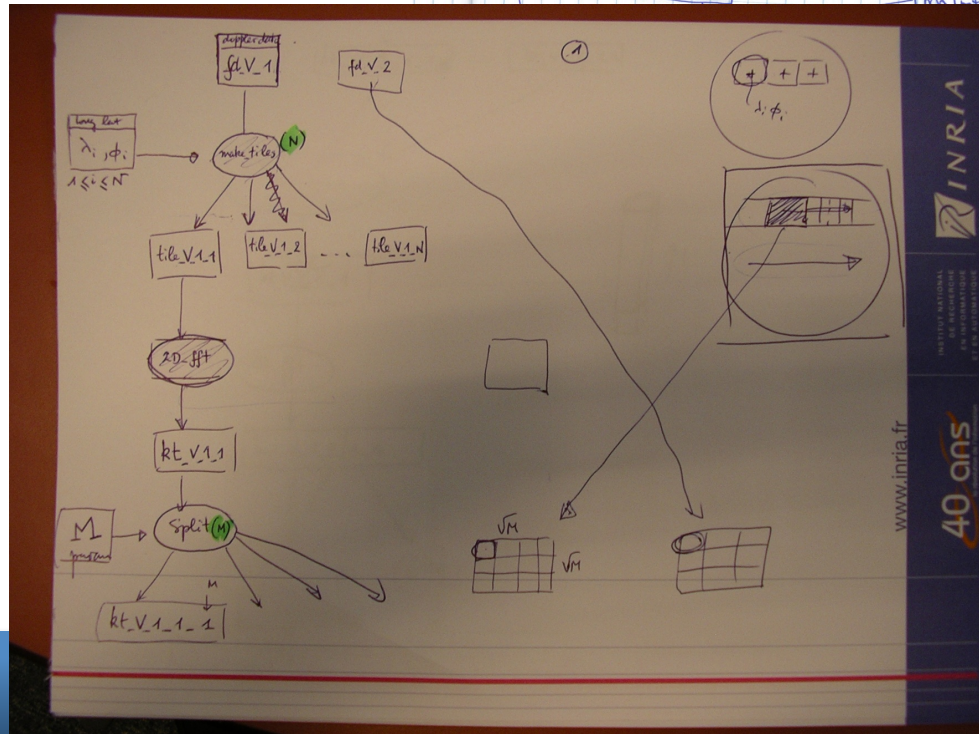# What is missing in workflow technologies?

*Ewa Deelman, Ph.D*
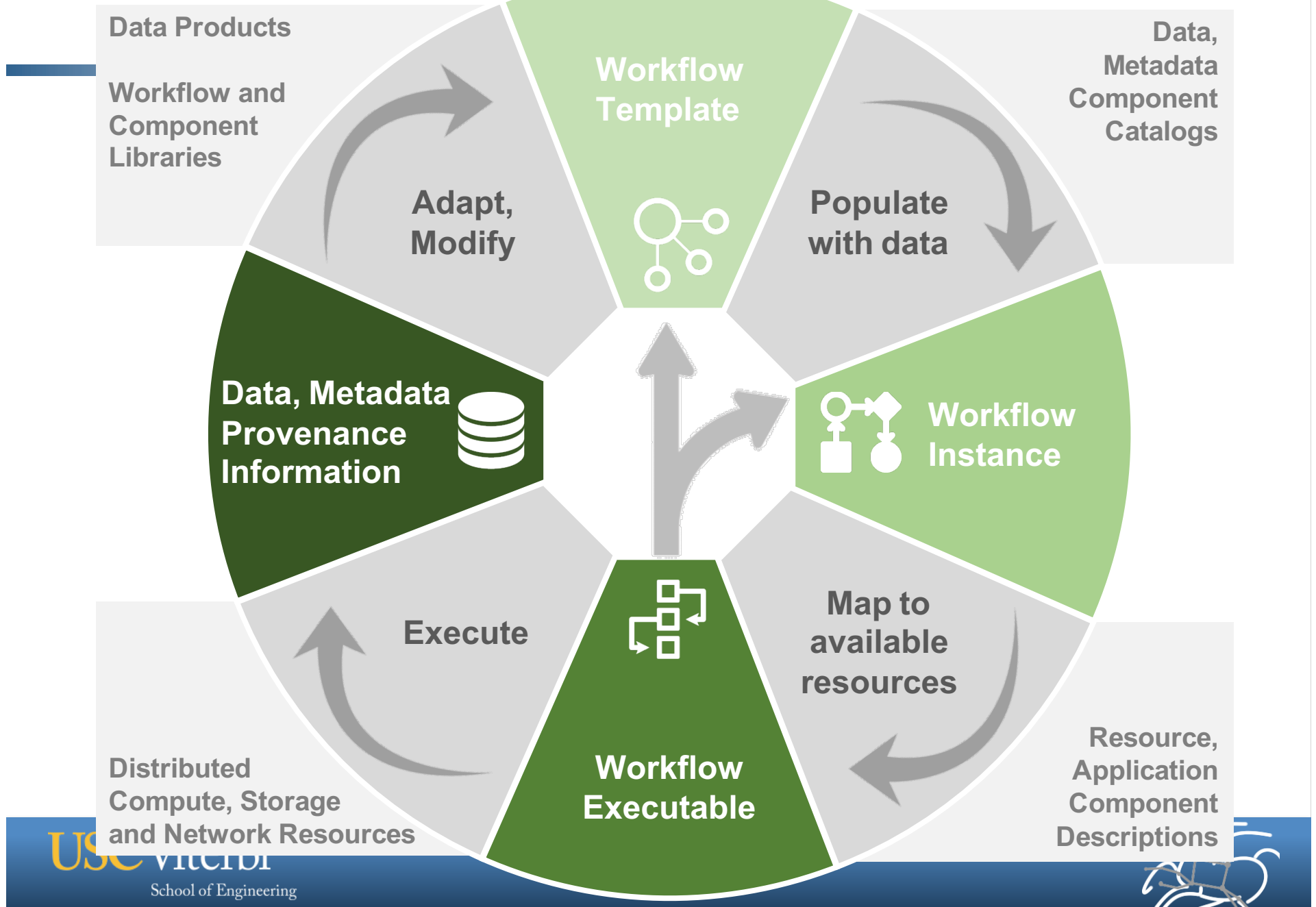
Science Automation Technologies Group

USC Information Sciences Institute
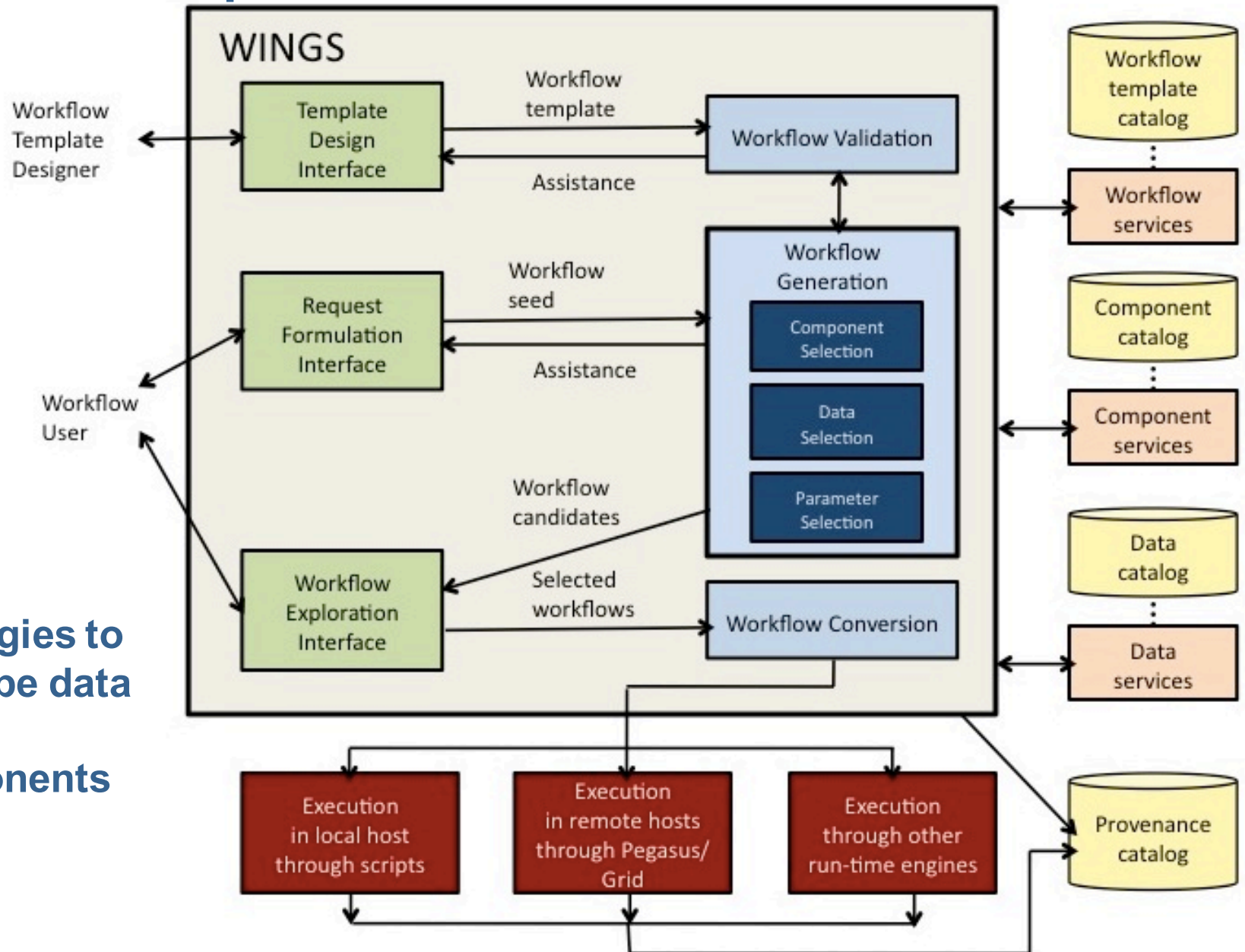
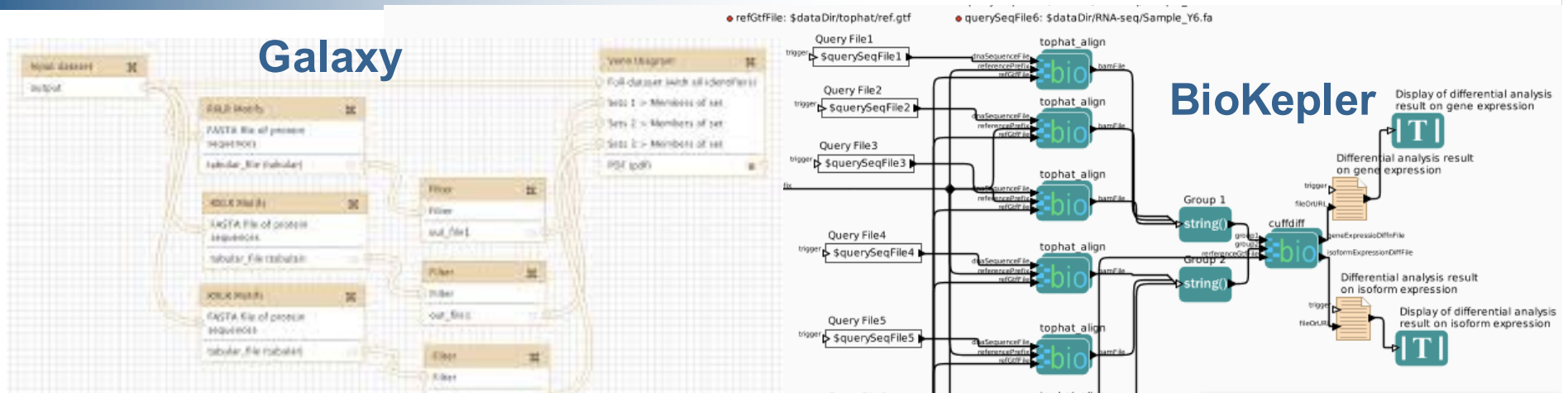*Funding from DOE, NSF, and NIH*

USC Viterbi
School of Engineering

http://deelman.isi.edu

Information Sciences Institute

# The challenge

# Generalized Workflow Lifecycle



- **Workflow Template**
- Populate with data
- **Workflow Instance**
- Map to available resources
- **Workflow Executable**
- Execute
- **Data, Metadata Provenance Information**
- Adapt, Modify

Data Products

Workflow and Component Libraries

Data, Metadata Component Catalogs

Resource, Application Component Descriptions

Distributed Compute, Storage and Network Resources

USC Viterbi
School of Engineering

# Workflow Templates



**Uses ontologies to describe data and components**

# A number of workflow composition frameworks

# Sometime the workflows is behind a portal



**No need to construct workflow**
**Helps with correctness**
**Helps with reproducibility**
**Hard to customize/change**

# Mapping of of Workflows onto Resources -- Provisioning

- **Traditionally resources were already provisioned: XSEDE clusters, DOE LCF systems**
  - Need to submit jobs to a queue and make sure the input data is there
  - Lack of storage provisioning capability or provisioning particular nodes

- **Opportunistic resources used with on-the-fly provisioning**
  - Open Science Grid:  HTC Condor Glideins, GlideinWMS
  - When you land on a resource, need to pull data and push out results before
  - Needs robust fault recovery, as resources can disappear
  - Lack of storage provisioning capability, lacks network provisioning

- **Cloud-based environments**
  - Need to know what to provision and for how long
  - Need to adjust the provisioning over time
  - Need to keep an eye on costs (stay within a budget, minimize the costs)
  - Need a fail-safe mechanism to deprovision resources when no longer needed or when things go wrong

USCViterbi
School of Engineering

# Determining the needed resources

Application Monitoring
- CPU, I/O, memory, perf counters
- Function interposition
- MPI and serial jobs
- Real-time reporting

Infrastructure Monitoring
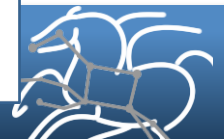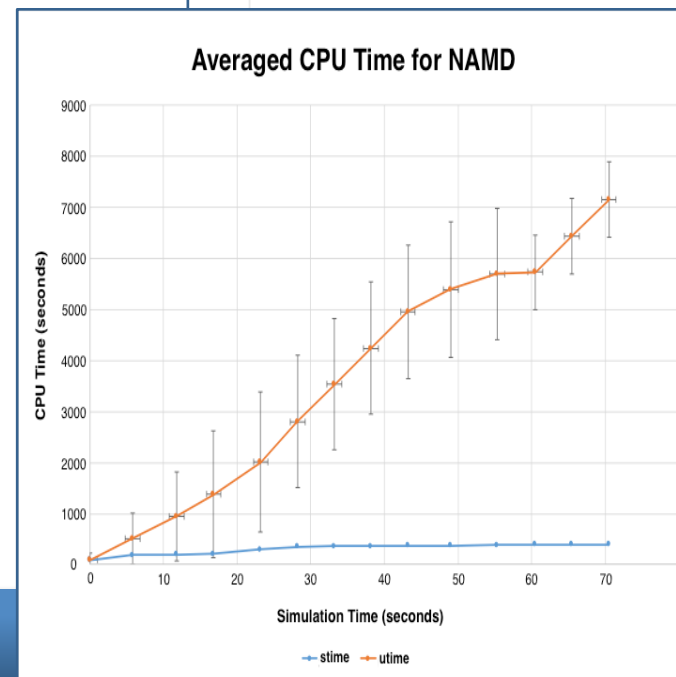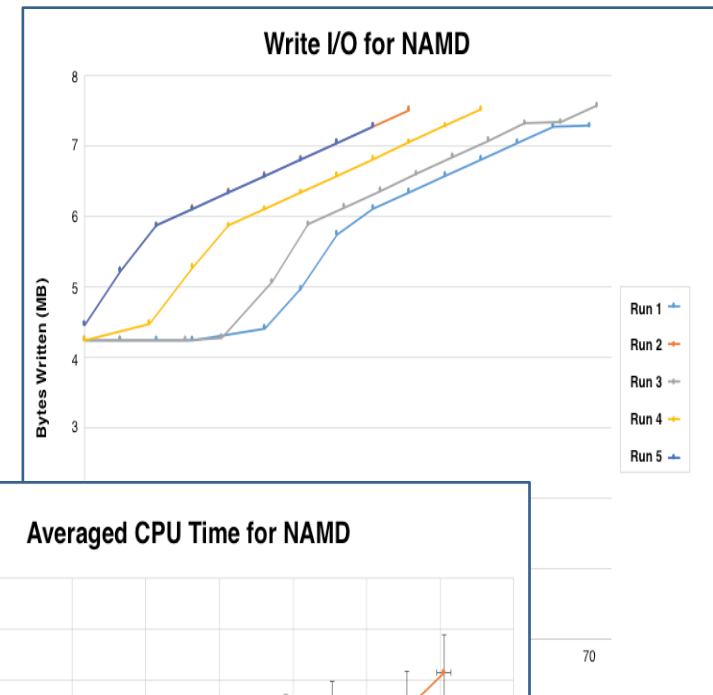- Load, disk I/O, network, etc.
- Standard tools
- Data stored in time series DB

# Analytical Modeling Example for SNS workflow with ASPEN

## MD template model

```
model NAMD_Template {
  // application parameters
  // (defined in the input file)
  param nAtoms    = 1e6
  param nTimeSteps = 100
  // solve for these parameters
  // (within the given ranges)
  param c = 1 in 1 .. 1e18
  param d = 1 in 1 .. 1e18
  // application behavior:
  // execution and control flow
  kernel main
  {
    iterate [nTimeSteps] {
      execute {
        loads [c * nAtoms^2]
        flops [d * nAtoms]
      }
    }
  }
}
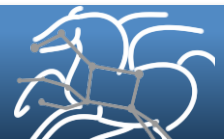```

## CSV data file with parameters and runtimes

+

| nAtoms | nTimeSteps | nCores | machine | runtime |
|--------|------------|--------|---------|---------|
| 1e6 | 100 | 144 | exogeni | 384.2 |
| 1e6 | 100 | 144 | hopper | 340.1 |
| 1e6 | 150 | 144 | hopper | 482.9 |

=

## Concrete NAMD model

```
model NAMD_Equilibrate {

  // NAMD input parameters
  param nAtoms    = 1e6
  param nTimeSteps = 100

  // calculation-specific constants
  param c = 402.1
  param d = 10.95

  // NAMD application behavior
  kernel main
  {
    iterate [nTimeSteps] {
      execute {
        loads [c * nAtoms^2]
        flops [d * nAtoms]
      }
    }
  }
}
```

nAtoms and nTimeSteps defined in template application model and CSV input data
nCores defined in machine models and CSV input data
solves for c and d, filling out a concrete application model for that problem
new predictions can still vary nAtoms, nTimeSteps, and nCores

USC Viterbi
School of Engineering

**Work with Jeff Vetter, ORANL**

# Interleaving Workflow Management and Provisioning



**Work Jeff Chase (ORCA),  Anirban Mandal, Paul Ruth, Ilya Baldin**

# ExoGENI Virtual SDX: Panorama Workflows

Software Defined Exchanges – meeting point of networks to exchange traffic, securely and with QoS, using SDN protocols



- Panorama modeling and simulation tools enable Pegasus to monitor and manipulate network connectivity & performance
- Virtual SDX **transparently arbitrates prioritized workflow data flows** communicated by Pegasus
- Advanced SDX capabilities can monitor and detect network anomalies, and take adaptation actions.

**Work Jeff Chase (ORCA),  Anirban Mandal, Paul Ruth, Ilya Baldin**

# Resource Selection Issues in Distributed Area

- Discover what resources (computation, data, software) are available (or what resources were provisioned)

- Select the appropriate resources based on a architecture, availability of software, performance, reliability, availability of cycles, storage,.. (or provision)
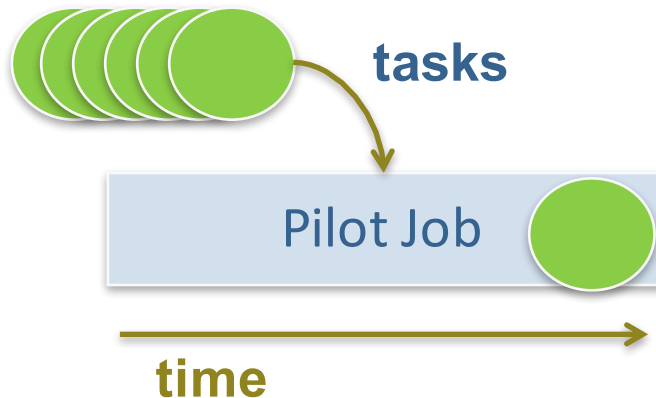
- Devise a plan:
  - What resources to use
  - How to best adapt the workflow to the resources
  - What protocols to use to access the data, to schedule jobs
  - What data to save

- Issues of compute "close" to the data, traditionally data moves
  - Extract subsets, compress, pre-compute some values at the source

- Issues of recompute vs retrieve the results

- Managing data access within a workflow and across workflow ensembles – supporting data reuse

USC Viterbi
School of Engineering

# Matching Workload to Target System

**Cluster tasks**



**Use "pilot" jobs to dynamically provision a number of resources at a time**

tasks

Pilot Job

time

**Partition the workflow into subworkflows and send them for execution to the target system**

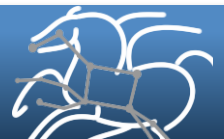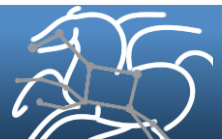# Make Data Flow Over Heterogeneous Fabric

**Typical Deployment in Clusters, sharing data through the file system**



**Tasks access data via Posix I/O**

# Variety of file system deployments: shared vs non-shared

## pegasus-transfer subsystem for various storage systems

- Command line tool used internally by Pegasus workflows

- Input is a list of source and destination URLs

- Transfers the data by calling out to tools – provided by the system (cp, wget, …) Pegasus (pegasus-gridftp, pegasus-s3) or third party (gsutil)

- Transfers are parallelized

- Transfers between non-compatible protocols are split up into two transfers using the local filesystem as a staging point
  - for example: GridFTP->GS becomes GridFTP->File and File->GS

**Supported protocols**

**GridFTP**
**SRM**
**iRods**
**S3**
**GS**
**SCP**
**HTTP**
**File**
**Symlink**

# What's needed in Mapping and Execution

- **Issues of efficiency (time, cost, energy)**

- **Security of data being managed, using trusted resources (new project with Von Welch and Ilya Baldin)**

- **Error reporting**
  - **Easy to interpret, maybe involve machine learning to better categorize errors**

- **Anomaly detection  (Panorama Project)**
  - **In prototype, but need the technologies to do it in production**

# Metadata, Provenance information



- **Some provenance capture standards**

- **Limited or no tools for provenance exploration**

- **Limited use of metadata when doing workflow composition**

- **Now more pressing issues of scientific collaboration**
  - **Transparency**
  - **Re-use**
  - **Reproducibility**

# Future Applications: Near real time feedback, human-in-the loop

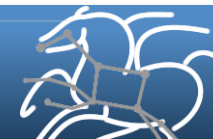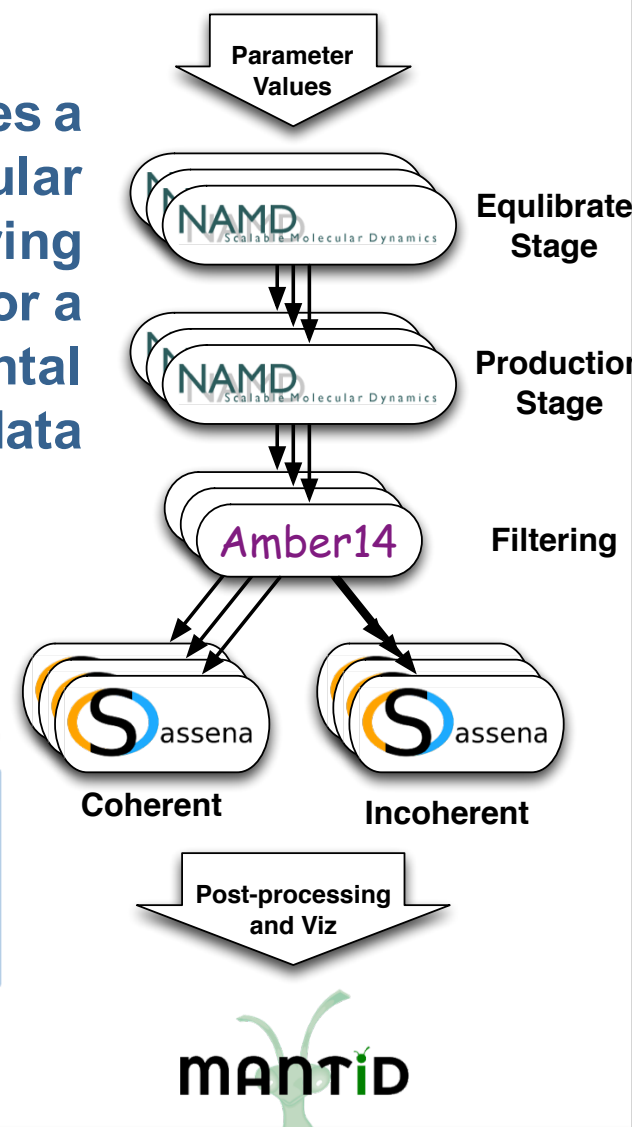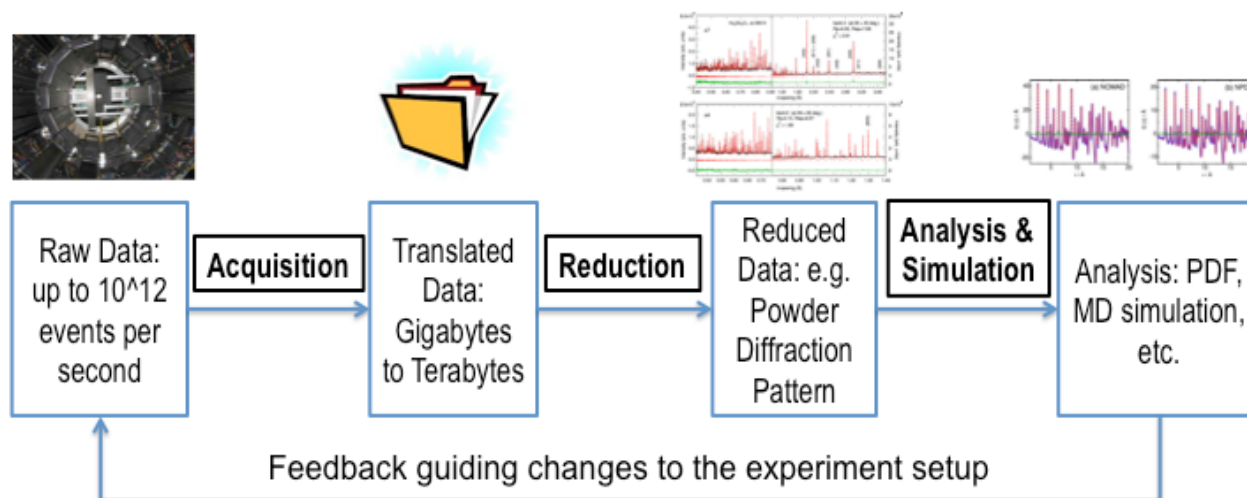**SNS refinement workflow executes a parameter sweep of molecular dynamics and neutron scattering simulations to optimize the value for a target parameter to the experimental data**

**More in-situ workflow management**



Parameter Values

NAMD — Equilibrate Stage

NAMD — Production Stage

Amber14 — Filtering

Sassena — Coherent

Sassena — Incoherent

Post-processing and Viz

MANTiD

Raw Data: up to 10^12 events per second → **Acquisition** → Translated Data: Gigabytes to Terabytes → **Reduction** → Reduced Data: e.g. Powder Diffraction Pattern → **Analysis & Simulation** → Analysis: PDF, MD simulation, etc.

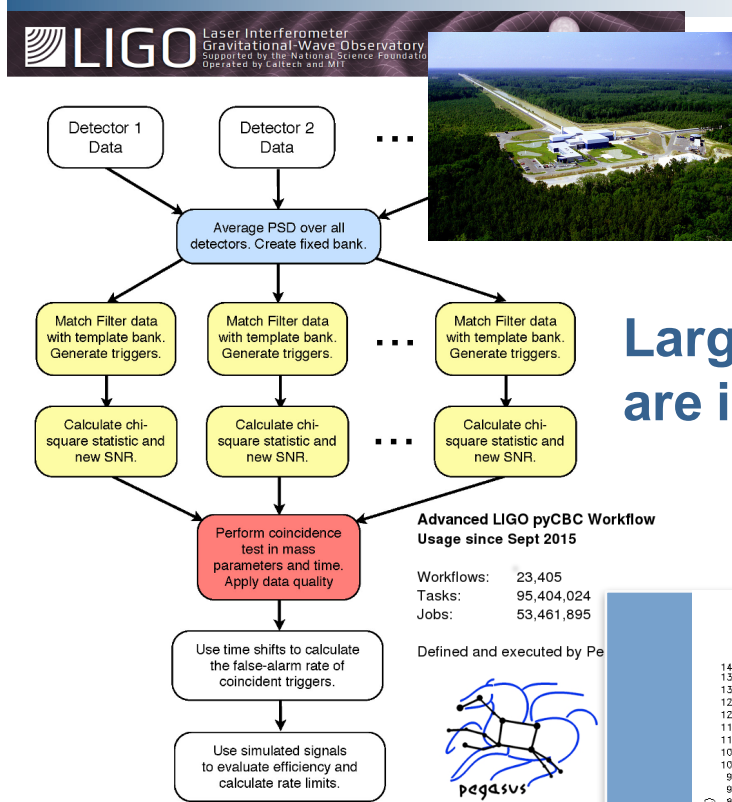Feedback guiding changes to the experiment setup

# What we learned in distributed area WMS
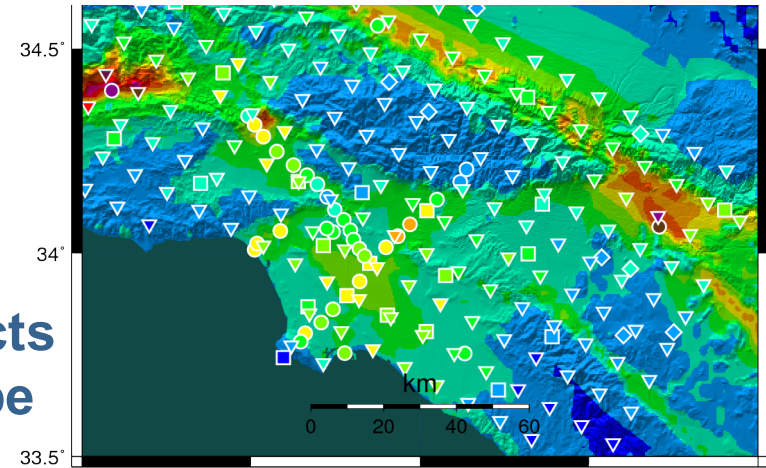# We can apply In-Situ workflow management

- **Provision HPC resources ahead of execution:**
  - Keep the resources for multiple tasks, exploit data locality
  - Support provisioning of storage/ incl. NVRAM, burst buffers
  - Alternatively explore interfaces between the WMS and the scheduler, support data-aware scheduling

- **Reliability: WMS deal with: task failures, problems accessing data, resource failures, and others.**
  - Investigate how data replication techniques can be used to improve fault tolerance, while minimizing the impact of energy consumption
  - Explore tradeoffs between data re-computation and data retrieval from DRAM/NVRAM/disk (time to solution and energy consumption)

- **Provenance Capture and Reproducibility: WMS capture provenance information about the creation, planning, and execution**
  - Provenance capture may need to adapt to the behavior of the application (coarse and fine levels of details, compression)
  - May want to automatically re-run parts of the computation and re-produce the results and a more detailed provenance trail on demand

# Developing solutions that impact science



**Large scale projects are in decent shape**

**Southern California Earthquake Center**

**Individuals are still struggling**

http://www.teachingboxes.org/avc/content/Cloud_Radar.htm