

Special Issue Paper

International Journal of HIGH PERFORMANCE COMPUTING APPLICATIONS

Orchestration of materials science workflows for heterogeneous resources at large scale

The International Journal of High Performance Computing Applications 2023, Vol. 37(3-4) 260–271 © The Author(s) 2023 Article reuse guidelines: sagepub.com/journals-permissions DOI: 10.1177/10943420231167800 journals.sagepub.com/home/hpc

S Sage

Naweiluo Zhou¹, Giorgio Scorzelli², Jakob Luettgau¹, Rahul R Kancharla³, Joshua J Kane³, Robert Wheeler⁴, Brendan P Croom⁵, Pania Newell², Valerio Pascucci² and Michela Taufer¹

Abstract

In the era of big data, materials science workflows need to handle large-scale data distribution, storage, and computation. Any of these areas can become a performance bottleneck. We present a framework for analyzing internal material structures (e.g., cracks) to mitigate these bottlenecks. We demonstrate the effectiveness of our framework for a workflow performing synchrotron X-ray computed tomography reconstruction and segmentation of a silica-based structure. Our framework provides a cloud-based, cutting-edge solution to challenges such as growing intermediate and output data and heavy resource demands during image reconstruction and segmentation. Specifically, our framework efficiently manages data storage, scaling up compute resources on the cloud. The multi-layer software structure of our framework includes three layers. A top layer uses Jupyter notebooks and serves as the user interface. A middle layer uses Ansible for resource deployment and managing the execution environment. A low layer is dedicated to resource management and provides resource management and job scheduling on heterogeneous nodes (i.e., GPU and CPU). At the core of this layer, Kubernetes supports resource management, and Dask enables large-scale job scheduling for heterogeneous resources. The broader impact of our work is four-fold: through our framework, we hide the complexity of the cloud's software stack to the user who otherwise is required to have expertise in cloud technologies; we manage job scheduling efficiently and in a scalable manner; we enable resource elasticity and workflow orchestration at a large scale; and we facilitate moving the study of nonporous structures, which has wide applications in engineering and scientific fields, to the cloud. While we demonstrate the capability of our framework for a specific materials science application, it can be adapted for other applications and domains because of its modular, multi-layer architecture.

Keywords

Scientific workflow, orchestration, resource management, job scheduling, materials science, cloud computing

Motivation and proposed solution

Materials science has evolved to adopt sophisticated scientific workflows to automate experiments and data analysis. Materials science workflows can express many applications, often consisting of multiple tasks that require resources beyond a single machine (Wu et al., 2015). A workflow usually comprises several stages, executing a task and taking input data from experimental tools or previous steps. Materials science workflows can generate and stream massive amounts of data in the order of terabytes or petabytes. Examples of these workflows include analyzing and visualizing internal material structure by performing synchrotron X-ray computed tomography reconstruction and segmentation of silica. A personal computer or a small

dedicated cluster is no longer sufficient to perform the reconstruction and segmentation operations on tomography data. Public cloud platforms can effectively provide

Corresponding author:

Michela Taufer, EECS, University of Tennessee Knoxville College of Engineering, 401 Min H. Kao Bldg. 1520 Middle Drive, Knoxville, TN 37996, USA.

Email: taufer@utk.edu

¹University of Tennessee Knoxville College of Engineering, Knoxville, TN, USA

²University of Utah, Salt Lake, Utah, USA

³Idaho National Laboratory, Idaho Falls, USA

⁴MicroTesting Solutions LLC, Hilliard, OH, USA

⁵Johns Hopkins University Applied Physics Laboratory, Laurel, MD, USA

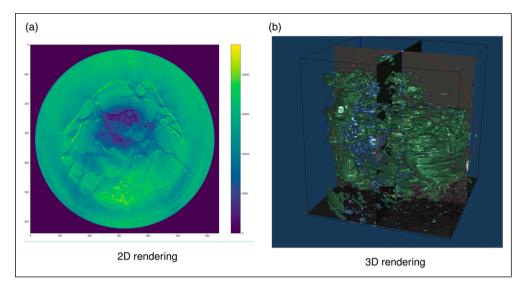


Figure 1. A single slice of the materials in the Z-axis view after reconstruction in 2D rendering (a) and the image with enhanced contrast and removal of artifacts after segmentation viewed in 3D rendering (b). Note that the images at both stages can be 3D rendered.

computing and storage resources, as users do not need to deploy an entire infrastructure that brings additional costs such as hardware purchase and maintenance (Cepuc et al., 2020). Furthermore, cloud resources are provisioned based on user demands with a pay-as-you-go model and can quickly adapt program execution environments to user needs. Still, efficiently deploying services on the cloud requires technical expertise.

In this work, we break through the expertise barrier by providing scientists in materials science with a framework that enables a scalable, easy use of the cloud for their studies. We leverage cutting-edge cloud technologies to design and implement our framework. To this end, our framework integrates a rich suite of software features needed for managing the heterogeneous hardware infrastructure (CPUs, GPUs, and storage) on the cloud and orchestrating the tasks composing the multi-stage materials science workflow, removing the need for a deep understanding of the cloud infrastructure or deploying workflow management systems (Deelman et al., 2015; Jain et al., 2015; Wilde et al., 2011). Furthermore, our framework enables interactive, visual exploration of large materials science data. Our framework comprises three layers: a top user interface, a middle layer for resource deployment, and a lower layer dedicated to resource management.

Using synchrotron X-ray computed tomography, we deploy our framework to study in situ mechanical testing of a nanoporous structure. Such study has various applications in engineering and scientific domains (from fillers and additives to adsorbents, catalysts, and recently therapeutic agents and vaccines in nanomedicine). The workflow studying the phenomena comprises six stages: downloading images from X-ray tomography scans (6.0 GB per image);

performing image segmentation to remove noises; performing image reconstruction to obtain a clear vision of the material structure; converting the data into OpenVisus formats Petruzza et al. (2020); streaming data to storage; and using the Open-Visus streaming services for user analysis. Figure 1 shows an example of the internal structure of a porous silica material at two different stages. Figure 1(a) shows a single slice of the materials in the Z-axis view in 2D rendering after reconstruction. Figure 1(b) shows the enhanced contrast and removal of artifacts after segmentation in 3D rendering.

We show how our framework provides a scalable solution to the study of the internal structure of a porous silica material in which an increase in data size requires storage resource elasticity provided by the cloud technologies building our framework. The contributions of this article are as follows:

- We present our multi-layer framework that provides a complete life cycle for a materials science workflow from its environment provision to final stage data analysis using cutting-edge cloud technologies.
- We use our framework to study the in situ mechanical testing of a nanoporous structure; nonporous structures have wide applications in engineering and scientific fields.
- We demonstrate the ability of our framework to handle large amounts of data at a magnitude of 10 TBytes in a scalable manner using a public cloud platform such as Chameleon (Keahey et al., 2020).

Although the framework capabilities are demonstrated for a specific study in the materials science domain, it can be easily adapted to support multi-stage, modular

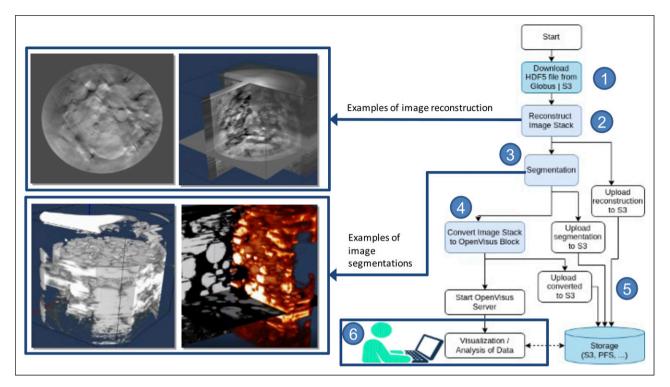


Figure 2. The six stages of a materials science workflow for in situ mechanical testing of a nanoporous structure using synchrotron X-ray computed tomography. Numbered circles identify the six stages. Examples of intermediate images from the reconstruction and segmentation stages are on the figure's left.

workflows in other domains (Singh et al., 2019). The modular implementation of our framework allows domain scientists to easily plug in their software modules as new workflow stages and run their workflow with the framework's Jupyter notebook interface.

Critical problem in materials science

Using synchrotron X-ray computed tomography, we leverage our framework to study a critical data and computation-intensive problem in materials science: the in situ mechanical testing of a nanoporous foam. Porous silica-based materials have wide applications in engineering and scientific domains, from fillers and additives to adsorbents, catalysts, and recently therapeutic agents and vaccines in nanomedicine (Lei et al., 2020). The internal structure of porous silica-based materials can be reconstructed from synchrotron X-ray tomography. Tomography combines the projections of a sample slice at different orientations taken by placing the slice between an X-ray source and detector (Maire et al., 2001). Analysis of the in situ tomography images can reveal the underlying mechanisms that drive macroscopic mechanical performance (Peukert et al., 2015), such as the collapse of pores, cracking and buckling of ligaments, and enable the design of higher-performing silica foams.

The tomography images obtained from X-ray scans require image processing to perform synchrotron X-ray computed tomography reconstruction and segmentation of silica. Reconstruction and segmentation remove blurriness, calculate the tomographic center of rotation, and eliminate scanning artifacts. Figure 2 illustrates the complete experiment process represented in a scientific workflow. The method comprises six stages performed on a sequence of images: downloading the images of X-ray tomography (Stage 1); performing the image reconstruction to obtain a clear vision of the materials structures (Stage 2); performing the image segmentation to remove background noise, enhance contrast, and reduce image size (Stage 3); converting the data into a networkfriendly multi-level format (Stage 4); streaming data to users (Stage 5); and visualizing and analyzing data with a Jupyter notebook (Stage 6) (Randles et al., 2017). Specifically, each X-ray tomography image is approximately 6.0 GB; the images are downloaded in HDF5 file Group (2000-2010). The stages of performing image reconstruction and segmentation (i.e., Stages 2 and 3) are timeconsuming, requiring massive computing resources. Intermediate stages (i.e., Stages 2, 3, and 4) yield a large amount of intermediate data reaching at least 30 times the input image size. Figure 2 shows examples of intermediate images after the segmentation and reconstruction

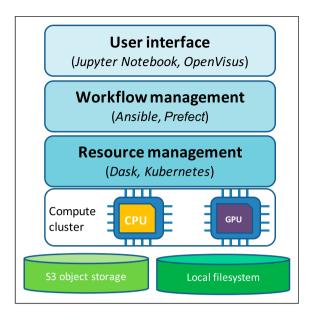


Figure 3. The software layers of our framework on top of the hardware infrastructure.

stages (images on the left). Preserving the intermediate data generated by each stage facilitates scientists to track data provenance (Wang et al., 2015), enabling the reusability of workflow building blocks and the reproducibility and trustworthiness of the results (Taufer et al., 2021a,b); Stodden et al., 2016).

Conceptual design of our framework

We need a software solution that works on cloud platforms to address the challenges of materials science workflows. Our software solution combines workflow management and resource provisioning with the concurrent execution of the workflow stages using the task parallelism paradigm. Data is streamed with the support of orchestration services from the object storage to the user interface in an interoperable, portable, and general manner.

Software solution and hardware infrastructures. Our software solution is a multi-layer framework composed of three layers: a top user interface, a middle layer for resource deployment, and a lower layer dedicated to resource management. Figure 3 illustrates the multi-layer organization of our framework.

Jupyter is a free, open-source, interactive tool that researchers can use to combine software code, computational output, explanatory text, and multimedia resources in a single document. A Jupyter notebook interface serves as the top layer and facilitates the domain scientist's interaction with the workflows, the interpretation of its execution, and the explainability of its results. We developed ad hoc notebooks that leverage the

OpenVisus Software Development Kit (SDK) (Pascucci et al., 2012) to:

- visualize framework outputs (i.e., 2D and 3D images);
- perform fast visualization of massive data; and
- run visualization services agnostic from the underneath hardware.

The middle layer uses the declarative markup language YAML to configure, describe, and run all the workflow stages modular and adaptable, making new stages easily identifiable and pluggable into the original workflow. For the setup, we run our configuration scripts on the Ansible (Sammons, 2016) platform, an IT automation system that handles configuration management, application deployment, cloud provisioning, ad hoc task execution, and Infrastructure-As-Code (IaC) through a suite of Python modules. For the material science part, we develop a domain-specific and modern software stack that provides services for high-level task definition; low-latency task coordination and communications (our dependencies are expressed as a parameterized Directed Acyclic Graph); automatic failure handling; and automatic data caching, central logging, observability tools.

The lower layer manages heterogeneous computing resources (i.e., GPU and CPU nodes) and schedules different jobs on these resources. The core of this layer uses Kubernetes and Dask: Kubernetes (Hightower et al., 2017) is an open-source system for automating deployment, scaling, and management of containerized applications; and Dask (Rocklin, 2015) is a flexible open-source Python library for parallel computing, scaling Python code from multi-core local machines to large distributed clusters or large-scale distributed cloud platform. For this layer, we intentionally architecture a hybrid resource management system (see Figure 4) where Dask performs large-scale job scheduling and resource management for GPU-intensive tasks, and Kubernetes orchestrates long-running CPU-based services such as the data-streaming services. We leverage three types of gateways (also known as entry points or control nodes) that provide indirect, monitored, and secured access:

- An Ansible Control Node orchestrates and performs cluster-wide deployments;
- A Kubernetes Login Node dispatches the commands to the Pods (i.e., atomic unit on the Kubernetes platform); and
- A Jupyter Serve Node launches on-demand multiuser Jupyter notebooks allowing interactive analysis.

The storage layer at the bottom of our software stack includes temporary local file storage and remote object

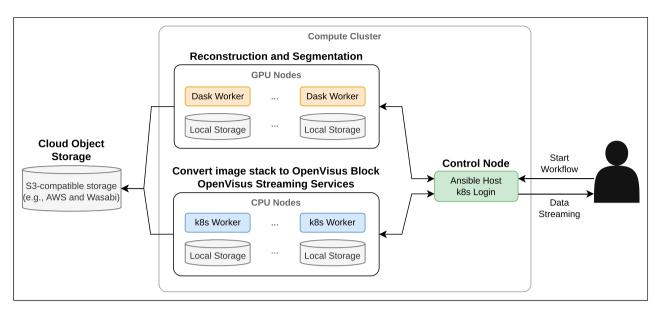


Figure 4. Architecture of our framework. The number of nodes in the figure is only for illustration.

storage. The local file system is used to download locally remote artifacts, such as the HDF5 files produced by the acquisition devices, temporary store artifacts produced by the reconstruction and segmentation algorithms, and support the artifacts conversion to Analysys-Ready Network Optimized (ARCO) file formats (Abernathey et al., 2021). Since most of the mentioned operations require random file access and low latency, we need a small amount of fast local storage in order of hundreds or thousands of Input Output operations per second (IOPS). The remote object storage, leveraging the S3 Application Programming Interface (API), allows the permanent storage of the results and their distribution and replication at the national level. Our framework is designed from the ground base to be compatible with any commercial or educational object storage solution. We chose Wasabi for this particular workflow, a vendor with a transparent policy of zero-EGRESS fees (i.e., zero cost for moving data outside) and zero API fees. Using Wasabi, we pay for the pure storage costs without incurring high and difficult-to-predict upload and download charges (Penny, 2020).

Workflow management and resource provisioning. To automate the deployment of our software stack, our framework solution uses Ansible (Sammons, 2016) that serves as an orchestration tool by offering an easy way to perform cluster-wide operations, such as configuration management, application deployment, cloud provisioning, ad hoc task execution, and multi-node orchestration. Ansible deploys thereof via its playbook written in yaml that describes the automation of jobs or tasks. yaml* (en-Kiki et al., 2009) is a high-level human-readable language in

```
1 # where to store local files
2 LOCAL: /tmp/in_situ_data
3 # remote address where to store results
4 REMOTE: s3://in_situ_data
6 #ansible inventory location
7 inventory: /ansible/inventory.ini
8 #the list of the hosts defined in inventory ini
9 group: dask_node_list
10 #dask worker directory
11 worker-local-dir: /tmp/nsdf/dask-workers
12
13 # what task to execute (preprocess | convert)
14 task: preprocess
16 # specify which stage of the workflow to execute
17 disable-reconstruction: true
18 disable-segmentation: false
20 #specify which file on s3 to download
21 - url: s3://in_situ_data/filename.h5
```

Figure 5. A snippet of the yaml script setting configurations for the workflow. The preprocess includes two stages: reconstruction (Line 17) and segmentation (Line 18). For simplicity, only a part of the code is shown herein.

configuration files. Kubernetes also utilizes yaml scripts to perform the deployments. Connections to nodes in Ansible are via ssh, and nodes managed by Ansible are grouped into an inventory file defined by users.

Ansible, our yaml scripts, and our Python scripts are essential in workflow executions. Ansible performs clusterwide software deployments such as:

 Provisioning execution environments by installing and switching the list of python libraries used at different stages of the workflow;

```
1 @task
2 def Reconstruction(d):
3 return CallReconstruction(d)
4 @task
5 def Segmentation(d):
6 return CallSegmentation(d)
```

Figure 6. A snippet of stage operations with parallelized task indicated by @task decorator. Some parts of the code are omitted for simplicity.

Deploying the Dask and Kubernetes software stacks.

Our yaml script defines all the configurations for our workflow (Figure 5). The script has these key steps:

- Give the local location to store results as shown in Line 2 of Figure 5, and the remote location on s3 to archive all the data as shown in Line 4. The link in Line 4 gives the bucket name of our object storage (Factor et al., 2005) on s3.
- Set the environment variables such as the Dask worker directory (Line 11). The environment settings are applied for all the Dask nodes.
- Define the stages of our workflow that scientists execute on the cloud. As an example, in our script shown in Figure 5, scientists insert the names of Stages 2, 3, and 4 representing the Python module names. Note that Stages 2 and 3 in this script are incorporated under one component Preprocess, as both stages perform heavy computation. The stages are scheduled to run on the Dask cluster. Convert represents Stage 4 in Figure 2. In this way, scientists only need to bring their Python modules that focus on the scientific functionalities, and our framework handles the rest (e.g., job scheduling or task parallelization); and
- Specify the series of input files. Line 21 gives the location of an input image filename.h5 on s3.

Our Python scripts pass the configurations defined in the yaml script to the layers of our framework.

Task parallelization. To execute workflows effectively, we introduce task parallelism in our framework. Figure 6 shows a simplified snapshot of the parallelization. The unit of parallelism is a stage of the workflow as shown in Figure 2; each unit is identified by the standard Python decorator @task (Line 1, Line 4, and Line 7). The defined tasks are formed into a pipeline (represented as a Prefect Flow) scheduled to GPU nodes by Dask. We use Prefect as our workflow management system on top of Dask.

```
1 # import Openvisus package
2 import OpenVisus as ov
3 # obtain the image named id_1265
4 im =ov.LoadDataset(
'http://xx.xx.xx.xx:32439/mod_visus?dataset=id_1265')
```

Figure 7. A code snippet on how to fetch the streamed dataset in a Jupyter notebook. The real IP address has been replaced by xx.xx.xx, and the port number, in this case, is 32,439.

Dask handles scheduling and resource management of tasks within a workflow. The Dask cluster consists of one scheduler and many worker nodes. The first node defined in the Ansible inventory is elected as the scheduler and the rest act as workers. A defined task is executed in parallel on all the worker nodes. Additionally, our framework records logs of task status on each Dask worker and passes them to the Dask scheduler during execution. The logs can be disabled to reduce communication costs between nodes.

Streaming service orchestration. To support streaming services on multiple platforms, the entire software stack of our framework is encapsulated inside a Docker container (Merkel, 2014), with the data being mounted inside the containers. The software stack comprising the three layers (Figure 3) is a collection of tools that are well-tested and have been used for general purposes: Kubernetes as a container orchestrator. Ansible as a cluster-wide deployment tool, OpenVisus as a data visualization tool. Kubernetes schedules the Docker instances and automatically relaunches them when failures occur. The data is exposed externally via Kubernetes services associated with a global IP address and a port number. Multiple users can access the data via a given link. Figure 7 illustrates how a user can access the dataset id 1265 inside a Jupyter notebook. The streaming service component in our architecture can be decoupled from the framework and deployed on a more permanent yet smaller cluster that can continuously provide long-running streaming services within the budget of small or median-size research laboratories or institutes.

Local filesystem and remote object storage. To manage the large volume of intermediate data, each computing node on the Dask and Kubernetes clusters has local storage (i.e., local filesystem) attached. The cloud object store s3 is accessible to all nodes on the compute cluster as illustrated in Figure 4. Data accessing to s3 is much slower than the local store. Hence s3 only serves as a complementary data store. The local filesystem enables quick data access during workflow execution. However, it has limited storage capacity and can only provide temporal storage. The hybrid

Table 1. Hardware specifications of the cloud infrastructures used for scalability evaluation. The Dask and Kubernetes clusters used for our experiments comprise three CPU and 16 GPU nodes.

CPU node	Num. cores 24	Memory size 187.0 GB	CPU vendor Intel Xeon Gold 6126, 2.60 GHz	Local storage 441 GB
GPU node	Num. GPU I	Memory size 23.4 GB	GPU vendor Quadro RTX 6000	Local storage 210 GB

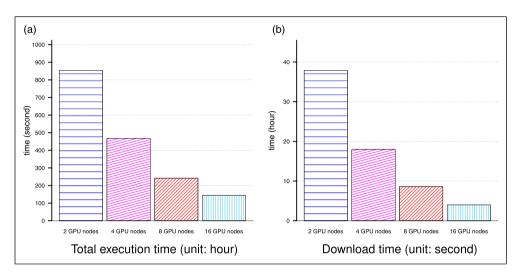


Figure 8. Scalability study of the two computation-intensive stages (i.e., Stages 2 and 3) in terms of workflow's total execution time (in hours) (a) and download time per node (in seconds) (b) on Chameleon.

data store allows our framework to handle a large volume of data during workflow execution and also offers long-term storage for data retrieval used by the Kubernetes streaming services.

The input HDF5 files are downloaded from s3 as shown on the left side of Figure 4 to the local disk. The desired images to be streamed are also fetched from s3. The intermediate data generated during execution on the local filesystem is freed once backed up to s3, performed in the background by multiple threads during computation.

Interoperability, portability, and generalizability. To sustain interoperable workflow executions, our software framework includes a Jupyter notebook as a user interface. Furthermore, the deployment layer provides portability, enabling our framework's application to run on different cloud environments. Last, the various parts of our software architecture are highly generalizable. The workflow configurations are defined in a single file written in yaml that users can easily customize without expertise in cloud technologies for specific use cases. In addition, our framework does not modify users' Python modules (e.g., the module performing image reconstruction); instead, we plug the module into our framework, which automatically manages parallelization, job scheduling, and resource management.

Performance evaluation and scalability

Materials scientists seek workflow execution speed-up and efficient management of the massive amount of generated data. We assess whether our framework meets these goals for our materials science workflow. To this end, we show the scalability of our framework by evaluating the performance of the two most computation-intensive stages (i.e., reconstruction and segmentation) with an increasing number of GPU nodes. We present the increase in intermediate data size at the end of each stage to demonstrate the importance of storage resource elasticity.

We run our workflow on Chameleon Keahey et al. (2020), a public cloud platform supported by the National Science Foundation with sites at the University of Chicago and the Texas Advanced Computing Center. Chameleon is built on the OpenStack cloud software platform. It provides researchers access to nodes of traditional KVM cloud architectures and reconfigurable baremetal nodes provisioned via OpenStack Ironic. Table 1 gives the Chameleon's hardware specifications of the 16 GPU nodes and three CPU nodes that consist of the Dask and Kubernetes clusters used for our testing, as shown in Figure 4. Figure 8.a presents our scalability study regarding the workflow's total execution time (in hours). The figure illustrates the average time of three executions with 2, 4, 8, and 16 GPU nodes. We can observe how the execution time is significantly shortened by 88.7% running with 16 nodes

Table 2. The variance and standard deviation of three workflow executions on a different number of GPU nodes. Unit: hour.

Number of nodes	2 GPU nodes	4 GPU nodes	8 GPU nodes	16 GPU nodes
Variance	0.65	0.04	0.08	0.0
Stand deviation	0.81	0.21	0.29	0.06

Table 3. The variance and standard deviation of the time spent downloading the HDF5 files on each node of three execution and its ratio concerning total time. Unit: second.

Number of nodes	2 GPU nodes	4 GPU nodes	8 GPU nodes	16 GPU nodes
Variance	775.5	2.7	8.8	38.4
Stand deviation	27.8	1.7	3.0	6.2

Table 4. Increase dataset size, stage by stage, when processing 16 input images.

	Stage I	Stage 2	Stage 3	Stage 4	Stage 5	Stage 6
Accumulative size (GB)	96.7	939.9	1003.9	1022.3	1022.3	1022.3

compared to the two nodes. Table 2 gives the measurements' variation and standard deviation, indicating that the execution time is stable on Chameleon.

The inputs in Stage 1 are 16 HDF5 files, with each having size of 6 GB. Computation is delayed until the files are downloaded to the local disk at each GPU node. The impact of file uploading to S3 is insignificant since it is performed by background threads when computation occurs. Figure 8(b) shows the average download time spent on each GPU node. Our experiments have experienced a stable network as given by the variance and standard deviation in Table 3. However, the network speed can fluctuate depending on cloud providers. For instance, the download speed can become 10 times slower than the network's peak speed. Notice that running with two nodes, the variance, and standard deviation are much higher than the rest. The experiment is more likely to undergo network speed variation on the cloud when its execution expands across several days. Consequently, the overall execution time of the two stages is significantly influenced. Nevertheless, compared with the total execution time in our experiment, such cost is still negligible. Table 3 gives the ratio of time cost in file downloading concerning the total time, which is within 1%.

Table 4 gives the data generated stage by stage cumulatively. The volume of streamed images at the end stage reaches approximately 10 times the input size. The experiment in this section only evaluates processing time for a small number of tomography of porous silica-based materials with around 100 GB, yet amounts to the size at a magnitude of terabyte at the final stage.

Figure 9 visualizes a set of output images generated with our framework and visualized with our Jupyter notebook interface (the top layer of our framework). Specifically, the figure shows silica-based nanopillar images generated with synchrotron-based X-ray computed microtomography during a compression test. The novelty aspect of this set of images is the first-of-its-kind visualization of the spherical silica nanoparticles (white to gray color) and void space (black color) in this type of testing.

Aspects of novelty and future directions

Our work is part of a broader effort to port scientific applications to the cloud. We present here a suite of other contributions, outlining for each what our framework introduces the novelty aspects. We also define the future direction to support the reusability, portability, and reproducibility of results. Jain et al. (2015) present a dynamic workflow system used in materials projects to calculate elastic tensors describing the response of materials to external forces within their elastic limits. This workflow system packs and schedules many small tasks for HPC systems that are otherwise not suitable for HPC queuing systems. Lingerfelt et al. (2016) presented a workflow system BEAM targeted for analysis and modeling of material characterization data on HPC systems. A web portal is given to the user for submitting, monitoring the multi-stage workflow execution, and visualizing the results. The algorithms are adapted to utilize MPI to scale on HPC. Our work provides a lightweight workflow system that enables moving applications to the cloud. The parallel tasks and the big data storage requirements benefit from the cloud Iaas (Manyi and Krishna Shyam, 2014) (Infrastructure As A Service) model at an economical cost.



Figure 9. Selected output of our framework, generated with the Jupyter notebook interface. The figure shows a first-of-its-kind visualization of the spherical silica nanoparticles (white to gray color) and void space (black color) in silica-based nanopillar images generated with synchrotron-based X-ray computed microtomography during a compression test).

Taghizadeh-Popp et al. (2020) describe a data-driven platform SciServer developed for astronomy. One key feature of this platform is to perform data filtering, analysis, and visualization close to data in the SQL database. The concept of bringing analysis close to data is fundamental in the era of big data, which is also one of the issues addressed in our architecture, where the inputs are fetched from remote storage to where computation occurs. Differently, our workflow system does not require a central database. Instead, a cloud object storage is used for long-term backup, and the current desired datasets are downloaded to the Kubernetes cluster providing data-streaming and visualization services.

Abernathey et al. (2021) prototype an architecture named Pangeo that is composed of a cloud object storage and compute cluster. The external cloud storage is dedicated to the easy retrieval of large-volume datasets. Kubernetes provision the compute cluster, and Dask enables computation parallelism. Differently, our architecture provides an

additional deployment layer on top of Dask and Kubernetes, which provisions native environments specified by users for their workflows rather than pulling existing container images, which can be inflexible.

There is a growing need for developing persistent scientific workflows to seamlessly connect and integrate software stacks and data services across cloud platforms supported by virtualization and data provenance (Bhatia et al., 2021). Containerization of scientific workflow enables reusability, portability, and reproducibility of results (Olaya et al., 2023) and ease of system maintenance efforts (Dusia et al., 2015); McDaniel et al., 2015); Monsalve et al., 2015); Herbein et al., 2016). Future directions point to the need for automatic containerization of complete working environments that include software dependencies (e.g., Python programs/modules) at all stages and the components of the Dask cluster, which are currently running natively in our workflow to maximize performance. Containerized Dask components can be deployed to the

corresponding GPU nodes through Ansible, automatically generating and scheduling the containers to connected nodes.

Conclusion

The software framework presented herein solves multistage workflows in scientific domains. The framework enables automatic workflow scheduling, scalable data management at runtime, and interactive data visualization. Domain scientists can easily plug in their Python modules as the workflow stages and run the workflows with a Jupyter notebook. In addition, this framework also offers a deployment layer to set up the software stacks on cloud platforms and provisions dependencies for programs according to their needs. We described a complex materials science workflow in the context of in situ mechanical testing of a nanoporous structure and explored the performance benefits brought by our framework. We presented the technical innovations that can extend the materials science workflow into a generalizable and scalable framework. Future work includes integrating our workflow into data commons, such as the Materials Commons and Digital Rock, which will collect the intermediate data enabling automatic indexing and metadata generation for data searchability.

Acknowledgments

This research was supported by the National Science Foundation (NSF) under grant numbers #1841758, #2028923, #2103845 and #2138811; the Extreme Science and Engineering Discovery Environment (XSEDE) under allocation TG-CIS210128; Chameleon Cloud under allocation CHI-210923; and IBM through a Shared University Research Award. Author BPC was supported by the JHU/APL Independent Research and Development Program.

Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This work was supported by the US National Science Foundation (#1841758, #2028923, #2103845 and #2138811).

Data Availability

The framework software and the Jupyter Notebook are both open access at https://github.com/nsdf-fabric/nsdf-materialscience.

ORCID iD

Michela Taufer https://orcid.org/0000-0002-0031-6377

Note

- *. yaml: yet another markup language.
- †. https://www.prefect.io/.

References

- Abernathey RP, Augspurger T, Banihirwe A, et al. (2021) Cloudnative repositories for big scientific data. *Computing in Science and Engineering* 23(2): 26–35.
- Bhatia H, Di Natale F, Moon JY, et al. (2021) Generalizable coordination of large multiscale workflows In: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC), St. Louis, MO, USA. 14-19 November 2021, 1–16, IEEE. DOI: 10. 1145/3458817.3476210.
- Cepuc A, Botez R, Craciun O, et al. (2020) Implementation of a continuous integration and deployment pipeline for containerized applications in amazon web services using jenkins, ansible and kubernetes In: Proceedings of the 2020 19th RoEduNet Conference: Networking in Education and Research (RoEduNet), Bucharest, Romania. 11-12 December 2020, 1–6. IEEE
- Deelman E, Vahi K, Juve G, et al. (2015) Pegasus, a workflow management system for science automation. *Future Generation Computer Systems* 46: 17–35.
- Dusia A, Yang Y and Taufer M (2015) Network quality of service in docker containers. In: Proceedings of the 2015
 IEEE International Conference on Cluster Computing. Chicago, IL. 08-11 September 2015, 527–528. IEEE
- en-Kiki O, Evans C and Ingerson B (2009) Yaml ain't markup language (yaml) version 1.1. Working Draft 2008.
- Factor M, Meth K, Naor D, et al. (2005) Object storage: the future building block for storage systems In: Proceedings of the 2005 IEEE International Symposium on Mass Storage Systems and Technology, Sardinia, Italy. 20-24 June 2005, pp. 119–123. IEEE
- Group H (2000-2010). Hierarchical data format version 5 http://www.hdfgroup.org/HDF5.
- Herbein S, Dusia A, Landwehr A, et al. (2016) Resource management for running hpc applications in container clouds. In: Kunkel J, Balaji P and Dongarra J (eds), *High Performance Computing. ISC High Performance 2016. Lecture Notes in Computer Science*, vol. 9697. Springer, Cham. https://doi.org/10.1007/978-3-319-41321-1 14.
- Hightower K, Burns B and Beda J (2017) *Kubernetes: Up and Running Dive into the Future of Infrastructure*. Sebastopol, CA, USA: O/'Reilly Media, Inc.
- Jain A, Ong SP, Chen W, et al. (2015) FireWorks: a dynamic workflow system designed for high-throughput applications. Concurrency and Computation: Practice and Experience 27(17): 5037–5059
- Keahey K, Anderson J, Zhen Z, et al. (2020) Lessons learned from the chameleon testbed In: *Proceedings of the 2020 USENIX*

- Conference on Usenix Annual Technical Conference. USE-NIX Association.
- Lei Q, Guo J, Noureddine A, et al. (2020) Sol-gel-based advanced porous silica materials for biomedical applications. *Advanced Functional Materials* 30(41): 1909539.
- Lingerfelt EJ, Belianinov A, Endeve E, et al. (2016) BEAM: A computational workflow system for managing and modeling material characterization data in hpc environments. *Procedia Computer Science* 80: 2276–2280. https://doi.org/10.1016/j.procs.2016.05.410.
- Maire E, Buffière JY, Salvo L, et al. (2001) On the application of x-ray microtomography in the field of materials science. *Advanced Engineering Materials* 3(8): 539–546.
- Manvi SS and Krishna Shyam G (2014) Resource management for infrastructure as a service (IaaS) in cloud computing: a survey. *Journal of Network and Computer Applications* 41: 424–440.
- McDaniel S, Herbein S and Taufer M (2015) A two-tiered approach to i/o quality of service in docker containers In: Proceedings of the 2015 IEEE International Conference on Cluster Computing, Chicago, IL. 08-11 September 2015, 490–491. IEEE
- Merkel D (2014) Docker: lightweight linux containers for consistent development and deployment. *Linux J* 2014(239).
- Monsalve J, Landwehr A and Taufer M (2015) Dynamic CPU resource allocation in containerized cloud environments. In: 2015 IEEE International Conference on Cluster Computing. Chicago, IL. 08-11 September 2015, 535–536, IEEE.
- Olaya P, Kennedy D, Llamas R, et al. (2023) Building trust in earth science findings through data traceability and results explainability. *IEEE Transactions on Parallel and Distributed Systems* 34(2): 704–717, IEEE.
- Pascucci V, Scorzelli G, Summa B, et al. (2012) The ViSUS Visualization Framework. In: Bethel EW, Childs H and Hansen CD (eds), High Performance Visualization—Enabling Extreme-Scale Scientific Insight, Chapman and Hall/CRC computational science series. Florida: CRC Press.
- Penny D (2020) NASA: We Forgot the \$30m Bandwidth Charges. https://www.linkedin.com/pulse/oops-we-forgot-30m-bandwidth-charges-david-penny/.
- Petruzza S, Venkat A, Morrical N, et al. (2020) The OpenVisus Framework for Extreme Data Management, Analysis and Visualization. https://github.com/sci-visus/OpenVisus.
- Peukert W, Segets D, Pflug L, et al. (2015) Unified design strategies for particulate products *Unified Design Strategies* for Particulate Products, 46, pp. 1–81.
- Randles BM, Pasquetto IV, Golshan MS, et al. (2017) Using the jupyter notebook as a tool for open science: an empirical study In: Proceedings of the 2017 ACM/IEEE Joint Conference on Digital Libraries. Toronto, ON, 19-23 June 2017, 1–2. JCDL
- Rocklin M (2015) Dask: parallel computation with blocked algorithms and task scheduling In: Proceedings of the 14th Python in Science Conference (SCIPY), Los Angeles, C. 126–132.

- Sammons G (2016) Exploring Ansible 2: Fast and Easy Guide. CreateSpace Independent Publishing Platform. https://books.google.com/books?id=vgQiMQAACAAJ.
- Singh R, Graves JA, Anantharaj V, et al. (2019) Evaluating scientific workflow engines for data and compute intensive discoveries In: Proceedings of the 2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA. 09-12 December 2019. 4553–4560. IEEE
- Stodden V, McNutt M, Bailey DH, et al. (2016) Enhancing reproducibility for computational methods. *Science* 354(6317): 1240–1241.
- Taghizadeh-Popp M, Kim J, Lemson G, et al. (2020) SciServer: a science platform for astronomy and beyond. *Astronomy and Computing* 33: 100412.
- Taufer M, Deelman E, da Silva RF, et al. (2021a) A roadmap to robust science for high-throughput applications: the scientists' perspective In: Proceedings of the 2021 IEEE 17th International Conference on eScience (eScience), Innsbruck, Austria, 20-23 September 2021. 247–248.IEEE
- Taufer M, Deelman E, Silva RF, et al. (2021b) A roadmap to robust science for high-throughput applications: the developers' perspective In: Proceedings of the 2021 IEEE International Conference on Cluster Computing (CLUSTER), Portland, OR, 07-10 September 2021. 807–808.IEEE
- Wang J, Crawl D, Purawat S, et al. (2015) Big data provenance: challenges, state of the art and opportunities In: Proceedings of the 2015 IEEE International Conference on Big Data (Big Data), Santa Clara, CA, 29 October 2015—01 November 2015, 2509–2516.IEEE
- Wilde M, Hategan M, Wozniak JM, et al. (2011) Swift: a language for distributed parallel scripting. *Parallel Computing* 37(9): 633–652.
- Wu F, Wu Q and Tan Y (2015) Workflow Scheduling in Cloud: a Survey. *The Journal of Supercomputing* 71(9): 3373–3418.

Naweiluo Zhou was a Postdoctoral Researcher at the University of Tennessee Knoxville.

Giorgio Scorzelli is the Director of Software Development at the Center For Extreme. Data Management Analysis And Visualization and a Research Scientist at the Scientific Computing Institute of the University of Utah. He earned his Computer Science Master's degree at the University of Roma Tre (2000). He was an Adjunct Professor and Research. Fellow at the University of Roma Tree (2001–2010). He worked as a Computer Scientist for the Lawrence Livermore National Laboratory (2002–2005). Scorzelli's research interests include big data analysis, real-time scientific visualization, and multi-resolution and streaming techniques. He authorizes over 40 peer-reviewed publications in international conferences and journals.

Jakob Luettgau is a Research Assistant Professor at the University of Tennessee, Knoxville. He earned his Ph.D. in

Computer Science in 2021 from the University of Hamburg (Germany) after graduating with his BSc and MSc from the same University in 2014 and 2016. From 2016 to 2019, he worked with The Centre of Excellence in Simulation of Weather and Climate in Europe (ESiWACE) at the German Climate Computing Center (DKRZ). He visited Argonne National Laboratory (ANL) as a Short-Term Scholar in 2018 and 2019 while holding a research position at DKRZ. In 2020, he started work as an AI Consultant for Helmholtz AI/DKRZ, advising researchers in the earth sciences throughout Germany. His research interests include high-performance storage, workflows leveraging hierarchical storage systems, machine learning and performance engineering, and instrumentation of parallel applications.

Rahul Reddy Kancharla is a Postdoctoral Research Associate at the Idaho National Laboratory, Idaho Falls, Idaho. He obtained his Bachelor of Technology (2015) in Computer Engineering from the Jawaharlal Nehru Technological University. He earned his Masters from the University of Wyoming College of Engineering and Physical Sciences (2017) and his Ph.D. from Boise State University (2021)

Joshua J. Kane is a staff scientist at Idaho National Laboratory with research interests in a wide range of carbon materials. His background is in Chemical Engineering and Materials Science. He obtained his Bachelor's (2009) in Chemical Engineering from the University of Idaho. He earned his Masters and Ph.D. (2013) from Boise State University in Materials Science and Engineering.

Brendan P. Croom is a senior materials scientist at JHU/ APL with diverse expertise in metal additive manufacturing, process monitoring, materials qualification, X-ray computed tomography and non-destructive evaluation. Previously, as an NRC postdoctoral researcher at AFRL and a graduate student at the University of Virginia, he developed novel methods to additively manufacture and evaluate the performance of structural composites for extreme environments. He is especially interested in machine learning methods to accelerate materials development, NDE, and build processing-structure-property-performance relationships for emerging materials systems. Robert Wheeler is the Chief Executive Manager of MicroTesting Solutions LLC, where he also directs the development of small-scale testing methods and instrumentation. At MTSL, he develops technologies for conducting in situ mechanical, thermal, and electrical studies within the SEM, FIB, and TEM. He obtained his Ph.D. from the University of Illinois at Urbana-Champaign. He held a postdoctoral position as a site electron microscopist for the NSF Center for High-Temperature Superconductivity at Argonne National Laboratory. He had a research scientist position at The Ohio State University before joining UES, Inc., where he concurrently acts as an on-sight scientist for the Air Force Research Laboratory.

Pania Newell is an Assistant Professor in the Department of Mechanical Engineering and holds an adjunct position at the School of Computing and the Department of Civil Engineering at the University of Utah. She obtained her M.S. and Ph.D. from the University of New Mexico and University of Colorado-Boulder. Her research interest lies at the interface of material sciences and mechanics. In particular, she is interested in multi-scale, multi-physics phenomena in heterogeneous porous materials through developing theoretical, computational, and experimental frameworks combined with data sciences.

Valerio Pascucci is the founding Director of the Center for Extreme Data Management Analysis and Visualization (CEDMAV) at the University of Utah. Pascucci is a Faculty of the Scientific Computing and Imaging Institute, a Professor of the School of Computing, University of Utah, and a DOE Laboratory Fellow, of the Pacific Northwest National Laboratory. Previously, Pascucci was a Group Leader and Project Leader in the Center for Applied Scientific Computing at the Lawrence Livermore National Laboratory and an Adjunct Professor of Computer Science at the University of California Davis. Pascucci earned a Ph.D. in computer science at Purdue University (2000) and a EE Laurea (Master) at the University "La Sapienza" in Roma, Italy (1993).

Michela Taufer is the Dongarra Professor in High-Performance Computing in the Department of Electrical Engineering and Computer Science at the University of Tennessee Knoxville (UTK). She earned her undergraduate degree in Computer Engineering from the University of Padova (Italy) and her doctoral degree in Computer Science from the Swiss Federal Institute of Technology or ETH (Switzerland). From 2003 to 2004, she was a La Jolla Interfaces in Science Training Program (LJIS) Postdoctoral Fellow at the University of California San Diego (UCSD) and The Scripps Research Institute (TSRI), where she worked on interdisciplinary projects in computer systems and computational chemistry. Taufer's research interests in high-performance computing include scientific applications, scheduling, reproducibility challenges, and big data analytics. She has nearly 100 publications and delivered 80 talks at various conferences and research institutes. She is a professional member of the IEEE and a Distinguished Scientist of the ACM.