

# Automating Science with Workflows at OLCF

Ketan Maheshwari  
Sean Wilkinson  
Rafael Ferreira da Silva

Data Lifecycle & Scalable Workflows  
National Center for Computational Sciences

ORNL is managed by UT-Battelle LLC for the US Department of Energy



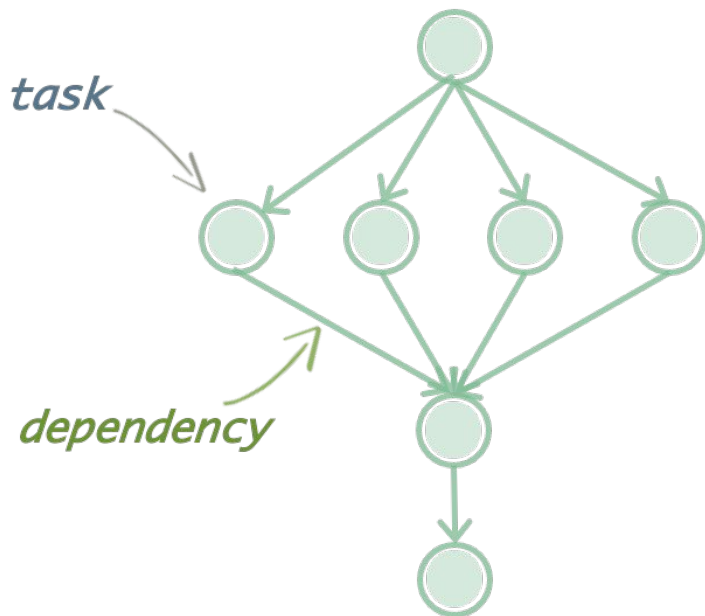
U.S. DEPARTMENT OF  
**ENERGY**

# What is involved in an experiment execution?



# Scientific Workflows

directed-acyclic graphs

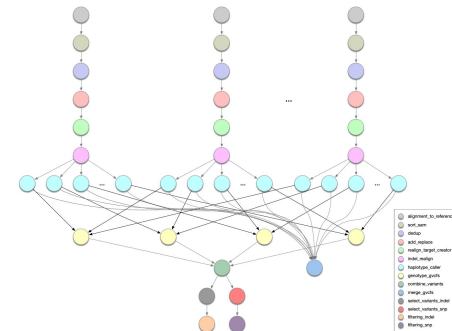
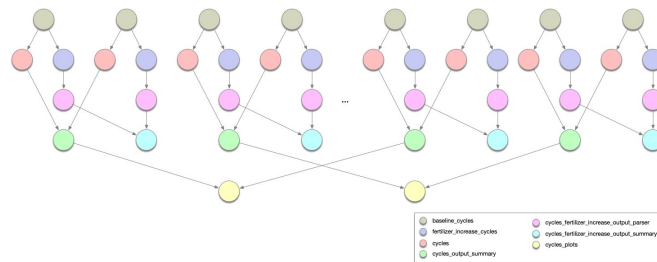


A task often represents a **program** (or script) written in any programming language (**closed box**)

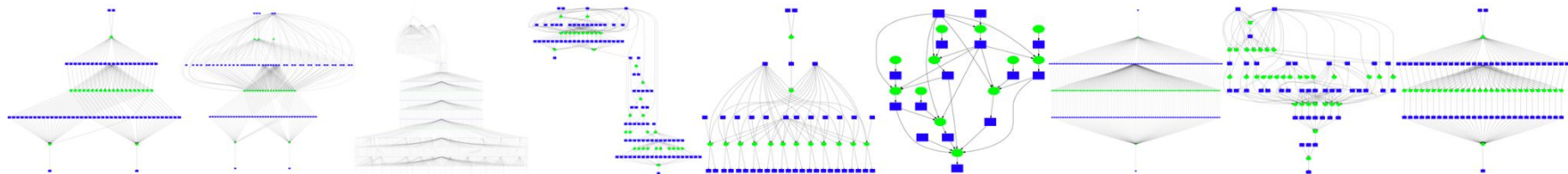
Dependencies are typically based on the **data flow**.

A dependency may also be expressed as **conditions**, **exceptions**, **user triggered action**, etc.

## A decorative graphic on the right side of the page. It features a vertical column of binary code (0s and 1s) in a light blue color. At the bottom of this column is a white hexagon with a blue outline.



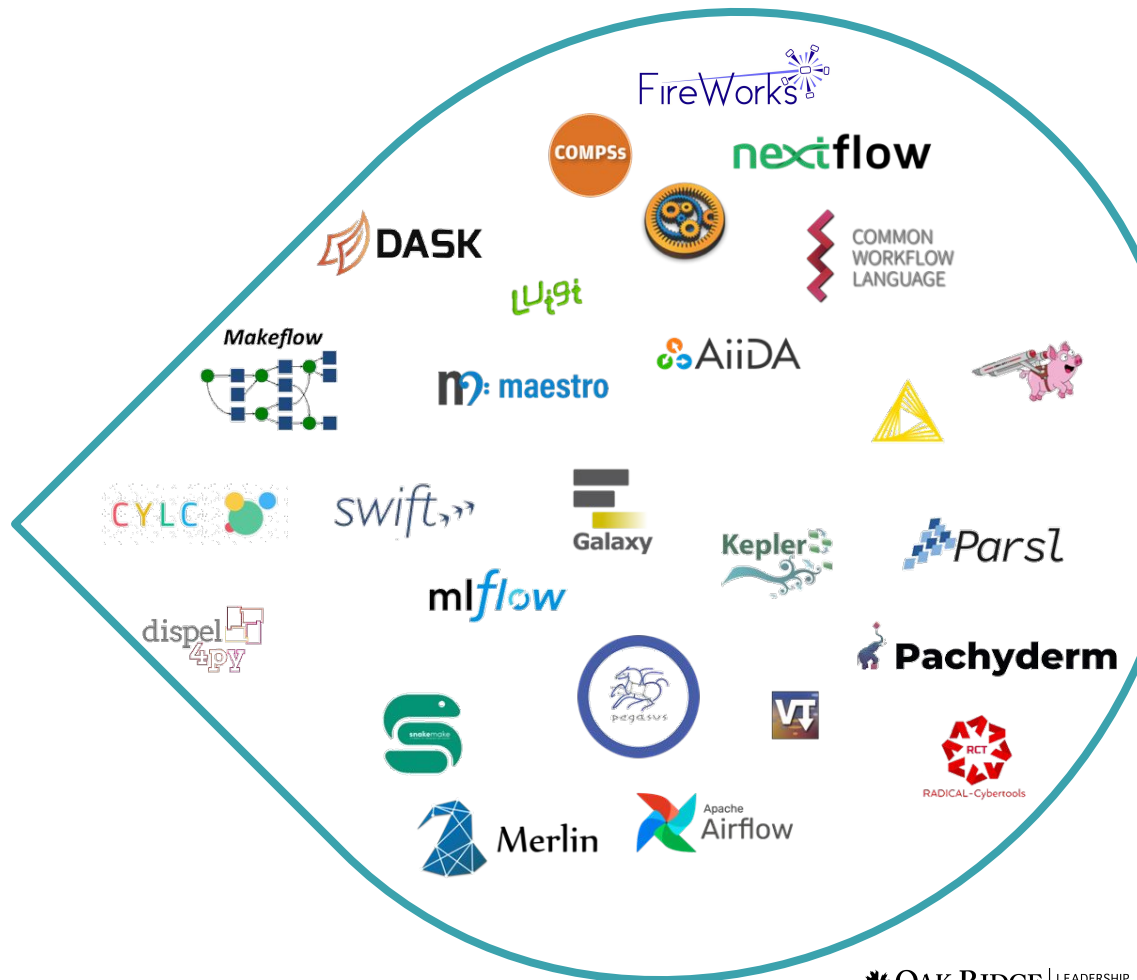
<https://github.com/wfcommons/pegasus-instances>



<https://github.com/cooperative-computing-lab/makeflow-examples>

# There is a myriad of workflow systems...

The workflow systems landscape is segmented and presents significant **barriers to entry** due to the hundreds of seemingly comparable, yet **incompatible**, systems that exist



<https://s.apache.org/existing-workflow-systems>

<https://github.com/pditommaso/awesome-pipeline>

# Characterization of Workflow Systems for Extreme-Scale Applications

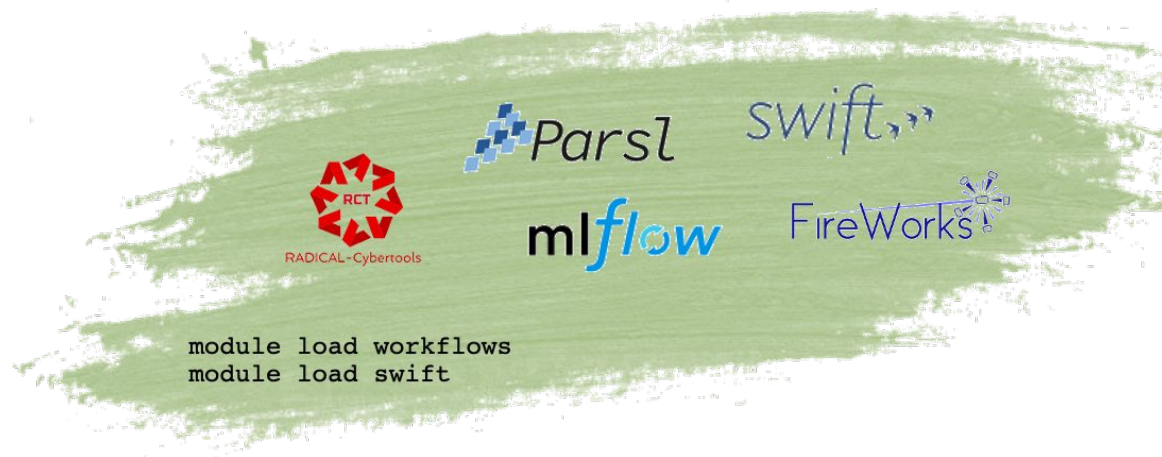
Workflow Properties	ADIOS	Airavata	Askalon	Bobolang	dispel4py	Fireworks	Galaxy	Kepler	Makeflow	Moteur	Nextflow	Pegasus	Swift	Taverna	Triana
<i>Workflow Execution Models</i>															
Sequential	✓	✓	✓	✗	✗	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓
Concurrent	✗	✗	✗	✓	✓	✗	✗	✗	✗	✗	✓	✗	✓	✗	✗
Iterative	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✗	✓	✓	✗
Tightly coupled	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
External steering	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✗	✗	✗	✗	✗
<i>Heterogeneous Computing Environments</i>															
Co-location	✗	✗	✗	✓	✓	✓	✗	✗	✗	✗	✓	✓	✓	✗	✗
External location	✗	✓	✓	✗	✗	✗	✓	✓	✓	✓	✗	✓	✓	✓	✓
In situ	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
<i>Data Access Methods</i>															
Memory	✓	✗	✗	✓	✓	✓	✗	✗	✗	✗	✓	✓	✓	✗	✗
Messages	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Local disk	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Shared file system	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Object store	✗	✗	✗	✗	✗	✗	✓	✗	✗	✗	✗	✓	✗	✓	✗
Other remote storage	✓	✗	✗	✗	✗	✗	✓	✗	✗	✗	✗	✓	✓	✓	✗

<https://doi.org/10.1016/j.future.2017.02.026>

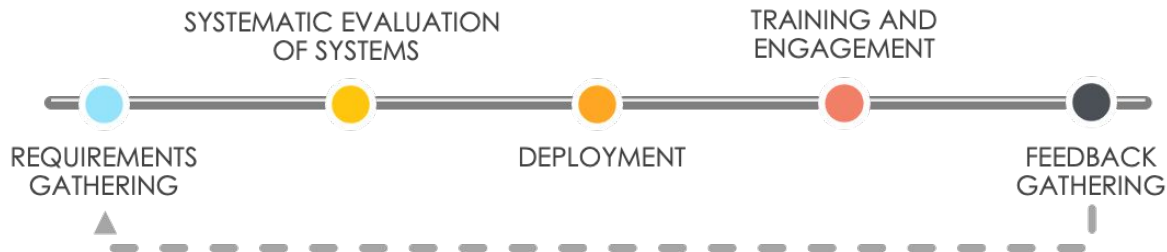


# Workflows@OLCF

We are constantly evaluating new systems and user requirements, and will deploy them as needed



## OLCF's process for deploying workflow systems



# Supported workflow systems at OLCF

- Why do we support multiple workflow systems?
  - Why do we support Fortran or Emacs? ;-)
- Important considerations:
  - Documentation
  - User community
  - “Paradigm”
  - Alignment/compatibility with your science and your tools



# Ensemble Toolkit (EnTK)

- Python-based
  - Workflows are Python programs that can manage external components.
- Launch workflows from login node with Python.
- Designed for ensemble-based applications such as
  - Molecular dynamics
  - Weather prediction models
- Depends on external services (MongoDB and RabbitMQ)

# MLflow

- Library agnostic
  - Workflows are Python, R, Java, or REST API programs that can manage external components.
- Launch workflows by submitting batch scripts.
- Designed with AI/ML workloads in mind
- Can be monitored using built-in web interface
- Lots and lots and lots of provenance tracking

# FireWorks

- CLI and Python API
  - Workflows are stored in a database.
  - Workflows can be defined with JSON or YAML files, or they can be Python programs that can manage external components.
- Launch workflows by submitting batch scripts.
- Can be monitored using built-in web interface
- Depends on an external service (MongoDB)

# Swift/T

- Swift language
  - Workflows are Swift programs that can manage external components.
- Launch workflows from login node with Swift.
- Designed to take advantage of MPI systems using Turbine and ADLB libraries
- Also available on Crusher and Andes

# Parsl

- Python-based
  - Workflows are Python programs that can manage external components.
- Launch workflows from login node with Python.
- Flexible enough to run “anywhere”
- Also available on Crusher and Andes

# Workflow Systems Modules on OLCF Summit

Workflow System	Module Command	Documentation / Examples
EnTK	module load workflows entk	radicalentk.readthedocs.io
MLflow	module load workflows mlflow	mlflow.org
FireWorks	module load workflows fireworks	materialsproject.github.io/fireworks
Swift*	module load workflows swift	swift-lang.github.io/swift-t/guide.html
Parsl*	module load workflows parsl	parsl-project.org

OLCF Documentation: [docs.olcf.ornl.gov/software/workflows](https://docs.olcf.ornl.gov/software/workflows)

\*ported and tested on Crusher and Andes



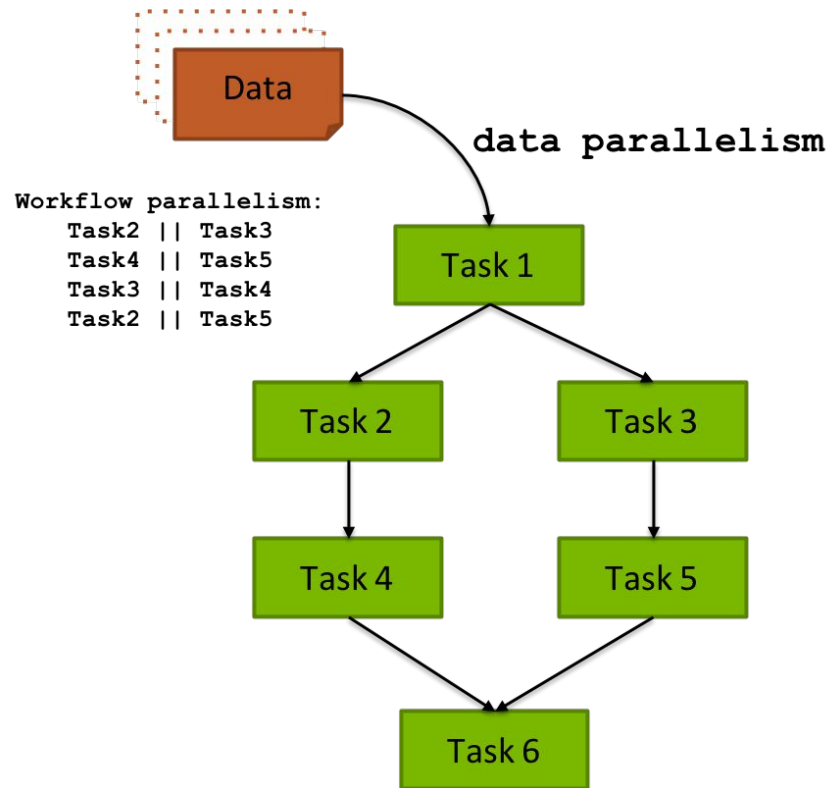
# A Quick Demo

Running "Hello World" with Parsl on Summit

Running a hypothetical "Crystal Workflow" with Swift/T on Summit

Running "Hello World" with MLFlow

Running "Hello World" with Fireworks



# Get in touch with us

As part of this deployment / support process, we would like to establish **close engagements** with users and applications

We kindly ask you to fill the following form so we can better plan our engagements:

<https://tinyurl.com/workflows-olcf>



Get in Touch

# Community Building

<https://workflows.community>

Provide a centralized source for  
resources, training,  
workshops, job  
opportunities, and news to  
scientists and developers working  
with workflows



managed by a nine-person leadership  
team, a steering committee, and a  
technical lead representing

25 companies and institutions from  
around the world

22

workflow  
systems

99

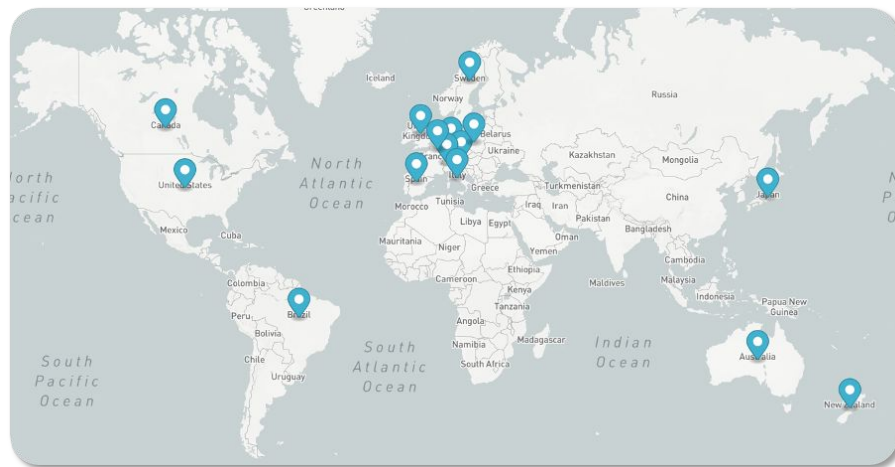
community  
members

5

working  
groups

4

research  
frameworks



# Community Summits

Over **100 participants** from a group of international researchers and developers, from **27 workflow systems** and user communities

## A Community Roadmap for Scientific Workflows Research and Development

Rafael Ferreira da Silva<sup>†</sup>, Henri Casanova<sup>‡</sup>, Kyle Chard<sup>§¶</sup>, Ilkay Altintas<sup>||</sup>, Rosa M Badia<sup>\*\*</sup>, Bartosz Balis<sup>††</sup>, Taina Coleman<sup>†</sup>, Frederik Coppens<sup>††</sup>, Frank Di Natale<sup>‡‡</sup>, Bjoern Enders<sup>§§</sup>, Thomas Fahringer<sup>§§</sup>, Rosa Filgueira<sup>||</sup>, Grigori Fursin<sup>§§</sup>, Daniel Garijo<sup>§§</sup>, Carole Goble<sup>§§</sup>, Dorran Howell<sup>§§</sup>, Shantenu Jha<sup>§§</sup>, Daniel S. Katz<sup>§§</sup>, Daniel Laney<sup>§§</sup>, Ulf Leser<sup>§§</sup>, Maciej Malawski<sup>||</sup>, Kshitij Mehta<sup>§§</sup>, Loic Potier<sup>||</sup>, Jonathan Ozik<sup>§§</sup>, J. Luc Peterson<sup>§§</sup>, Lavanya Ramakrishnan<sup>§§</sup>, Stian Soiland-Reyes<sup>§§</sup>, Douglas Thain<sup>§§</sup>, Matthew Wolf<sup>§§</sup>  
<sup>\*</sup>Oak Ridge National Laboratory, Oak Ridge, TN, USA <sup>†</sup>University of Southern California, Marina Del Rey, CA, USA  
<sup>‡</sup>University of Hawaii, Honolulu, HI, USA <sup>§</sup>Argonne National Laboratory, Lemont, IL, USA  
<sup>¶</sup>The University of Chicago, Chicago, IL, USA <sup>||</sup>University of California, San Diego, La Jolla, CA, USA  
<sup>\*\*</sup>Barcelona Supercomputing Center, Spain <sup>††</sup>AGH University of Science and Technology, Krakow, Poland  
<sup>‡‡</sup>Ghent University, Ghent, Belgium <sup>§§</sup>VIB Center for Plant Systems Biology, Belgium  
<sup>§§</sup>Lawrence Livermore National Lab, Livermore, CA, USA <sup>§§</sup>University of Innsbruck, Innsbruck, Austria  
<sup>§§</sup>Heriot-Watt University, Edinburgh, UK <sup>§§</sup>OctoML, USA <sup>§§</sup>Universidad Politécnica de Madrid, Spain  
<sup>§§</sup>The University of Manchester, Manchester, UK <sup>§§</sup>Tweag, Zürich, Switzerland  
<sup>§§</sup>Brookhaven National Laboratory, Upton, NY, 11973 <sup>§§</sup>University of Illinois at Urbana-Champaign, USA  
<sup>§§</sup>Humboldt-Universität zu Berlin, Berlin, Germany <sup>§§</sup>Lawrence Berkeley National Lab, Berkeley, CA, USA  
<sup>§§</sup>University of Amsterdam, Amsterdam, The Netherlands <sup>§§</sup>University of Notre Dame, Indiana, USA



arXiv:2110.02168

<https://doi.org/10.1109/WORKS54523.2021.00016>

# Automating Science with Workflows at OLCF

**Ketan Maheshwari**  
**Sean Wilkinson**  
**Rafael Ferreira da Silva**

Data Lifecycle & Scalable Workflows  
National Center for Computational Sciences



**Thank you!**  
**Questions?**

<https://tinyurl.com/workflows-olcf>