



Workshop on the Future of Scientific Workflows

April 20-21, 2015, Rockville, MD

Co-organizers

Ewa Deelman (USC) and Tom Peterka (ANL)

Program committee

Ilkay Altintas (SDSC), Chris Carothers (RPI), Ken Moreland (SNL),
Manish Parashar (Rutgers), Lavanya Ramakrishnan (LBNL), Jeff Vetter (ORNL),
Kerstin Kleese van Dam (BNL), Michela Taufer (UD)

Program managers

Lucy Nowell (DOE) and Rich Carlson (DOE)



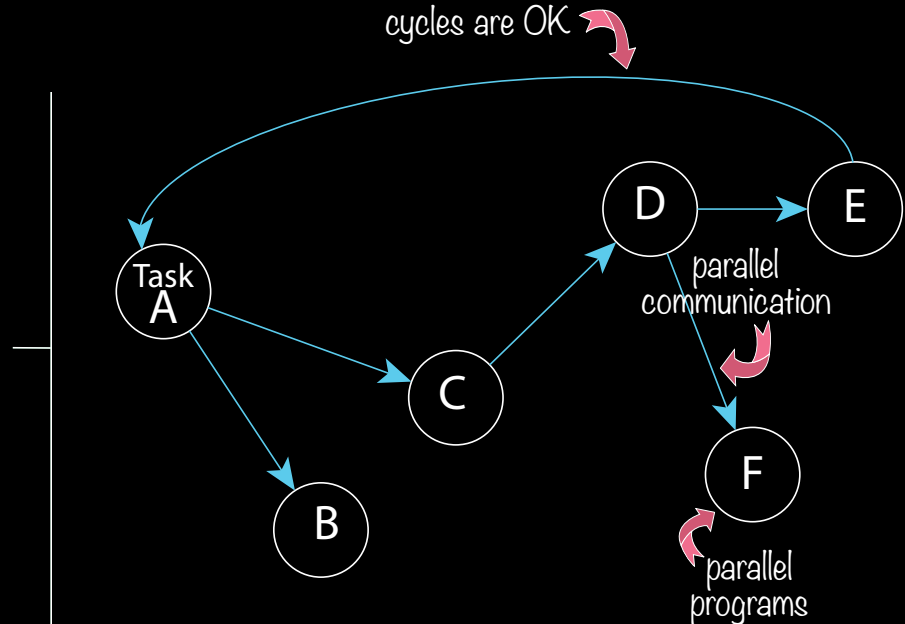
Workflows 101

Definitions

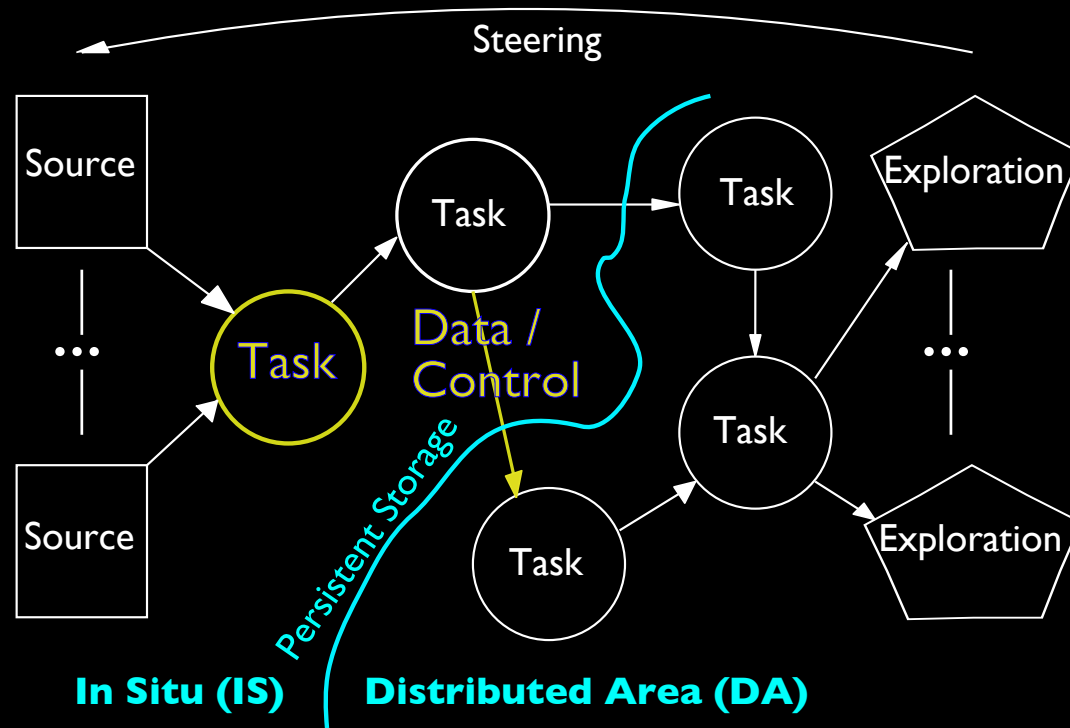
- **Workflow:** Sequencing and orchestrating operations, along with the attendant tasks of, for example, moving data between workflow processing stages.
- **Workflow management systems:** Aiding in the automation and capture the provenance of these processes, freeing the scientist from the details of the process.
 - Manage the execution of constituent tasks
 - Manage the information exchanged between them

Graph Model

- Directed graph of tasks and communication
- Graph nodes are the tasks
- Graph links are the communication
- Graph does not have to be acyclic
- “Large tasks” (programs), not “small tasks” (threads)
- Nodes and links can be parallel



In Situ (IS) and Distributed Area (DA)



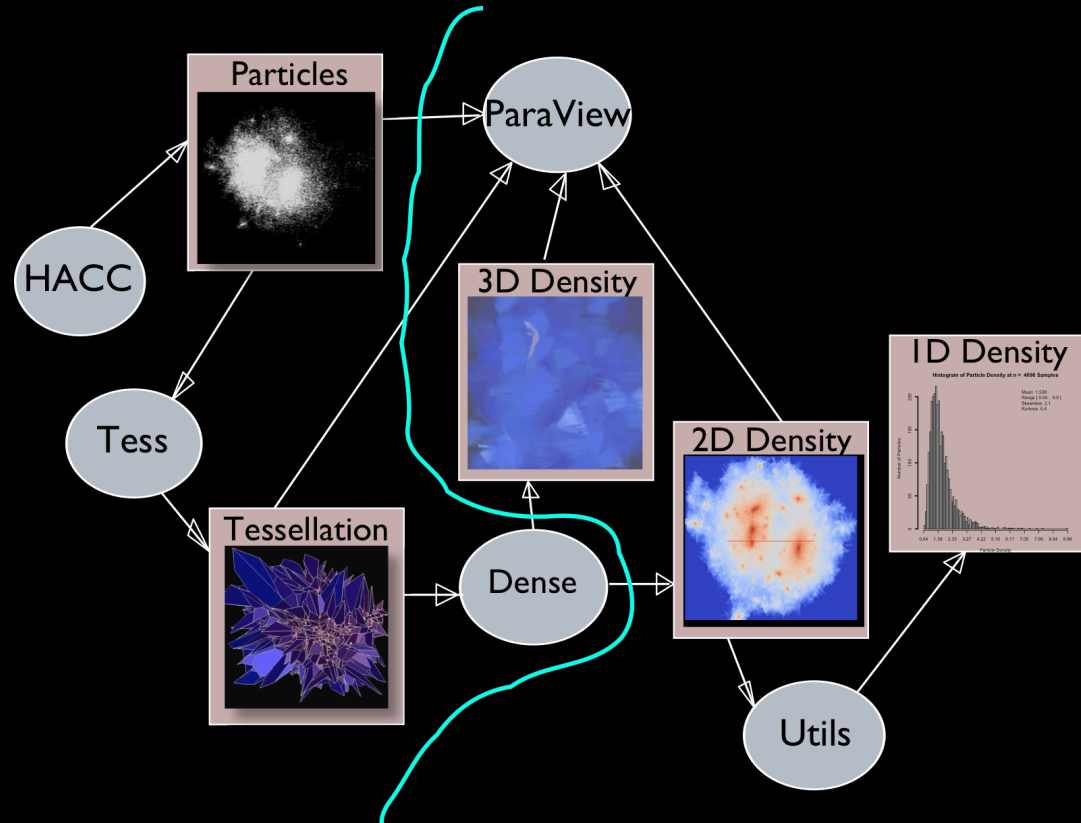
- **In Situ (IS):** Within an HPC system (synonyms: in situ, in transit, coprocessing, run-time, online)
- **Distributed Area (DA):** Across systems, potentially geographically distributed

Examples

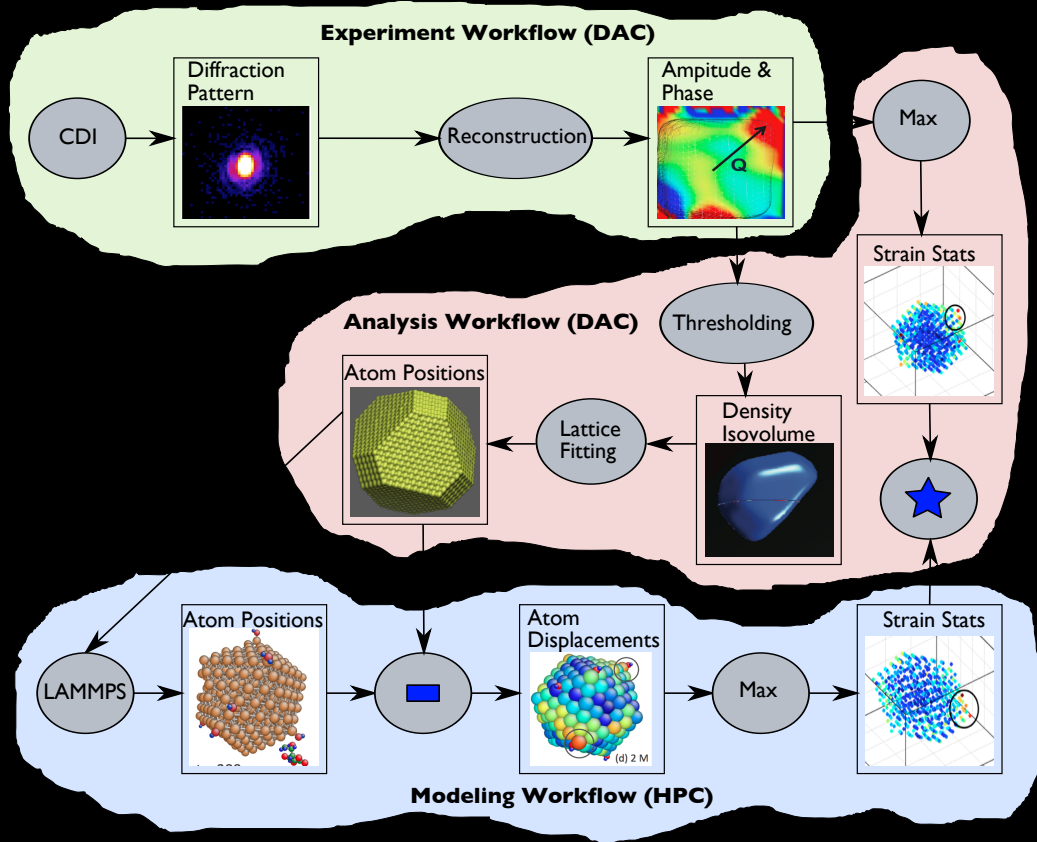
Cosmology: HPC Simulation



- Computational workflow is one small part of a complete cosmology campaign
- Converts dark matter particles to an unstructured mesh
- Converts an unstructured mesh to a regular grid
- Computes statistics over the grid and visualizes the results



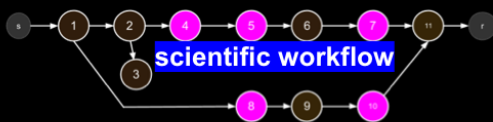
Materials Science: Simulation + Experiment



Science workflow for the comparison of a molecular dynamics simulation with a high-energy X-ray microscopy of the same material system includes three interrelated in situ and distributed area experimental workflows.

Bioinformatics: High Throughput + Cloud

Scalable multicloud analysis for Ultra-High Throughput Bioinformatics



High-throughput workflows often feature many small tasks, web services, cloud architectures, and big data tools.



control



data



software



software



Future work:

- Extend to NERSC
- Tie into more workflow systems (e.g. SWIFT or Pegasus)



Shock object store: Versatile data storage for data; native support for scientific metadata; multi-cloud capable

Skyport: Linux containers provide an isolated execution environment for workflow software. Skyport orchestrates the deployment of workflow software onto AWE worker machines using Docker images stored in Shock.

Gerlach et al, 2015 IEEE International Conference on Cloud Engineering
Tang et al, 2013 IEEE International Conference on Big Data

The Workshop

Mission

Develop requirements for workflow methods and tools in a combined IS and DA environment to enable science applications to better manage their end-to-end data flow.

Objectives

- **Identify** the workflows of representative **science use cases** in IS and DA settings
- **Understand** the state of the art in **existing workflow technologies**, including creation, execution, provenance, (re)usability, and reproducibility
- **Address** emerging **hardware and software trends**, both in centralized and distributed environments, as they relate to workflows
- **Bridge** the gap between **IS and DA workflows**

Findings and Research Priorities Related to SSIO

- Applications: Lagging I/O bandwidth motivates in situ workflows
- Hardware: New storage technology, heterogeneity, hierarchy
- System software: Provisioning storage resources and services as first-class citizens
- Dataflow: Storage for communication
- Provenance: Fast storage and retrieval of log data (e.g., performance profiles)
- Validation: Predicting performance and validate results (e.g., storage systems)

Report

http://science.energy.gov/~media/ascr/pdf/programdocuments/docs/workflows_final_report.pdf

Journal paper

Ewa Deelman, Tom Peterka, Ilkay Altintas, Christopher Carothers, Kerstin Kleese van Dam, Kenneth Moreland, Manish Parashar, Lavanya Ramakrishnan, Michela Taufer, Jeffrey Vetter: The Future of Scientific Workflows. International Journal of High Performance Computing Applications, 2017.

Acknowledgments

Meeting Organizers

Ewa Deelman (USC)
Tom Peterka (ANL)
Ilkay Altintas (SDSC)
Christopher Carothers (RPI)
Kerstin Kleese van Dam (PNNL)
Kenneth Moreland (SNL)
Manish Parashar (RU)
Lavanya Ramakrishnan (LBNL)
Michela Taufer (UD)
Jeffrey Vetter (ORNL)

Meeting Participants

Greg Abram (UT)
Gagan Agrawal (OSU)
Jim Ahrens (LANL)
Pavan Balaji (ANL)
Ilya Baldin (RENCI)
Jim Belak (LLNL)
Amber Boehnlein (SLAC)
Shreyas Cholia (NERSC)
Alok Choudhary (NU)
Constantinos Evangelinos (IBM)
Ian Foster (ANL)
Geoffrey Fox (IU)
Foss Friedman-Hill (SNL)
Jonathan Gallmeier (AMD)
Al Geist (ORNL)
Berk Geveci (Kitware)
Gary Grider (LANL)
Mary Hall (UU)

Ming Jiang (LLNL)
Elizabeth Jurrus (UU)
Gideon Juve (USC)
Larry Kaplan (Cray)
Dimitri Katramatos (BNL)
Dan Katz (UC)
Darren Kerbyson (PNNL)
Scott Klasky (ORNL)
Jim Kowalkowski (FNAL)
Dan Laney (LLNL)
Miron Livny (UW)
Allen Malony (UO)
Anirban Mandal (RENCI)
Folker Meyer (ANL)
Dave Montoya (LANL)
Peter Nugent (LBNL)
Valerio Pascucci (UU)
Wilfred Pinfold (Intel)
Thomas Proffen (ORNL)

Rob Ross (ANL)
Erich Strohmaier (LBNL)
Christine Sweeney (LANL)
Douglas Thain (UND)
Craig Tull (LBNL)
Tom Uram (ANL)
Dean Williams (LLNL)
Matt Wolf (GT)
John Wright (MIT)
John Wu (LBNL)
Frank Wuerthwein (UCSD)
Dantong Yu (BNL)

DOE Program Managers

Richard Carlson
Lucy Nowell

http://science.energy.gov/~media/ascr/pdf/programdocuments/docs/workflows_final_report.pdf Tom Peterka