

Flux: Overcoming Scheduling Challenges for Exascale Workflows

Dong H. Ahn*, Ned Bass*, Albert Chu*, Jim Garlick*, Mark Grondona*, Stephen Herbein*, Joseph Koning*, Tapasya Patki*, Thomas R. W. Scogland*, Becky Springmeyer* Michela Taufer†

*Lawrence Livermore National Laboratory, 7000 East Ave. Livermore, CA

{ahn1, bass6, chu11, garlick1, grondona1, herbein1, koning1, patki1, scogland1, springmeyer1}@llnl.gov

†University of Tennessee, Knoxville. Knoxville, TN
taufer@utk.edu

Abstract—Many emerging scientific workflows that target high-end HPC systems require complex interplay with the resource and job management software (RJMS). However, portable, efficient and easy-to-use scheduling and execution of these workflows is still an unsolved problem. We present Flux, a novel, hierarchical RJMS infrastructure that addresses the key scheduling challenges of modern workflows in a scalable, easy-to-use, and portable manner. At the heart of Flux lies its ability to be nested seamlessly within batch allocations created by other schedulers as well as itself. Once a hierarchy of Flux instance is created within each allocation, its consistent and rich set of well-defined APIs portably and efficiently support those workflows that can often feature non-traditional execution patterns such as requirements for complex co-scheduling, massive ensembles of small jobs and coordination among jobs in an ensemble.

I. INTRODUCTION

Scientific workflows continue to become more complex, and their execution patterns are also drastically changing. In order to exploit the ever-growing compute power of systems and upcoming exascale platforms, modern workflows increasingly employ multiple types of simulation applications coupled with in-situ visualization, data analytics, data stores and machine learning [1], [2], [3], [4]. The current push towards rigorous verification and validation (V&V) and uncertainty quantification (UQ) [5] approaches often features simulations that involve enormously large numbers of short-running jobs (e.g., reduced models and 1-D simulations), straying away from traditional long-running execution.

These trends have become ever more apparent on some of the most massive high performance computing (HPC) systems, such as the Sierra [6] and Summit [7] machines, which are the new pre-exascale systems being fielded by the world’s largest supercomputing centers. Three major early science applications running on Lawrence Livermore National Laboratory (LLNL)’s Sierra, including [2], for instance, now embrace non-traditional workflows. Additionally, our recent analysis on other large *production* clusters at LLNL shows that 48.1% of jobs involved the submission of at least 100 identical jobs by the same user with 27.8% submitted within one minute of each other, a pattern typically associated with V&V and UQ. Such workflows, often referred to as *ensemble-based*, are quickly becoming a norm.

Resource and job management software (RJMS) is central to enabling efficient execution of applications on HPC systems, and therefore is also the main interface for scheduling and executing these complex workflows. However, recent trends of complex workflows with new execution patterns, significantly complicate efficient (co-)scheduling and execution of their tasks. In particular, traditional *centralized* techniques implemented within RJMS such as SLURM [8], LSF [9], MOAB [10], or PBS Pro [11] no longer work well as they are fundamentally designed for the traditional paradigm: a few large, long-running, homogeneous jobs rather than ensembles composed of many, and often small, short-running heterogeneous elements.

These limitations are already presenting greater technical challenges for exascale workflows, which will only worsen if not met properly. Four such key challenges are listed below.

- 1) *Throughput Challenge*: Large ensemble simulations require massive numbers of jobs that cannot comfortably be ingested and scheduled by the traditional approach;
- 2) *Co-scheduling Challenge*: Complex coupling requires sophisticated co-scheduling that the existing centralized approaches cannot easily provide;
- 3) *Job coordination and communication challenge*: Intimate interactions with RJMS is required to keep track of the overall progress of the ensemble execution, and existing approaches lack well-defined interfaces;
- 4) *Portability Challenge*: There has been a proliferation of ad hoc implementations of user-level schedulers as an attempt to tackle the above challenges. They are often non-portable and come with a myriad of side effects (e.g., millions of small files just to coordinate the current state of an ensemble).

In this paper, we present Flux, a novel resource management and scheduling infrastructure that overcomes these challenges in a scalable, easy-to-use, portable, and cost-effective manner. At the core of Flux lies its ability to be seamlessly nested within allocations created by other resource managers or itself, along with allowing for user-level customization of policies and parameters. This *fully hierarchical* approach allows the target workflows to submit fewer jobs that resemble the traditional execution pattern to the low-level schedulers, most notably the native system scheduler, while

more fine-grained scheduling is performed by a hierarchy of nested instances running within each allocation. Each level also allows customizable scheduling policies and parameters, addressing both the *throughput and co-scheduling challenges*.

In addition, Flux is designed from the ground up as a software framework with a rich set of well-defined APIs such as: job submission, job status and control, messaging, as well as input and output streaming. Workflows can use any of these to *facilitate communication and coordination of various tasks* to be executed within and across ensembles. Finally, to address *portability challenges*, its APIs are consistent across different platforms. Creating an instance requires only the lower-level resource manager to provide the Process Management Interface (PMI), the de facto standard for MPI bootstrapping, or the user to provide a configuration.

Specifically, this paper makes the following contributions:

- Identification and discussions of specific exascale workflow scheduling challenges based on emerging practices at LLNL, one of the world’s largest supercomputer centers;
- Novel hierarchical approaches for providing resource management and scheduling infrastructure at the user level to address the above challenges;
- Performance evaluations of our hierarchical approaches on up to one million short-running jobs using both synthetic and real simulation codes;
- Case studies and lessons learned from integrating our approaches to three distinct real-world workflow management systems targeting exascale computing;
- Discussions on techniques needed to address the remaining challenges.

Our evaluation with three recent workflow efforts at LLNL shows that Flux significantly overcomes all of the stated challenges. Our performance measurements on synthetic and real ensemble-based workflows suggest that our hierarchical scheduling approach can improve the job throughput performance of these workflows by $48\times$. Further, our case study on the Cancer Moonshot Pilot2 project shows that Flux can efficiently co-schedule a new workflow that employs machine learning to couple a large continuum model-based application with an ensemble of thousands of MD simulations starting and stopping during a run at high speed. Finally, our integration with Merlin, a workflow management system designed to support next-generation machine learning on HPC, shows that Flux significantly enables not only co-scheduling of various task types within each ensemble but also its needs for high portability and task communication and coordination.

II. COMPLEX EXECUTION PATTERNS EXEMPLIFIED BY CANCER MOONSHOT PILOT2

To motivate the need for our technology, we consider the Cancer Moonshot Pilot2 workflow as our motivational example, an early science application being run on LLNL’s Sierra system. This workflow features non-traditional co-scheduling and execution patterns that the existing system scheduler could not reasonably provide.

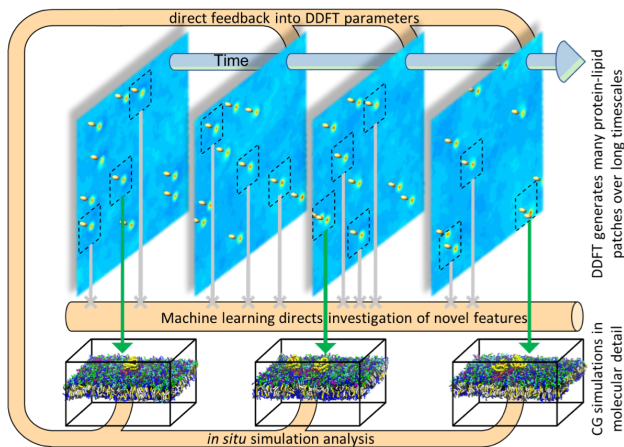
The Cancer Moonshot project aims at furthering cancer research and advanced drug discovery through HPC simulations. It was discovered over forty years ago that 30% of human cancers are caused by the RAS family of cancer-causing genes [12]. Yet, there are still no drugs targeting RAS because both computational techniques cannot explore molecular interactions at high resolutions with the right sizes and time scales. The Pilot2 project seeks to develop an effective HPC simulation method to uncover the detailed characterizations of the behavior of RAS in cellular membranes.

The Pilot2 project combines continuum model-based and MD simulations to bring the best from both worlds. Here, the novel continuum model coupled with a machine-learning module drives the sampling of *patches*—small neighborhoods around a molecule of RAS. These patches are then used to instantiate and run corresponding MD simulations. Additionally, several *in situ* processing capabilities need to be connected in order to control how long a particular MD simulation is run and to provide feedback to the continuum model for parameter refinement. This process is depicted in Figure 1a.

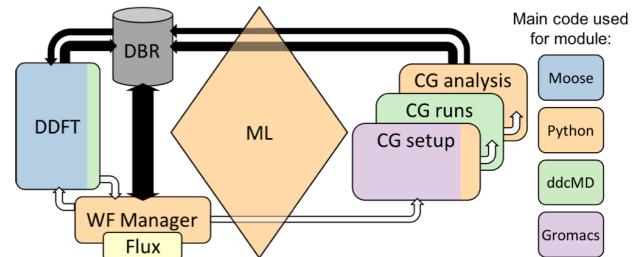
Figure 1b presents the framework in detail. The current workflow is coordinated through the IBM DataBroker (DBR), which provides cross-machine, shared access to storage for data and message exchange. At the macroscale level, to simulate a membrane at biologically relevant and experimentally accessible time and length scales, the continuum model is used with the finite element solver MOOSE [13]. This DDFT simulation is then coupled to a Langevin particle model running on ddcMD [14] that allows the evolution of discrete RAS proteins on the membrane.

At each time step of the continuum model, 300 patches are extracted (one centered on each RAS protein) and compared to all patches previously explored using MD simulations. Whenever computing resources become available, the most unusual new patch (i.e., the patch with the largest distance to its neighbors in latent space) is taken and a new corresponding MD simulation is started. The framework discussed above crucially relies on the ability to automatically instantiate MD simulations, monitor them on-the-fly, and provide feedback to the continuum model. RAS orientations are selected from pre-constructed libraries and pulled to the membrane surface.

The CG setup module is Python-based and uses the GROMACS MD package v5.1.2 [15] before ddcMD evolves the given patch on a GPU. A new ddcMD version that implements the Martini force field with a new strategy for domain decomposition, and an atom padding technique has been developed in CUDA to leverage the benefits of GPUs. This implementation offloads the entire computation to the GPU. Every non-constant time calculation step necessary for Martini now runs on the GPU via CUDA kernels, including the integrator and constraint solver, such that particles are only communicated back to the host for I/O purposes and never for calculations of the particle forces or movement. This leaves the CPUs tasked with only managing the order and launches of the aforementioned kernels. While the MD simulations are running, analysis modules are executed every two seconds of



(a) By implementing an adaptive multiscale model, the Cancer Moonshot Pilot2 project directly couples molecular detail to a cellular-scale continuum simulation. Machine learning directs instigation and investigation of coarse-grained (CG) particle simulations from only the continuum (DDFT) simulation patches with novel features, allowing for intelligently sampling of the simulation space far more efficiently, resulting in a scope of exploration that is not achievable using only brute force calculations. Furthermore, in situ analysis of the CG simulations and feedback allows for the DDFT simulation parameters to evolve in real time, incorporating the vast sampling carried out at the particle level.



(b) Multiscale code framework. The WorkFlow (WF) Manager connects two scales: DDFT and CG. Frames resulting from the DDFT simulation are decomposed into patches, and the WF Manager feeds them to the machine learning (ML) infrastructure, which maintains a priority queue of candidate patches. When new resources become available, the WF Manager picks top candidates and uses the Flux resource manager to start new CG simulations. Data transfer and messaging are handled through the DataBroker (DBR), which implements a fast, system-wide key-value store. Thickness of black arrows represents the bandwidth of data flow to and from the DBR.

Fig. 1: Multiscale Simulation for RAS Initiation of Cancer

wall time to accumulate data of interest continuously. When running on 3,500 nodes of Sierra, the workload needed to run a single 1,000 node continuum model, a single node machine learning and workflow management system, the data broker, and GROMACS simulations on CPUs of all 3,500 nodes. While those were running, four separate ddcMD simulations were run on each node using the GPUs, running at least 5 logically separate items on each node. In order for this to work well, the job execution system needed to manage at least 7,500 simultaneous jobs and continually re-schedule work as microscale jobs completed.

Overall, this workflow exemplifies the many (co-)scheduling and execution challenges faced by emerging workflows. They include co-scheduling of coupled simulations at different scales (i.e., continuum models-based simulations with several thousand MD simulations, coordination between CPU and GPU runs), the use of a machine learning module to schedule (or de-schedule) and execute simulations dynamically at a high rate, and the use of data store to coordinate the data flow between different tasks. We will further characterize the key scheduling and execution challenges such as the ones shown in the Cancer Moonshot Pilot2 workflow in the next section.

III. CHALLENGES IN WORKFLOW SCHEDULING

This section characterizes the workflow scheduling and execution challenges based on our analysis on some of the emerging workflow management practices at LLNL. Our analysis is based on our direct interactions with three distinct workflow management software development teams at LLNL,

namely the Cancer Moonshot Pilot2 workflow, Uncertainty Quantification Pipeline (UQP) [16], and the Merlin workflow that supports extreme-scale machine learning, as well as interviews with developers of other workflow management software such as PSUADE UQ framework [17] and end users who have created *ad hoc* schedulers for their workflows. While each of these workflows often addresses entirely different domains of science, they exhibit common scheduling issues. As briefly highlighted in Section I, they are referred to as throughput, co-scheduling, job coordination/communication, and portability challenges.

A. Throughput Challenge

Many workflows feature large ensembles of small, short-running jobs, which can create thousands or even millions of jobs that need to be rapidly ingested and scheduled. For the Cancer Moonshot Pilot2 example presented in the previous section, several thousand MD simulations need to be run successfully with a quick turnaround time to facilitate the refinement of parameters in the continuum model and produce microscale results. In the case of the UQP, building a surrogate model can require tens to hundreds of thousands of simulation executions to adequately sample the simulation's input parameter space. Such ensemble workloads are becoming a norm rather than an exception on high-end HPC systems.

Traditional RJMS in most cloud and HPC centers today are based on centralized designs. Cloud schedulers such as Swarm and Kubernetes [18], [19] and HPC schedulers such

as SLURM, MOAB, PBSPro and LSF [8], [10], [9], [11] are implemented using this model. This model often fails to cope with rapid job ingestion, and because of this, a site imposes a cap on the number of jobs submitted at once and allowed in the scheduler. The cap then requires workflow managers to throttle the rate of their job submissions to match the ingestion rate, artificially decreasing the job throughput of the workload.

Furthermore, this pattern can also lead to shared resource thrashing and exhaustion. For example, the Sequoia supercomputer at LLNL, which has 1.6 million cores, encountered several scale-up problems when users tried to run about 1500 small UQ jobs (1-4 MPI tasks each) at the same point in time in 2014. While SLURM and IBM's control software managed to expand their limits to about 3-5K simultaneously executing jobs after fine-tuning various configuration parameters for some cases, several rare errors still kept cropping up. Eventually, LLNL created a temporary solution by building CRAM, a library that packs many small jobs into a single large job [20]. Unfortunately, libraries such as CRAM are not the panacea for centralized schedulers, and even well-engineered centralized solutions can suffer from several scalability and resiliency issues.

B. Co-scheduling Challenge

Coupling in complex workflows requires co-scheduling of different components. In the example we presented earlier, the CPU and GPU workloads need to be co-scheduled effectively. Additionally, data needs to be communicated to the host when necessary, and support for *in situ* analysis, as well as online techniques, require other jobs to be active on the node. More specifically, the Pilot2 workflow needs to schedule four different kinds of jobs on CPUs only, and an additional type of job on some CPUs and GPUs. One of the four needs to be on every single node, along with an instance of the fifth job with GPUs also on every node. Moreover, this decision is dynamically determined by a machine learning module, a completely new execution pattern.

Most traditional schedulers do not allow for such customization, making it challenging to utilize resources well. Co-scheduling can offer several utilization and job throughput benefits, as well as allow for customization of application kernels and efficient co-existence of several workflow components. Current schedulers offer little or no support for sharing multiple kinds of jobs within an allocation or for customizing resource allocations such as cores or GPUs (or others, such as burst buffers). If at all, only fixed mechanisms for requesting allocations exist, and users cannot tune these from one application to the next or leverage their domain knowledge about the resource utilization of their application.

C. Job coordination and communication Challenge

Modern scientific workflows depend on data transfer between various components of a framework. For example, as we showed in Figure 1b, the information about unusual patches triggers additional MD simulations, which in turn are used for further parameter refinement. Multiple such simulations

need to be analyzed to understand an unusual scenario, which requires regular coordination and communication between jobs as well as within the job.

Existing schedulers have limited support for ingesting, storing/retrieving job output or job status information, often requiring inefficient communication through the file system. Many workflow managers, such as UQP circumvent these issues by having jobs create an empty file whenever they start or complete. This allows UQP to track the state of every job in the workflow, but it is at the cost of creating a large and unnecessary metadata load on the target file system, infringing on the performance of both the workflow itself and the entire system.

D. Portability Challenge

One of the common problems with emerging workflow management systems is that they have to be ported to a wide range of RJMS. With no common infrastructure for supporting their scheduling, the task of porting m workflows to n environments amounts to an $m \times n$ effort. Often, those point solutions are non-portable, and even if a solution is ported on a new platform, they can often come with a multitude of side-effects (e.g., creating too many files for ensemble status checking). The more complex the target workflow is, the more difficult porting would become because a new scheduler may not provide all of the advanced features that the workflow might have used in its previously tested schedulers.

Scientists and developers often need to rewrite their scripts from scratch in order to adapt to a new environment, potentially introducing several scripting/setup bugs, requiring additional testing, underutilizing resource allocations and reducing overall productivity. For example, in the RAS multiscale simulation, moving from a cluster that uses IBM's LSF and `jsrun` to another that relies on SLURM can be challenging regarding setup cost. Also, being able to leverage different heterogeneous resources, including GPUs and burst buffers, often requires new flags and configuration parameters to be specified. This often results in ad hoc solutions for application scheduling.

IV. FLUX

The Flux framework is a suite of projects, tools and libraries that can be used to provide site-customable resource managers and schedulers for large HPC centers. Flux supports a fully hierarchical architecture that allows for seamless nesting in a highly scalable, customizable, and resilient manner. The main foundation of the Flux framework is an overlay network underpinned by a communications broker that supports various messaging idioms (such as publish-subscribe and remote procedure calls) and asynchronous event handling, referred to as *flux-core*. The job scheduling component, *flux-sched*, consists of an engine that handles all the functionality common to scheduling. The engine has the ability to load one or more scheduling plugins that provide specific scheduling behavior. These can be user-defined or administrative, providing for a

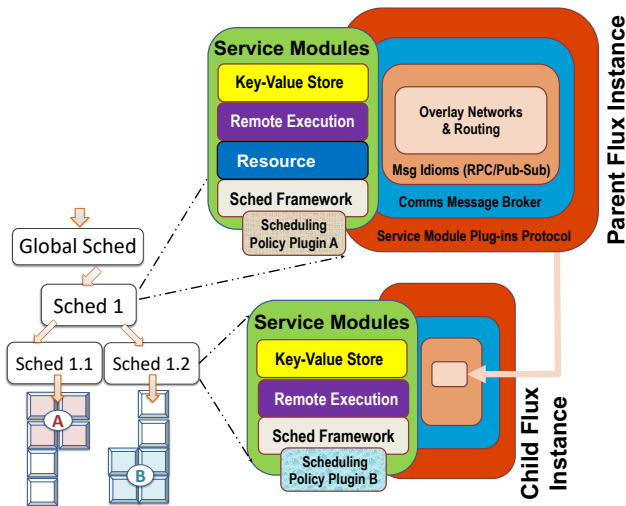


Fig. 2: Flux framework

truly customizable infrastructure. Figure 2¹ shows the modular architecture of Flux, and also depicts how the Flux network can be structured to manage two schedulers at different levels of the hierarchy, with a parent Flux instance and a child Flux instance. We discuss Flux’s fully hierarchical scheduling model in detail in the subsections below.

A. Scheduler Parallelism for Throughput Challenge

The hierarchical design of Flux provides ample parallelism to overcome the job throughput challenge present in traditional scheduling techniques. Under the hierarchical design of Flux, any Flux scheduler instance can spawn child instances to aid in scheduling, launching, and managing jobs. The parent Flux instance grants a subset of its jobs and resources to each child. This parent-child relationship, depicted in Figure 2, can extend to an arbitrary depth and width, creating a limitless opportunity for parallelization while avoiding the high communication overhead of other distributed schedulers (e.g., fully connected graphs of schedulers and all-to-all communications). Parent and child instances communicate using the Flux communication overlay network described further in Section IV-C.

Our current implementation of hierarchical Flux consists of three main design points: the scheduler hierarchy, the resource assignment, and the job distribution. For the scheduler hierarchy, our implementation supports a hierarchy of schedulers with a fixed size and shape. Ensemble workflow managers or users specify the exact hierarchy size and shape using JSON, which our implementation parses and uses to launch the corresponding scheduler hierarchy automatically. For the resource assignment, by default, our implementation assigns a uniform number of resources to schedulers at each level in the hierarchy (e.g., all of the leaf schedulers are allocated the same number of cores). Non-uniform

¹flux-core implements most of the core services and infrastructure depicted in this figure, while flux-sched provides scheduler-related services such as resource and sched framework services.

assignments of resources are possible but require careful consideration when distributing jobs.

To minimize the changes required for workflows to leverage hierarchical Flux, the workflow manager submits each job in the ensemble individually at runtime to the root scheduler instance (as it would with a traditional scheduler), and then, the jobs are distributed automatically across the hierarchy. In this configuration, it may seem that the root instance will become a bottleneck, but the work required to map and send a job to a child scheduler is significantly less than the work required to schedule and launch a job. After a job is submitted, the root instance in the hierarchy must only consider tens to hundreds of children, while a traditional scheduler must consider thousands of cores as well as all other jobs in the queue. Additionally, the job distribution at the root instance can overlap with the scheduling and launching of jobs at the leaf instances. For the job distribution, our implementation, by default, uses round-robin to distribute jobs uniformly across the scheduler hierarchy, but other distribution policies are supported and can be implemented by users.

B. Scheduler Specialization Solves Co-scheduling Challenge

Flux’s user-driven, customizable approach to scheduling provides inherent support for co-scheduling. Flux’s flexible design allows users to decide whether or not co-scheduling should be configured and also lets users choose their own scheduling policies within the scope of an instance. With the help of the job submission API, several tasks can efficiently coexist on a single node without any restrictions on their number, type, or resource requirements. This allows for submission and tuning at all possible levels of heterogeneity within a node (and across nodes), including individual cores, a set of cores, sockets, GPUs, or burst buffers.

Users can also choose a policy within their Flux instance. These can be simple policies, such as first-come-first-serve/backfilling, and the infrastructure can be easily extended to incorporate complex policies for advanced management of resources such as IO or power or multiple constraints. Traditional resource managers do not provide any such capability or extensible design to users, resulting in underutilized resources and limited throughput. While some workflows need exclusive scheduling per node, other workflows may need co-scheduling or different distributions of jobs between the resources available on a node. Traditional RJMS software has no support for user-level scheduling, which Flux addresses by design, giving users the freedom to adapt to their instance to the needs and characteristics of their particular application.

C. Rich APIs for Easy Job Coordination and Communication

Flux provides various communication idioms and APIs to help solve the job coordination and communication challenge. To support coordination within and across both Flux instances and jobs, Flux provides primitives that encapsulate the *publish-subscribe* (pub/sub), *request-reply*, and *push-pull* communication patterns. These primitives allow individual jobs within a workflow to synchronize without the use of ad

hoc methods like empty file creation on a POSIX-compliant file system. Flux also provides several high-level services that jobs and workflows can leverage: an in-memory key-value store (KVS) and a job status/control (JSC) API.

The KVS provided by Flux enables jobs and workflow managers to scalably retrieve and store information. One example KVS use-case for workflows is accessing job provenance data. All of a job’s metadata is stored in Flux’s KVS, including the resources requested, the environment variables used, and the contents of `stdout` and `stderr`. The storage of the `stdout` and `stderr` enables workflow managers to easily inspect a job’s output without requiring expensive file system accesses. A specific feature of Flux’s KVS, watcher callbacks, enables workflow managers to ingest and analyze a job’s output efficiently as it is being generated. Advanced workflows can leverage this real-time output analysis to detect job failures as they happen and take corrective actions, such as re-submitting the job for execution.

Traditional schedulers provide limited access to job status information, most commonly through a slow and cumbersome command line interface (CLI). Many workflow managers work around this interface by tracking job states via extraneous file creation. Flux’s JSC provides a fast, programmatic way to receive job status updates, eliminating the use of the slow CLI and tracking via the file system. JSC users can subscribe to real-time job status updates, which are sent whenever a job changes its state (e.g., from *running* to *completed*). This allows workflow managers to stay up-to-date on the state of their jobs with minimal overhead and without degrading file system performance.

D. Consistent API Set for High Portability

To serve as the common, portable scheduling infrastructure, Flux offers two main characteristics: 1) its APIs are consistent across different platforms and 2) the porting and optimization effort of Flux itself for a new environment is small. Creating a Flux instance on a given environment only requires the lower-level resource manager to provide the Process Management Interface (PMI), or the user to provide a configuration. Because PMI is the de facto standard for MPI bootstrapping, the system resource managers (including Flux itself) on a majority of HPC systems directly offer this interface or else provide other variant interfaces such as PMIx on top of which PMI can be easily implemented.

V. EVALUATING PERFORMANCE AND SCALABILITY

To demonstrate how Flux, with its fully hierarchical design, addresses the throughput challenge, we measure the scheduler throughput on real-world and stress-test ensemble workflows. We measure throughput as the average number of jobs ingested, scheduled, and launched per second (the higher, the better). We schedule the workflows using three different hierarchies: depth-1, depth-2, and depth-3². The depth-1 hierarchy only has a single scheduler instance that schedules

²Our model supports additional levels. In our evaluation, we use a one-to-one mapping between hardware and scheduler levels.

every job in the workflow, similar to existing schedulers like SLURM and Moab. For the depth-2 hierarchy, we create a root scheduler with one child scheduler for every node allocated to the workflow, and we distribute the jobs equally among the lowest level of schedulers (i.e., the leaf schedulers). For the depth-3 hierarchy, we extend the hierarchy by adding one scheduler for every core allocated to the workflow, and as with the previous hierarchy, we distribute the workflow’s jobs equally among the leaf schedulers. Our throughput evaluations on both workflows use 32 nodes of an Intel Xeon E5-2695v4 cluster, each node with 36 physical cores and 128 GB of memory.

To demonstrate the effects of hierarchical Flux on a real-world workflow, we generated an ensemble workflow with the Uncertainty Quantification Pipeline (UQP) [16]. Our UQ ensemble simulates a semi-analytical inertial confinement fusion (ICF) stagnation model that predicts the results of full ICF simulations [21], [22], [23]. UQ ensembles with this semi-analytical model typically consist of tens of thousands of runs, but the scientists’ goal is to execute millions of jobs.

Figure 3a shows the scheduler throughput of the three hierarchies when applied to variably sized real-world UQ ensemble workflows. For each ensemble size, we perform the test three times and present the min, max, and median job throughput values. As we increase the ensemble size, the throughput of the depth-1 scheduler plateaus at 10 jobs/sec, artificially limiting the overall performance of the ensemble workflow and creating idle resources. By adding additional levels to the scheduler hierarchy (i.e., depth-2 and depth-3) and thus increasing the scheduler parallelism, we can improve the peak job throughput by an order of magnitude. With a job throughput of 100 jobs/sec, the scheduler is no longer on the critical path of the workflow and the compute resources are 100% utilized. After the scheduler throughput enhancements provided by hierarchical scheduling, the ensemble workflow’s critical path now consists primarily of the ensemble application’s runtime.

To demonstrate the throughput capabilities of hierarchical Flux, unrestrained by the workflow application’s runtime, we created a *stress-test ensemble workflow* in which each job exits immediately after it launches (i.e., has a negligible runtime). Figure 3b shows the throughput of Flux on this stress-test workflow. As before, for each ensemble size, we perform the test three times and present the min, max, and median job throughput values. No longer limited by the workflow application’s runtime, the depth-2 and depth-3 hierarchies achieve a peak throughput of 370 jobs/sec and 760 jobs/sec, respectively. These represent a 23.5 \times and 48 \times increase over the job throughput achieved by the traditional, depth-1 scheduler.

VI. ENABLING EMERGING WORKFLOW MANAGEMENT WITH FLUX

In this section, we describe how we improve the scheduling and execution of real-world production workflows using Flux. Our study targets both the Cancer Moonshot Pilot2 workflow already described in Section II and the Merlin workflow.

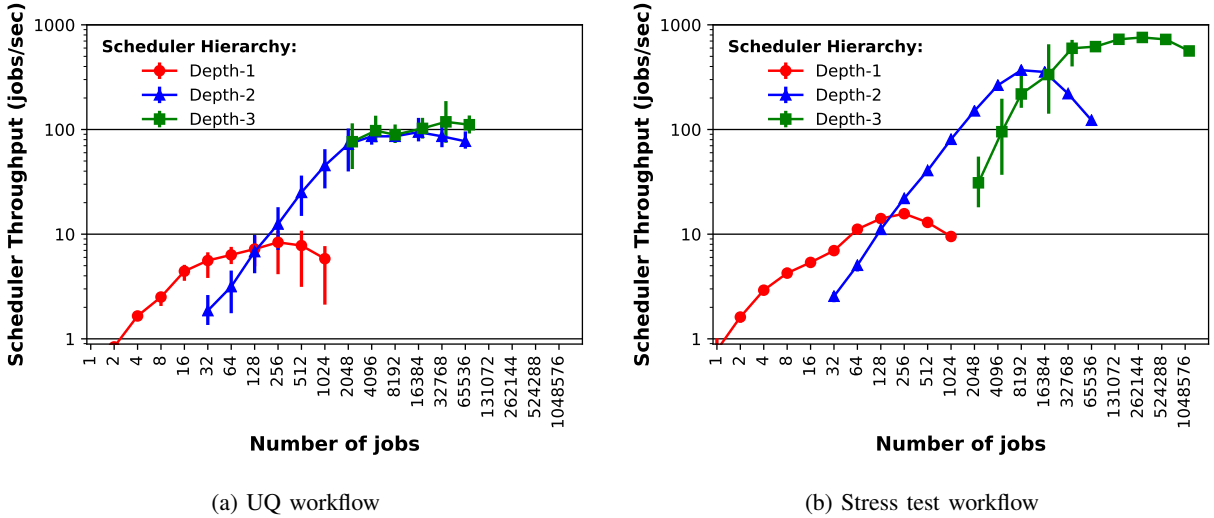


Fig. 3: Job throughput (in jobs/sec, on a logarithmic scale) for the depth-1, depth-2, and depth-3 scheduler hierarchies for fixed-size clusters and differing numbers of total jobs (on a logarithmic scale)

A. Easy, Scalable Interaction with the Scheduler for Merlin

The Merlin workflow is a component of the Machine Learning Strategic Initiative (MLSI) [2] at LLNL. Merlin’s goal is to provide a python based workflow that is adaptable and efficient. This workflow runs an ensemble of simulations and records the results while concurrently running machine learning on the results as they become available. The machine-learned model then helps steer the ensemble of simulations as it improves with more data.

The workflow executes a variety of tasks to generate and analyze the data. The first of these is defining the ensemble of simulations. This ensemble consists of a set of samples spanning the domain needed for creating a unique set of data describing the domain. A simulation executable task will accept the sample set as input parameters and produce data for the machine-learning model. The simulation can range from a simple ODE to a massively parallel MPI rad-hydro simulation. These simulations may also be run on a heterogeneous set of compute resources where scheduling and launching the simulations in a general manner becomes difficult.

The first version of the Merlin MPI parallel launch used a simple python-based subprocess call to take a set of MPI parameters such as the number of nodes and tasks and map them onto the commands needed for a SLURM or LSF launch. This became a maintenance issue when each new batch system required a set of runtime parameters that do not map one to one between the various launch systems. In the case of `jsrun` for LSF, the system did not handle nested launches where there was one `jsrun` call for the allocation and a subsequent `jsrun` call for the simulation. Some parallel runs need GPU support and few CPU cores, while others require only CPU cores. This requirement puts the onus on the workflow to schedule resources for the various types of parallel jobs.

In Merlin, Flux scheduler solves both the nesting issue and co-scheduling issue through the use of a single instance. Jobs can be co-scheduled because the single flux instance is tracking all of the resources. Nesting is not an issue due to

this single instance.

In the Flux-based launch system for Merlin, the python subprocess call was replaced with a Flux `rpc_send` with a `job.submit` command that includes the environment and resource request for this job. This Flux instance can be augmented with a callback function that will be invoked on each status change of the submitted job so the workflow can be informed on all stages of the job submission: `submitted`, `completed`, `anceled`, and `failed`. This information can be sent back through the Merlin workflow to inform the system on the state of the simulation task. This Flux interface is independent of the native job launcher and provides a single interface for the user to configure a simulation launch.

In addition, we also evaluated the two main characteristics described in Section IV-D by integrating Flux into Merlin across two environments: LLNL’s large clusters with Intel Xeon E5-2695v4, RedHat Enterprise Linux (RHEL) 7, and SLURM being resource manager and scheduler; and LLNL’s Sierra pre-exascale system with a completely different environment: IBM POWER little endian, RHEL7, IBM JSM being the resource manager and LSF the system scheduler.

We first designed and implemented our initial integration, on one of LLNL’s Intel Xeon E5-2695v4 Linux clusters. Then, we ported and customized Flux on Sierra while Merlin is being tested on the Intel Xeon Linux systems. Flux’s porting efforts on Sierra are mainly threefold.

- Port a PMI library to PMIX because the PMI library, though a de facto standard, was not bundled with IBM’s Spectrum MPI distribution;
- Compile our own `libhwloc` library to ensure GPUs are correctly discovered and used in our scheduling (The system-provided `libhwloc` was misconfigured such that its discovered GPU was not marked as Co-Processor, an attribute required for any scheduler to identify its element as a schedulable compute entity);
- Create an MPI plug-in within Flux for IBM Spectrum MPI to hide the passing of various environment variables

to each MPI job in order to assist its bootstrapping.

While they require some communications with IBM, once the proper porting path is set, implementing these required changes was trivial.

Once Flux has been ported, porting Merlin code to Flux on the new platform required only minimal changes. While Merlin still uses Sierra’s resource manager specific launcher (`jsrun`) to bootstrap a Flux instance per each batch allocation, once the instance is bootstrapped, Merlin uses the same Flux API and commands to perform its workflow. Further, Flux has been installed in public locations on both environments to further assist other workflows with portability.

B. Scheduling Specialization Addresses Challenges in Pilot2

Figure 1b shows how the Flux infrastructure interacts with the rest of the Cancer Moonshot Pilot2 workflow. Its workflow manager instantiates the machine-learning infrastructure, which implements the latent space, and uses Maestro [24] to start and stop jobs accordingly. To handle the volume of jobs and the required co-scheduling of resources, the team developed a Maestro adapter to Flux.

The workflow manager is closely coupled to the continuum simulation and constantly receives all RAS patches. Each patch is transformed into the latent space and the workflow manager maintains a priority queue of the top n candidate patches as these appear. When new resources become available, the queue is re-evaluated and a set of three new, interdependent jobs are scheduled.

This means that the primary scheduling objective they require from Flux is a simple first-come, first-served (FCFS) policy tailored for high-throughput workload. Leveraging Flux’s ability to specialize the scheduling policy and parameters for each instance, we instantiate the preexisting FCFS scheduling plugin with a scheduling parameter that further optimizes the scheduler performance for high-throughput workload. We specifically set the depth of the queue to one so that the scheduler does not have to look ahead for later jobs to schedule, an optimization of the FCFS policy that can improve resource utilization without having to break the definition of FCFS. (If the blocked highest priority job requires a compute node without GPU while the next job requires a node with GPU, the latter job can be scheduled without affecting the schedulability of the first job.)

Production runs of the Pilot2 workflow have used this flexibility to tune the scheduler for higher performance for the particular workload. It is tuned to consider fewer jobs when making a decision of what to run, which would be inappropriate for a center-wide scheduler that must maintain fairness, but that provides significant performance benefits for this application and its high quantity of concurrent jobs.

Furthermore, as described before, Pilot2 must schedule four different kinds of jobs on CPUs only, and an additional type of job on some CPUs and GPUs. One of the four needs to be on every single node, along with an instance of the fifth job with GPUs also on every node. Upon being instantiated, Flux automatically discovers CPUs and GPUs using `libhwloc`, and uses them for scheduling. So setting the

scheduling granularity to CPU/GPU-level instead of exclusive node-level was all that is needed to support this co-scheduling requirement.

Overall, the simulation reaches a steady-state in about 1 hour and 30 minutes, at which point, all resources are occupied. The steady-state utilizes all 14000 GPUs and 154000 CPU cores on all available nodes.

VII. RELATED WORK

This section presents a summary of the existing system and user-level solutions to workflow scheduling.

A. System-level Solutions

System-level solutions can be broken down into centralized, limited hierarchical, and decentralized schedulers. Centralized schedulers use a single, global scheduler that maintains and tracks the full knowledge of jobs and resources to make scheduling decisions. This scheduling model is simple and effective for moderate-size clusters, making it the state of the practice in most cloud and HPC centers today. Cloud schedulers such as Swarm [18] and Kubernetes [19] and HPC schedulers such as SLURM [8], MOAB [10], LSF [9], and PBSPro [11] are centralized. While simple, these centralized schedulers are capped at tens of jobs/sec [25], provide limited to no support for co-scheduling of heterogeneous tasks [26], have limited APIs, and cannot be easily nested within other system schedulers.

Limited hierarchical scheduling has emerged predominantly in grid and cloud computing. This scheduling model uses a fixed-depth scheduler hierarchy that typically consists of two levels. The scheduling levels consist of independent scheduling frameworks stacked together, relying on custom-made interfaces to make them interoperable. Example implementations include the cloud computing schedulers Mesos [27] and YARN [28] as well as the grid schedulers Globus [29] and HTCondor [30]. Implementations such as HTCondor provide a multilevel scheduling called pilot job systems in which resources are first acquired by an application so that the application can schedule work into those resources directly, a variant of limited hierarchical scheduling. Efforts to achieve better scalability in HPC have resulted in this model’s implementation in some large HPC centers. For example, at LLNL multiple clusters are managed by a limited hierarchical scheduler that uses the MOAB grid scheduler on top of several SLURM schedulers, each of which manages a single cluster [31]. While this solution increases throughput over centralized scheduling, it’s ultimately limited by its shallow hierarchy and the capabilities of the scheduling frameworks used at the lowest levels. In the case of LLNL example, all of the co-scheduling, coordination, and portability limitations of SLURM still apply.

Decentralized scheduling is the state-of-the-art in theoretical and academic efforts, but, contrary to centralized scheduling, it has not gained traction. To the best of our knowledge, decentralized schedulers are not in use in any production environment. Sparrow [32], in cloud computing, and SLURM++ [33] and Swift/T [34], in HPC, are existing

decentralized schedulers. In decentralized scheduling, multiple schedulers each manage a disjoint subset of jobs and resources. The schedulers are fully connected and thus can communicate with every other scheduler. In this model, a scheduler communicates with other schedulers when performing work stealing and when allocating resources outside of its resource set (i.e., resources managed by another scheduler). Despite providing higher job throughput, decentralized schedulers suffer from many of the same problems as centralized schedulers: little to no support for co-scheduling of heterogeneous tasks and limited APIs. Additionally, cloud schedulers commonly make assumptions about the types of applications being run to improve performance. For example, Sparrow assumes that a common computational framework, such as Hadoop or Spark, is used by most of the jobs, enabling the use of long-running framework processes and lightweight tasks over short-lived processes and large application binaries [32].

B. User-level Solutions

User-level solutions can be broken down into application-level runtimes and workflow managers. Application-level runtimes work by offloading a majority of the task ingestion, scheduling, and launching from the batch job scheduler onto a user-level runtime. These application-level runtimes are typically much simpler and less sophisticated than the complex system-level schedulers described in VII-A but in exchange provide extremely high throughput. For example, CRAM provides no support for scheduling (i.e., once a task completes, the resources remain idle until all other tasks have completed), tasks requiring GPUs, or an API to query the status of tasks, but it can launch ~1.5 million tasks in ~19 minutes, resulting in an average job throughput of ~1,200 jobs/sec [35].

Workflow managers are designed to ease the composition and execution of complex workflows on various computing infrastructures, including HPC, grid, and cloud resources [36]. Example workflow managers include Pegasus [37], DAGMan [38], and the UQ Pipeline [16]. Workflows can be represented as a directed acyclic graph (DAG), as is the case with Pegasus and DAGMan, or a parameter sweep, as is the case with the UQP. Once the workflow has been specified by the user, the workflow manager handles moving data between and submitting the tasks to the various computing resources. Workflow managers provide an interface for users to track the status of their workflow, and provide portability across many types of computing infrastructures. While the use of a workflow manager can improve the overall workflow throughput by taking advantage of multiple, independent computing resources (e.g., clusters), they do not improve the job throughput or co-scheduling capabilities of any individual computing resource. Additionally, to submit and manage jobs in a portable way, many workflow managers incur expensive side-effects, such as the creation of millions of job status files [39].

VIII. CONCLUSION

Emerging scientific workflows present several system-level challenges. These include, but are not limited to, throughput, co-scheduling, job coordination/communication and portability across HPC systems. In this paper, we took a deep dive into upcoming workflows and described these four specific challenges that are becoming increasingly commonplace across modern workflows. Specifically, we show three workflow examples, the Cancer Moonshot Pilot2, the Uncertainty Quantification Pipeline, and the MLSI Merlin workflow. We then presented Flux, a hierarchical and open-source resource management and scheduling framework, as a common infrastructure that can address these challenges flexibly and efficiently. The core of Flux lies in its ability to be nested seamlessly within batch allocations created by other schedulers as well as itself. Once a hierarchy of Flux instance is created within each allocation, the rich set of well-defined, platform-independent APIs efficiently support advanced workflows that can often feature non-traditional execution patterns. Our results show the performance and functionality benefits of our approach as applied to various exascale workflow challenges. Future work involves performing diverse explorations in the directions of the workflow challenges that we presented in this paper, which includes developing a deeper understanding on the effect of scheduling specialization on more diverse sets of workflows, as well as enriching our scheduling infrastructure to support heterogeneous and multi-constraint resources with the help of an advanced data model.

ACKNOWLEDGMENT

This work was performed under the auspices of the U.S. Department of Energy by LLNL under contract DE-AC52-07NA27344 (LLNL-CONF-756663).

REFERENCES

- [1] S. H. Langer, B. Spears, J. L. Peterson, J. E. Field, R. Nora, and S. Brandon, "A hydra uq workflow for nif ignition experiments," in *Proceedings of the 2Nd Workshop on In Situ Infrastructures for Enabling Extreme-scale Analysis and Visualization*, ser. ISAV '16. Piscataway, NJ, USA: IEEE Press, 2016, pp. 1–6. [Online]. Available: <https://doi.org/10.1109/ISAV.2016.6>
- [2] J. L. Peterson, "Machine learning aided discovery of a new nif design," Lawrence Livermore National Laboratory, August 2018.
- [3] D. Wang1, X. Luo2, F. Yuan1, and N. Podhorszki, "A data analysis framework for earth system simulation within an in-situ infrastructure," *Journal of Computer and Communications*, vol. 5, no. 14, pp. 76–85, Dec. 2017. [Online]. Available: <http://www.scirp.org/journal/doi.aspx?DOI=10.4236/jcc.2017.514007>
- [4] M. Dorier, J. M. Wozniak, and R. Ross, "Supporting task-level fault-tolerance in hpc workflows by launching mpi jobs inside mpi jobs," in *Proceedings of the 12th Workshop on Workflows in Support of Large-Scale Science*, ser. WORKS '17. New York, NY, USA: ACM, 2017, pp. 5:1–5:11. [Online]. Available: <http://doi.acm.org/10.1145/3150994.3151001>
- [5] D. Higdon, R. Klein, M. Anderson, M. Berliner, C. Covey, O. Ghattas, C. Graziani, S. Habib, M. Seager, J. Sefcik, P. Stark, and J. Stewart, "Uncertainty quantification and error analysis," U.S. Department of Energy, Office of National Nuclear Security Administration, and the Office of Advanced Scientific Computing Research, Tech. Rep., Jan 2010.
- [6] L. L. N. Laboratory, "Sierra," <https://hpc.llnl.gov/hardware/platforms/sierra>, Lawrence Livermore National Laboratory, August 2018, retrieved July 30, 2018.

- [7] O. R. N. Laboratory, "Summit," <https://www.olcf.ornl.gov/summit/>, Oak Ridge National Laboratory, August 2018, retrieved July 30, 2018.
- [8] A. B. Yoo, M. A. Jette, and M. Grondona, "SLURM: Simple linux utility for resource management," in *Proceedings of the 9th International Workshop on Job Scheduling Strategies for Parallel Processing (JSSPP)*, June 2003.
- [9] "IBM spectrum LSF," <https://www.ibm.com/>, 2017, retrieved April 03, 2017.
- [10] "The Moab workload manager," <http://www.adaptivecomputing.com/>, 2017, retrieved April 03, 2017.
- [11] "PBSPPro: An HPC workload manager and job scheduler for desktops, clusters, and clouds," <https://github.com/PBSPPro/pbspro>, Altair, 2018, retrieved August 8, 2018.
- [12] I. A. Prior, P. D. Lewis, and C. Mattos, "A comprehensive survey of ras mutations in cancer," *Cancer Research*, vol. 72, no. 10, pp. 2457–2467, 2012.
- [13] "MOOSE," <https://moose.inl.gov/SitePages/Home.aspx>.
- [14] J. N. Glosli, D. F. Richards, K. J. Caspersen, R. E. Rudd, J. A. Gunnels, and F. H. Streitz, "Extending stability beyond cpu millennium: A micron-scale atomistic simulation of kelvin-helmholtz instability," in *Proceedings of the 2007 ACM/IEEE Conference on Supercomputing*, ser. SC '07, 2007.
- [15] M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, and E. Lindahl, "Gromacs: High performance molecular simulations through multi-level parallelism from laptops to supercomputers," *SoftwareX*, vol. 1-2, pp. 19 – 25, 2015.
- [16] T. L. Dahlgren, D. Domyancic, S. Brandon, T. Gamblin, J. Gyllenhaal, R. Nimmakayala, and R. Klein, "Poster: Scaling uncertainty quantification studies to millions of jobs," in *Proceedings of the 27th ACM/IEEE International Conference for High Performance Computing and Communications Conference (SC)*, November 2015.
- [17] L. L. N. Laboratory, "Non-intrusive uncertainty quantification: Psuade," <https://computation.llnl.gov/projects/psuade-uncertainty-quantification/>, Lawrence Livermore National Laboratory, August 2018, retrieved August 3, 2018.
- [18] "Swarm: a docker-native clustering system," <https://github.com/docker/swarm>, Docker Inc., 2017, retrieved April 03, 2017.
- [19] "Kubernetes by Google," <http://kubernetes.io>, 2017, retrieved April 03, 2017.
- [20] J. Gyllenhaal, T. Gamblin, A. Bertsch, and R. Musselman, "Enabling high job throughput for uncertainty quantification on bg/q," in *IBM HPC Systems Scientific Computing User Group*, ser. ScicomP'14, Chicago, IL, 2014.
- [21] J. Gaffney, P. Springer, and G. Collins, "Thermodynamic modeling of uncertainties in NIF ICF implosions due to underlying microphysics models," *Bulletin of the American Physical Society.*, October 2014.
- [22] J. Gaffney, D. Casey, D. Callahan, E. Hartouni, T. Ma, and B. Spears, "Data driven models of the performance and repeatability of NIF high foot implosions," *Bulletin of the American Physical Society.*, November 2015.
- [23] "Inertial confinement fusion," https://en.wikipedia.org/wiki/Inertial_confinement_fusion, Wikipedia, 2017, retrieved August 22, 2017.
- [24] F. D. Natale, "Maestro workflow conductor (maestrowf)," <https://github.com/LLNL/maestrowf>, Lawrence Livermore National Laboratory, August 2018, retrieved Aug 11, 2018.
- [25] K. Wang, "Slurm++: A distributed workload manager for extreme-scale high-performance computing systems," <http://www.cs.iit.edu/~iraicu/teaching/CS554-S15/lecture06-SLURM++.pdf>, Feb 2015.
- [26] "SLURM heterogeneous jobs: Limitations," https://slurm.schedmd.com/heterogeneous_jobs.html#limitations, SchedMD, Dec 2017, retrieved August 8, 2018.
- [27] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A. D. Joseph, R. Katz, S. Shenker, and I. Stoica, "Mesos: A Platform for Fine-grained Resource Sharing in the Data Center," in *Proc. of the 8th USENIX Conference on Networked Systems Design and Implementation*, ser. NSDI'11. Berkeley, CA, USA: USENIX Association, 2011, pp. 295–308. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1972457.1972488>
- [28] V. K. Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, B. Saha, C. Curino, O. O'Malley, S. Radia, B. Reed, and E. Baldeschwieler, "Apache Hadoop YARN: Yet Another Resource Negotiator," in *Proceedings of the 4th Annual Symposium on Cloud Computing*, ser. SOCC '13. New York, NY, USA: ACM, 2013, pp. 5:1–5:16. [Online]. Available: <http://doi.acm.org/10.1145/2523616.2523633>
- [29] *Applications*, vol. 11, no. 2, pp. 115–128, Jun. 1997. [Online]. Available: <http://dx.doi.org/10.1177/109434209701100205>
- [30] T. Tannenbaum, D. Wright, K. Miller, and M. Livny, "Condor – a distributed job scheduler," in *Beowulf Cluster Computing with Linux*, T. Sterling, Ed. MIT Press, October 2001.
- [31] B. Barney, "Slurm and moab," <https://computing.llnl.gov/tutorials/moab>, Lawrence Livermore National Laboratory, August 2017, retrieved August 22, 2017.
- [32] K. Ousterhout, P. Wendell, M. Zaharia, and I. Stoica, "Sparrow: Distributed, low latency scheduling," in *Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP)*, November 2013.
- [33] X. Zhou, H. Chen, K. Wang, M. Lang, and I. Raicu, "Exploring distributed resource allocation techniques in the SLURM job management system," Illinois Institute of Technology, Department of Computer Science, Tech. Rep., 2013.
- [34] J. M. Wozniak, T. G. Armstrong, M. Wilde, D. S. Katz, E. Lusk, and I. T. Foster, "Swift/t: Large-scale application composition via distributed-memory dataflow processing," in *Proceedings of the 13th IEEE/ACM International Symposium on Cluster, Cloud, and Grid Computing*, ser. CCGrid, May 2013, pp. 95–102.
- [35] J. Gyllenhaal, T. Gamblin, A. Bertsch, and R. Musselman, "Enabling high job throughput for uncertainty quantification on BG/Q," in *IBM HPC Systems Scientific Computing User Group (ScicomP)*, May 2014.
- [36] J. Yu and R. Buyya, "A taxonomy of workflow management systems for grid computing," *Journal of Grid Computing*, vol. 3, no. 3, pp. 171–200, Sep 2005. [Online]. Available: <https://doi.org/10.1007/s10723-005-9010-8>
- [37] E. Deelman, G. Singh, M. hui Su, J. Blythe, Y. Gil, C. Kesselman, G. Mehta, K. Vahi, G. B. Berriman, J. Good, A. Laity, J. C. Jacob, and D. S. Katz, "Pegasus: a framework for mapping complex scientific workflows onto distributed systems," *Scientific Programming*, vol. 13, no. 3, pp. 219–237, Dec 2005.
- [38] P. Couvares, T. Kosar, A. Roy, J. Weber, and K. Wenger, *Workflow Management in Condor*. London: Springer London, 2007, pp. 357–375. [Online]. Available: https://doi.org/10.1007/978-1-84628-757-2_22
- [39] S. Hebrein, T. Patki, D. H. Ahn, D. Lipari, T. Dahlgren, D. Domyancic, and M. Tauffer, "Poster: Fully hierarchical scheduling: Paving the way to exascale workloads," in *Proceedings of the 29th ACM/IEEE International Conference for High Performance Computing and Communications Conference (SC)*, November 2017.