

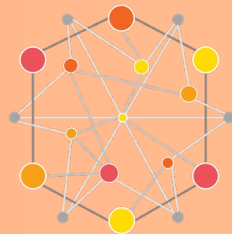
Exascale Workflows Community: A Post-ECP Roadmap for Workflow Systems and Applications

Rafael Ferreira da Silva – ORNL

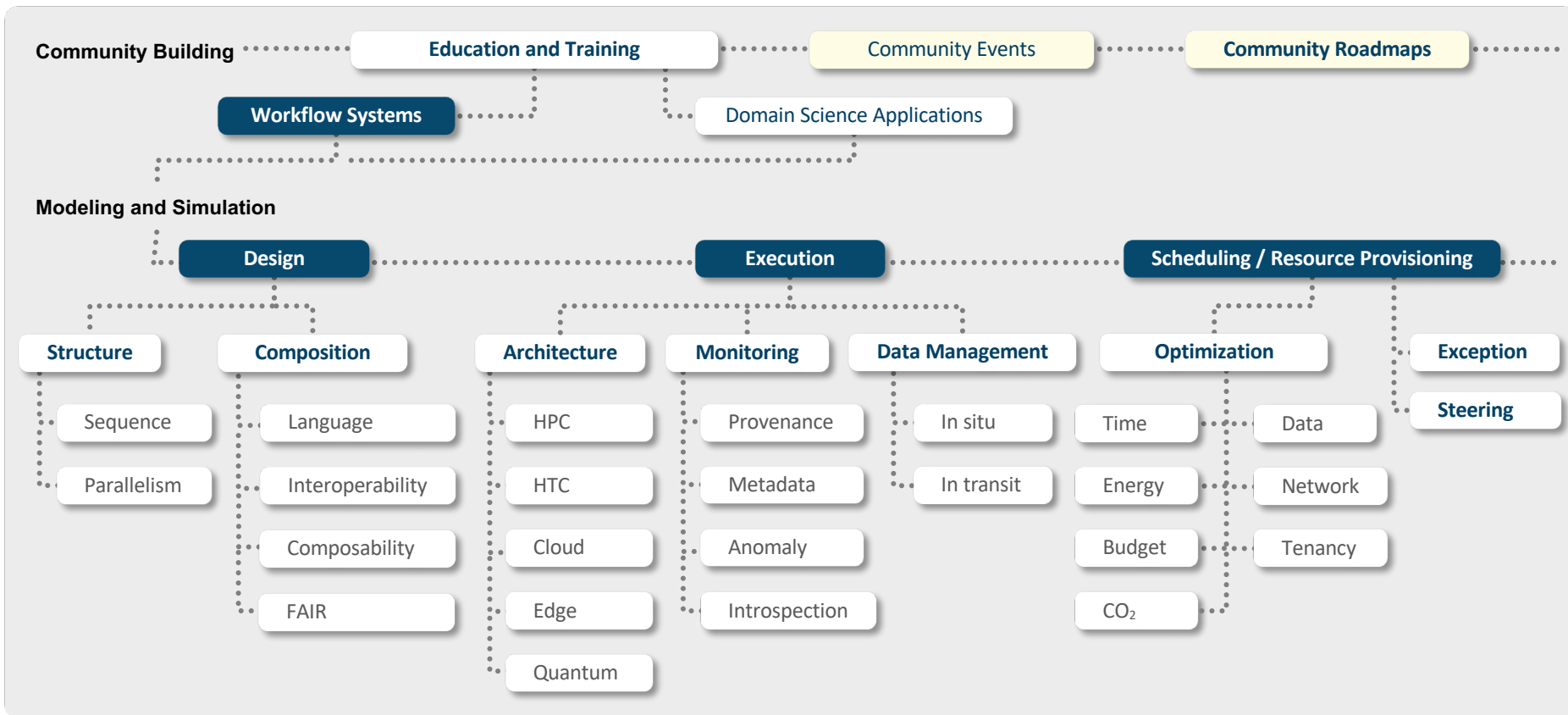
Shantenu Jha – BNL

Daniel Laney – LLNL

Kyle Chard – ANL



Overview of Scientific Workflows Research Challenges



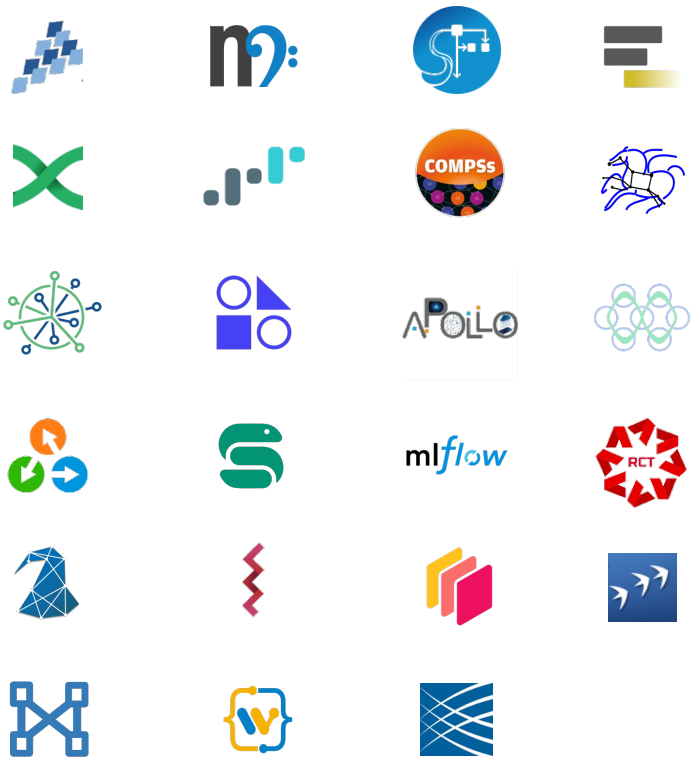
The Universe of Workflow Systems

Forcing applications to fit a particular system

Needs a developer to be “shipped” with the workflow system

The term “workflow” is overloaded

<https://s.apache.org/existing-workflow-systems>
324 workflow systems as of Dec 12, 2022



The 2021 Community Roadmap



arXiv.org > cs > arXiv:2110.02168

Search...
Help | Advance

Computer Science > Distributed, Parallel, and Cluster Computing

[Submitted on 5 Oct 2021 (v1), last revised 8 Oct 2021 (this version, v2)]

A Community Roadmap for Scientific Workflows Research and Development

Rafael Ferreira da Silva, Henri Casanova, Kyle Chard, Ilkay Altintas, Rosa M Badia, Bartosz Balis, Tainã Coleman, Frederik Coppens, Frank Di Natale, Bjoern Enders, Thomas Fahringer, Rosa Filgueira, Grigori Fursin, Daniel Garijo, Carole Goble, Dorrán Howell, Shantenu Jha, Daniel S. Katz, Daniel Laney, Ulf Leser, Maciej Malawski, Kshitij Mehta, Loïc Pottier, Jonathan Ozik, J. Luc Peterson, Lavanya Ramakrishnan, Stian Soiland-Reyes, Douglas Thain, Matthew Wolf

The landscape of workflow systems for scientific applications is notoriously convoluted with hundreds of seemingly equivalent workflow systems, many isolated research claims, and a steep learning curve. To address some of these challenges and lay the groundwork for transforming workflows research and development, the WorkflowsRI and ExaWorks projects partnered to bring the international workflows community together. This paper reports on discussions and findings from two virtual "Workflows Community Summits" (January and April, 2021). The overarching goals of these workshops were to develop a view of the state of the art, identify crucial research challenges in the workflows community, articulate a vision for potential community efforts, and discuss technical approaches for realizing this vision. To this end, participants identified six broad themes: FAIR computational workflows; AI workflows; exascale challenges; APIs, interoperability, reuse, and standards; training and education; and building a workflows community. We summarize discussions and recommendations for each of these themes.

Comments: arXiv admin note: substantial text overlap with arXiv:2103.09181
Subjects: **Distributed, Parallel, and Cluster Computing (cs.DC)**
Cite as: arXiv:2110.02168 [cs.DC]
(or arXiv:2110.02168v2 [cs.DC] for this version)

We summarize the discussions and findings by presenting a consolidated view of the **state of the art, challenges**, and potential efforts, which we eventually synthesize into a **community roadmap**

<https://arxiv.org/abs/2110.02168>

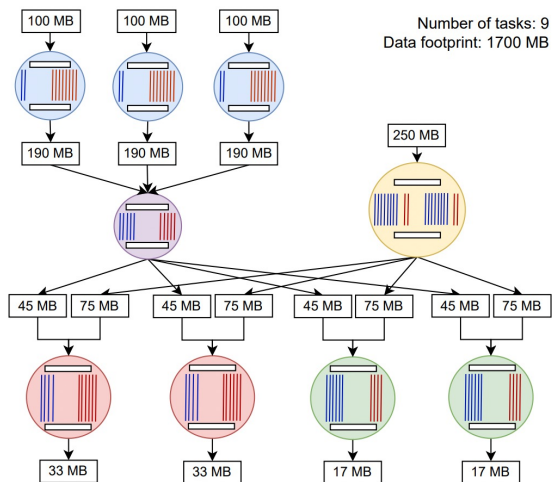
National Academy of Sciences, Engineering, and Medicine Report



*The needs and demands placed on science to address a range of urgent problems are growing. The world is faced with complex, interrelated challenges in which the way forward lies hidden or dispersed across disciplines and organizations. For centuries, scientific research has progressed through iteration of a **workflow built on experimentation or observation and analysis of the resulting data**. While computers and automation technologies have played a central role in research workflows for decades to acquire, process, and analyze data, these same computing and **automation technologies can now also control the acquisition of data**, for example, through the design of new experiments or decision making about new observations.*

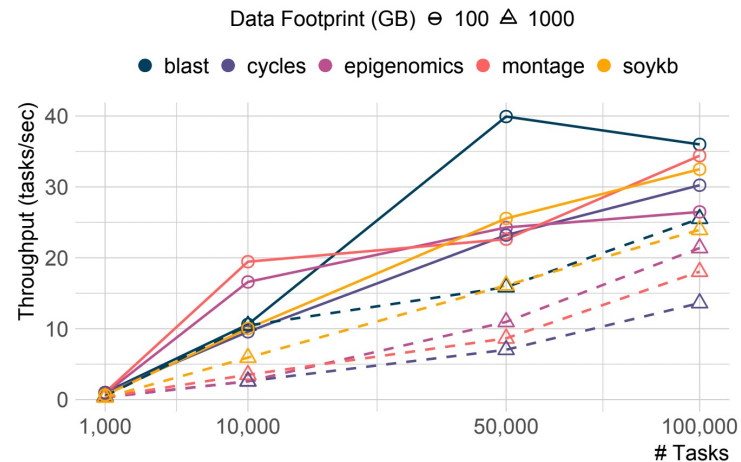
<https://doi.org/10.17226/26532>

WfBench: Workflow Benchmarks



Representative
tasks and
workflow
benchmarks

Analysis of
workflow system
overhead



↑ data footprint ↓ throughput
↑ #tasks ↑ throughput

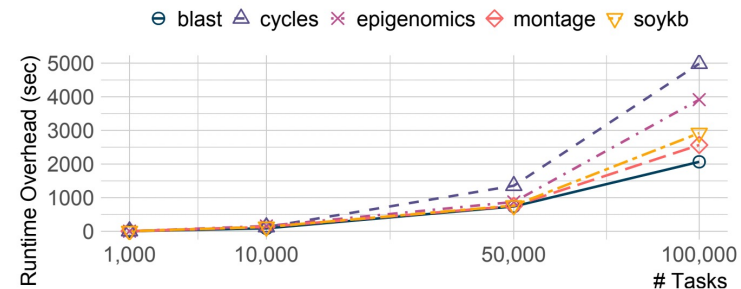


Fig. 10. Workflow execution time (or total overhead) vs. number of tasks.

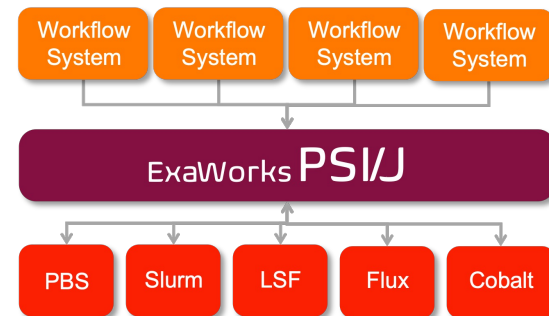
PSI/J, a Portable Submission interface for jobs

Motivated by a sequence of **interviews and community meetings** identifying job launch/management as low hanging fruit

Community generated a **light-weight user-space API Specification** via a public process hosted on GitHub

ECP ExaWorks team created an Initial **Python reference implementation**

ECP ExaWorks creating a public dashboard and CI infrastructure for multiple DOE compute centers that could be extended by community



```

import jpsi

jex = jpsi.JobExecutor.get_instance('slurm')

def make_job():
    job = jpsi.Job()
    spec = jpsi.JobSpec()
    spec.executable = '/bin/sleep'
    spec.arguments = ['10']
    job.spec = spec
    return job

jobs = []
for i in range(N):
    job = make_job()
    jobs.append(job)
    jex.submit(job)

for i in range(N):
    jobs[i].wait()
  
```

Workflows Community Initiative

The goal of the *Workflows Community Initiative (WCI)* is to bring the workflows community together (**users, developers, researchers, and facilities**) to provide community resources and capabilities to enable scientists and workflow systems developers to discover software products, related efforts, events, technical reports, etc. and engage in community-wide efforts to tackle workflows grand challenges.

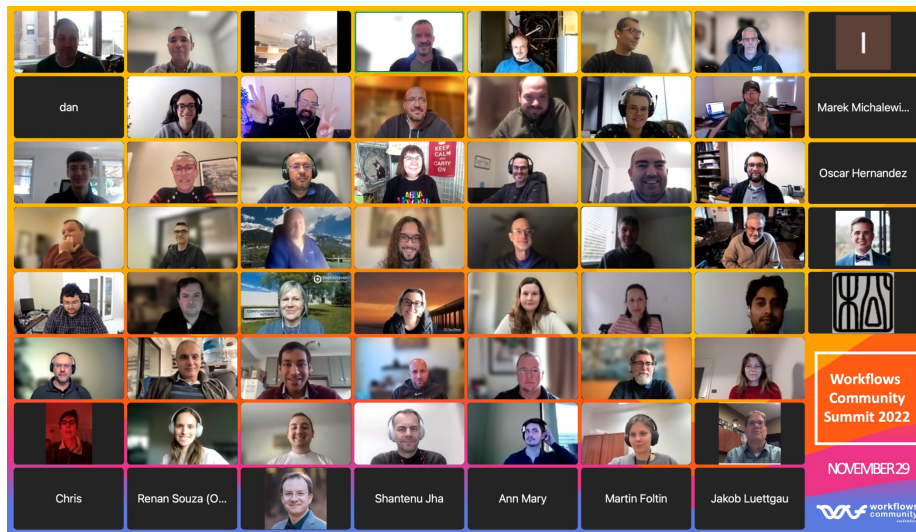


We are always looking for volunteers...

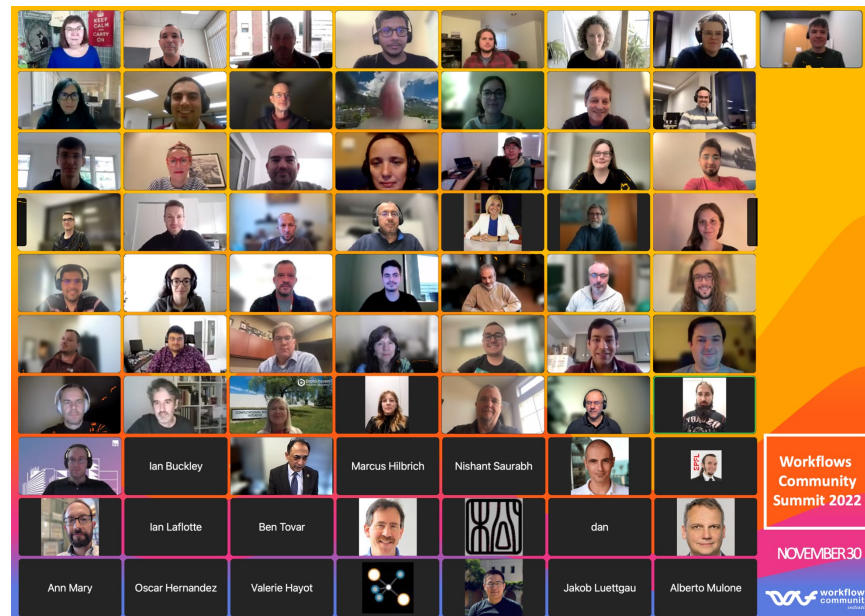
- 24 workflow systems
- 148 community members
- 6 working groups
- 4 research frameworks

Workflows Community Summit 2022

Nov 29 and 30, 2022



75+ international participants



Revised Roadmap: Initial Thoughts

Specifications, standards, and APIs

Standards are often constrained or **hard to implement**

Define **common terms**, building blocks, and concepts

Define some standard for giving **data** from user to workflow systems/operators

AI workflows

Unreliability of models – training processes often have a **human-in-the-loop**

Types of workflows – *workflows to **create AI*** and *workflows that **use AI***

Challenges: **random access** to datasets in training and management of **small files**

High performance data management and in situ workflows

Edge to cloud continuum and data exchange through data objects

Data usage is more **fine-grain** than a typical HPC code

Adaptive **compression** to reduce data necessary to represent the problem domain

Revised Roadmap: Initial Thoughts

HPC and Quantum workflows

Community do not know how to transfer information to the QC system
Limited resources (hard to access), long queues & expensive
Heterogeneity in quantum devices (vendor specific APIs)

FAIR workflows

Standards for expressing the **inputs** of the workflow
Limited availability of **metadata**
FAIR data and FAIR workflows are **intertwined**

Workflows for continuum and cross-facility computing

Describe **aggregate I/O** needs of a workflow
Coordinate **communication** between sites (different **security** domains)
People who design **experiment facilities** are not necessarily computing experts

Panel

Panelists / Questions



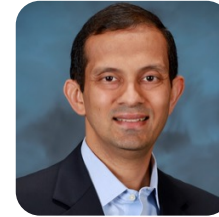
Katie Antypas
NERSC/LBNL



Todd Gamblin
LLNL



Andrew Gallo
GE Research



Arjun Shankar
ORNL



Shantenu Jha
BNL (moderator)

In the era of exascale and ML/AI, what are the emerging and future crucial challenges for post-ECP workflows?

Considering workflows sustainability, what are the key constraints and opportunities to attain sustainability?

Which software/technology are essential for enabling sustainable workflows post-ECP?



slido.com
#1557 064

Thank you!