

I/O-Aware Batch Scheduling for Petascale Computing Systems

Zhou Zhou,^{*} Xu Yang,^{*} Dongfang Zhao,^{*} Paul Rich,[†] Wei Tang,[‡] Jia Wang,^{*} Zhiling Lan,^{*}

^{*} *Illinois Institute of Technology, Chicago, IL 60616, USA*

{zzhou1, xyang56, dzhao8}@hawk.iit.edu, lan@iit.edu, jwang@ece.iit.edu

[†] *Argonne National Laboratory, Argonne, IL 60439, USA*

[†]richp@alcf.anl.gov

[‡] *Google Inc, New York, NY 10018 USA*

[‡]weitang@google.com

Abstract—In the Big Data era, the gap between the storage performance and an application’s I/O requirement is increasing. I/O congestion caused by concurrent storage accesses from multiple applications is inevitable and severely harms the performance. Conventional approaches either focus on optimizing an application’s access pattern individually or handle I/O requests on a low-level storage layer without any knowledge from the upper-level applications. In this paper, we present a novel I/O-aware batch scheduling framework to coordinate ongoing I/O requests on petascale computing systems. The motivation behind this innovation is that the batch scheduler has a holistic view of both the system state and jobs’ activities and can control the jobs’ status on the fly during their execution. We treat a job’s I/O requests as periodical subjobs within its lifecycle and transform the I/O congestion issue into a classical scheduling problem. We design two scheduling policies with different scheduling objectives either on user-oriented metrics or system performance. We conduct extensive trace-based simulations using real job traces and I/O traces from a production IBM Blue Gene/Q system. Experimental results demonstrate that our design can improve job performance by more than 30%, as well as increasing system performance.

I. INTRODUCTION

We have already entered the age of petascale computing, and the insatiable demand for more computing power continues to drive the deployment of ever-growing high-performance computing (HPC) systems [1]. Today’s production systems already comprise hundreds of thousands processors [2][3] and are predicted to have millions with exascale computing by 2018 [4][5]. Along with the rapid evolution of micro-processors, the explosion of the amount of data must be considered. In this so-called Big Data era, the datasets generated by scientific applications are increasing exponentially in both volume and complexity [6][7][8][9].

While the computing systems can leverage more parallelism at an exponential rate in order to improve computing performance, the storage infrastructure performance is improving at a significantly lower rate. For example, the IBM Blue Gene/Q (BG/Q) system *Mira* at the Argonne Leadership Facility has a peak performance of 10 petaflops, which is 20 times faster than that of its predecessor *Intrepid* (0.5 petaflops), an IBM Blue Gene/P system [2][10][11]. But *Mira*’s I/O throughput increases only 3 times compared

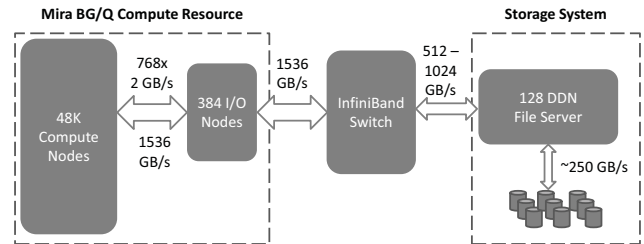


Figure 1. Storage system architecture of the 10-petaflop Mira at Argonne National Laboratory

with that of *Intrepid* [11]. As shown in Figure 1, *Mira*’s computing nodes can drive the I/O network at a full speed of 1536 GB/s whereas its storage system can deliver only 250 GB/s. This situation implies that one-quarter of the computing system can saturate the bandwidth of its storage system and incur I/O congestion, potentially degrading application performance. Consequently, the increasingly large gap between an application’s I/O requirement and the storage performance makes big data processing one of the most challenging problems facing the computer science community.

To bridge this gap, researchers have conducted numerous studies from different perspectives. A well-known approach is to boost I/O performance through I/O middleware (e.g., PLFS [12], DataStager [13], Damaris [14], IOOrchestrator [15]) or an I/O library (e.g., HDF5 [16], NetCDF [17]). As SSD becomes more cost effective, memory-class storage is proposed to cache the temporary data from compute nodes in order to mitigate I/O bandwidth contention [18][19]. Additionally, researchers propose new architectural changes to push the I/O handling closer to the compute resource by installing a burst buffer on I/O nodes [20].

These studies typically focus on application-level optimization or are implemented on I/O nodes or file servers. From the application’s perspective, some techniques are dedicated to individually optimizing this application’s I/O access pattern. The I/O interference among multiple applications is not taken into account. On the other hand, from the system view, most approaches handle lower-level I/O requests in an uncoordinated manner without any knowledge from the

upper-level applications [21]. For instance, when processing I/O accesses, the storage server attempts to establish a fair share of throughput among multiple concurrent applications, which may lead to unexpected application performance degradation. Moreover, because of the lack of a global view of all ongoing I/O activities from different applications, important systemwide performance metrics—such as average job turnaround time and system utilization—are often overlooked.

In this work, we address the I/O congestion problem at the level of batch scheduling. More specifically, this work aims to answer the following question: If a batch scheduler is aware of ongoing I/O requests from multiple jobs, will it be able to mitigate the I/O congestion issue by coordinating different I/O requests?

We argue that batch scheduler is a good candidate to handle I/O congestion for the following reasons.

- 1) The batch scheduler is a global controller of all user jobs. It has high-level knowledge of user jobs and can initiate, suspend, terminate, or restart user jobs once they are submitted to the system [22][23].
- 2) The batch scheduler often contains a monitoring component to collect abundant information about the system and user job status (e.g., node utilization, bandwidth usage, sensor data) [24][25][26][27].

In this paper we present a novel I/O-aware scheduling framework to coordinate concurrent I/O requests for petascale computing systems. The new batch scheduler not only selects queued user jobs for execution but also coordinates job I/O activities during their execution. In the case of I/O congestion, our batch scheduler suspends certain running jobs from conducting I/O. The decision is made based on a holistic view of all the running jobs and their respective I/O activities. More specifically, this paper makes the following contributions:

- 1) We present an I/O-aware batch scheduling framework for petascale computing that encompasses job abstraction, machine model, I/O congestion model, and new I/O-aware scheduling policies. In particular, two I/O-aware scheduling policies, *conservative* and *adaptive*, are designed for the framework, each for achieving different scheduling objectives (e.g., fairness, lower job slowdown, higher system utilization).
- 2) We conduct extensive trace-based simulations by using real job traces and I/O traces from the production 10-petaflop Mira system. Experimental results demonstrate that our design improves job performance by more than 30% and also improves system performance.

The remainder of this paper is organized as follows. Section II introduces background information about the target platform. Section III describes our design of the I/O-aware batch scheduling. Section IV presents experimental

results of the scheduling study. Section V discusses related work. Section VI concludes this paper and points out future work.

II. BACKGROUND

In this section we present introduce the IBM Blue Gene/Q system at Argonne and its scheduling system.

A. Mira: The IBM Blue Gene/Q at Argonne

Mira is a 10-petaflops (peak) Blue Gene/Q system operated by Argonne National Laboratory for the U.S. Department of Energy [1]. It is a 48-rack system, arranged in three rows of sixteen racks. Each rack contains 1,024 sixteen-core nodes, for a total of 16,384 cores per rack. Mira has a hierarchical structure: nodes are grouped into midplanes, each midplane contains 512 nodes, and each rack has two such midplanes. Each node has 16 cores, giving a total of 786,432 cores. Mira was ranked 5th in the latest Top500 list [1]. Mira is a capability system, with single jobs frequently occupying substantial fractions of the system. The smallest production job on Mira occupies 512 nodes; 8,192-node and 16,384-node jobs are common on the system; and larger jobs occur frequently. Jobs up to the full size of Mira run without administrator assistance. Time on Mira is awarded primarily through the DOE Innovative and Novel Computational Impact on theory and Experiment (INCITE) program [28] and the ASCR Leadership Computing Challenge (ALCC) program [29].

B. I/O Infrastructure of Mira

As shown in Figure 1, the IBM Blue Gene/Q system comprises three major components. The first component is the compute resources that comprise compute and I/O nodes. Compute nodes are responsible only for application execution and do not directly communicate with the outside. I/O nodes, on the other hand, are not involved in applications and only send/receive data to/from the second component—a large number of high-performance switches. The third component, which is connected by the second component of storage resources, comprises file servers where the application's data are persistently stored.

The I/O bandwidth between three components is also illustrated: both compute and storage resources can transfer more than 1 TB/S data via the high-performance InfiniBand switches. Yet, the aggregate disk I/O throughput within the file servers is only 250 GB/S, which is 4X slower than the outside I/O bandwidth. Therefore, the file server disk throttles the overall performance of data-intensive applications.

C. Cobalt: Batch Scheduler on Mira

Our work is based on Cobalt [26][30], a production batch scheduling system used on Mira, as well as its previous generations of HPC systems at Argonne National Laboratory. It uses a scheduling policy called “WFP” to order the

jobs in the queue [30]. WFP favors large and old jobs, adjusting their priorities based on the ratio of their wait times to their requested runtimes. On Mira, Cobalt uses a partition-based resource allocation scheme that allocates computing resource to jobs in an exclusive manner [30]. Cobalt components correspond to pieces of functionality in resource management systems, such as scheduling, queue management, hardware resource management, and process management. Its component architecture allows easy replacement of key software functionality.

In this work, we present an I/O-aware framework for Cobalt extension (see Section III). While we target Mira in this work, the idea can be easily extended to other petascale systems and their batch schedulers.

III. METHODOLOGY

In this section, we present our methodology including problem formulation, I/O-aware scheduling framework, and two new scheduling policies.

A. Problem Formulation

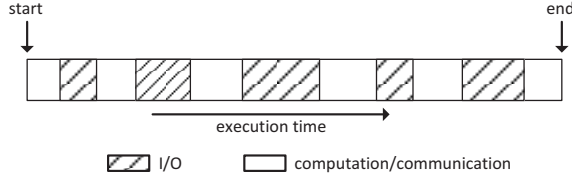


Figure 2. Lifecycle of an HPC scientific application

1) **Job abstraction:** A typical HPC job often contains three types of activity: computation, communication, and I/O. Figure 2 shows the lifecycle of a sample job. The job's runtime is split into phases for different activities, which repeatedly run during the job lifecycle. The horizontal axis represents the execution time. The height usually represents the job's resource requirement (e.g., nodes, memory, bandwidth). This abstraction provides a high-level view of the application behavior. For example, the I/O chunk in Figure 2 may contain several consecutive I/O calls (e.g., a loop of write accesses). In our job abstraction, this series of I/O calls is regarded as a single I/O request. Notice that our target platform is the IBM BG/Q system, which uses a partition-based resource allocation scheme [30]. The computation and network resource in a partition are dedicated to serve the job running on it. Hence, once a job is scheduled to run, the time consumed in computation and communication is fixed. As such, we unify computation and communication into one type of activity in Figure 2. The time for I/O activity is substantially variable because of potential I/O congestion.

In Figure 3 we depict the runtime phases of a job. Each running job J_i is associated with two basic attributes: *start time* as t_i^{start} and *size* (in nodes) as N_i .

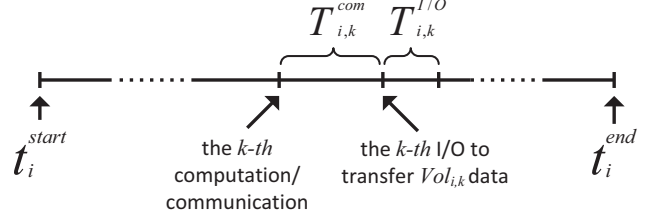


Figure 3. Runtime phases of the i th job J_i requiring N_i nodes

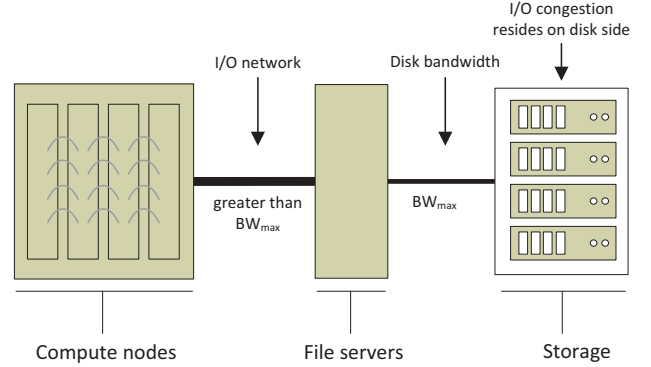


Figure 4. Simplified I/O infrastructure model of petascale HPC systems

A job shows a periodical running pattern that computation/communication interleaves with I/O. We assume that each computation/communication is followed by an I/O request that transfers a chunk of data from or to the storage system via the I/O network. Therefore, each job J_i consists of n_i computation/communication and n_i I/O activities. The k th computation/communication activity of job J_i requires time $T_{i,k}^{com}$ to finish, and the k th I/O needs to transfer $Vol_{i,k}$ of data (in gigabytes).

2) **System model:** As shown in Figure 4, we assume an HPC system composed of N compute nodes, each installed with the same number of identical uniform-speed processors. Each compute node is connected to the I/O nodes and can transfer data to the storage system across the I/O network at maximum bandwidth b (in gigabytes per second). The aggregate bandwidth of the whole system $b \times N$ is guaranteed to be less than the maximum capacity of the I/O network. Inside the storage system, we assume a group of file servers connected to a centralized parallel file system with a maximum bandwidth BW_{max} . In many systems such as Intrepid and Mira, this maximum bandwidth BW_{max} is always less than the compute nodes' aggregate bandwidth $b \times N$ and is usually subject to the upper-bound maximum speed of the hard drives. More specifically, the rotation speed of the disk heads dominates the actual access speed of the storage system. Therefore, in this model we assume that the I/O congestion takes place at the storage side if the aggregated bandwidth of all active compute nodes $b \times N_{active}$ exceeds

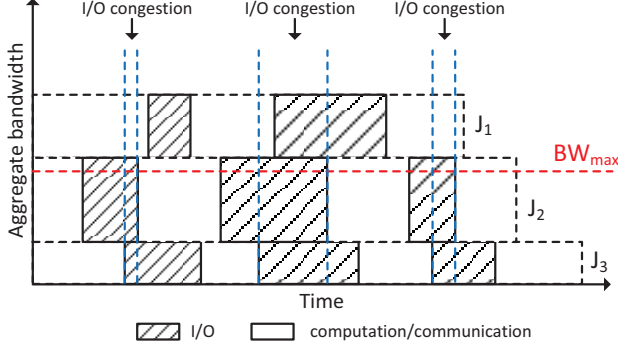


Figure 5. Scheduling model

the total bandwidth capacity BW_{max} of the storage system.

3) **Scheduling model:** Figure 5 shows the scheduling model regarding I/O congestion based on the aforementioned job abstraction and the system model. The x-axis represents the execution time, the y-axis the aggregated bandwidth. Suppose three jobs J_1 , J_2 , and J_3 are running concurrently on separated sets of nodes and each job has its own I/O pattern along with different bandwidth usage, illustrated as the height of each rectangle. Naturally a job's I/O operations may overlap with those of others. At this point, multiple jobs are accessing the storage system simultaneously, and inevitably they have to compete for I/O bandwidth with one another. If the aggregate bandwidth of these three jobs exceeds the maximum bandwidth BW_{max} denoted by the red dashed line, I/O congestion occurs and impacts these jobs' I/O performance if no congestion control is taken.

This scheduling model is simple, but it reflects the motivation of this work: to utilize the global job scheduler in order to mitigate I/O congestion by monitoring and controlling the jobs' I/O operations on the fly. After polling information about all the ongoing I/O operations, the job scheduler determines the bandwidth usage among jobs within a certain length of time. If we treat all I/O operations as subjobs in a normal job, this I/O controlling issue can be transformed into a classical job scheduling problem that can be solved by using practical approaches.

4) **I/O congestion model:** In the system model of Figure 4, we assume that the I/O network bandwidth is much larger than that of the storage system. Thus, the I/O congestion actually locates at the hard disk side, which is limited by the disk rotation speed. The lower-level scheduler in the storage server may apply a simple, independent "first-in-first-out" policy on incoming I/O requests or a more elaborated strategy, such as minimizing disk-head movements by aiming at better data locality. From the system view, the storage system is accessed in a fair-sharing mode. We define a general I/O congestion model on HPC systems as follows: If $b \times N_{active} \geq B$, the actual bandwidth allocated to each compute node is $\frac{B}{N_{active}}$. Although extra overhead such as

the time for disk-head movements should be considered, it is hard to precisely calculate the overhead in an analytical model. Therefore, in this work we set up a fair-sharing storage model with no extra overhead.

B. I/O-Aware Batch Scheduling

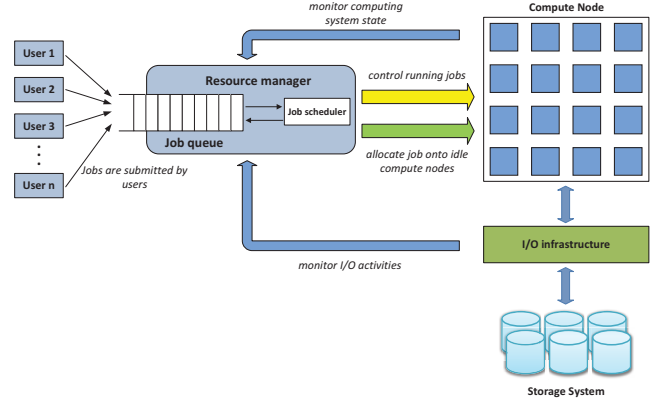


Figure 6. Our I/O-aware scheduling framework with two new features: runtime I/O monitoring (blue arrow) and dynamic control of running jobs (yellow arrow)

To mitigate I/O congestion and improve system performance, we designed a novel I/O-aware batch scheduling framework to coordinate concurrent I/O requests. As shown in Figure 6, the new scheduling framework extends the current HPC scheduling framework by integrating the runtime monitoring of system state and I/O activities and runtime I/O control. The key idea of our solution works as follows:

- 1) Every time a new I/O request is forwarded to the I/O server, the batch scheduler monitors all I/O requests being processed and checks whether the aggregate bandwidth requirement exceeds the maximum value BW_{max} .
- 2) In the case of I/O congestion, the I/O-aware batch scheduler dynamically acts on the runtime information to coordinate these concurrent I/O requests.
- 3) The runtime decision is made based on certain scheduling objectives that will be specified in the next subsection. The basic idea is to dynamically coordinate concurrent jobs' I/O activities (e.g., running or suspending) for avoiding I/O congestion.

C. I/O-Aware Scheduling Policies

Suppose K jobs in the system are performing I/O or are ready to perform I/O. Our scheduler selects a subset from the K jobs to perform I/O. Jobs that are not selected are suspended and wait until the next scheduling cycle. In this work, we propose two scheduling policies, each focusing on different objectives, either user oriented or system oriented. The selection of the subset of jobs involves searching for the optimal solution to maximize or minimize a certain objective.

1) *Job performance quantification*: First we present two user-oriented metrics to quantify job performance regarding I/O congestion.

- **InstantSlowdown (InstSld)**: Suppose at time t , job J_i is performing the k th I/O operation. It has a total of $Vol_{i,k}$ data to transfer and has already transferred $W_{i,k}$ data. Let $t_{i,k}^{I/O}$ be the start time of the k th I/O operation. We define $InstSld_{(i,k)}(t)$, the instant slowdown of the k th job J_i at time t as

$$InstSld_{(i,k)}(t) = \frac{b \times N_i \times (t - t_{i,k}^{I/O})}{W_{i,k}}, \quad (1)$$

where $t - t_{i,k}^{I/O}$ is the elapsed time from the start time of current I/O operation to now and $b \times N_i \times (t - t_{i,k}^{I/O})$ is the theoretically maximum total size of data job J_i could transfer by now within this I/O operation, in other words, an ideal case without any I/O congestion. The ratio of this maximum size of data to the already transferred $W_{i,k}$ data represents the slowdown of transferring data caused by I/O congestion. Obviously, $InstSld_{(i,k)}(t) \geq 1$ with $InstSld_{(i,k)}(t) = 1$ indicating the data transferring is not interfered with by any I/O congestion.

- **AggregateSlowdown (AggrSld)**: Again, we assume job J_i is performing its k th I/O operation. We define $AggrSld_{i,k}(t)$, the aggregate slowdown of job J_i at time t as

$$AggrSld_{i,k}(t) = \frac{t - t_i^{start}}{\sum_{j \leq k} T_{i,j}^{com} + \sum_{j < k} T_{i,j}^{I/O}}, \quad (2)$$

where $T_{i,j}^{com}$ is the time spent on the j th computation/communication operation and $T_{i,j}^{I/O}$ the time of the j th I/O operation without I/O congestion. Here $t - t_i^{start}$ equals the total elapsed time between the job start and the present time; and $j \leq k$ is the number of computation/communication or I/O the job has executed at time t since the start time t_i^{start} . The $AggrSld$ reflects the extent of an application's delay due to I/O congestion.

2) *Scheduling policies*: We propose two types of scheduling policies: *conservative* and *adaptive*.

Conservative: A conservative scheduling policy always obeys a basic principle: I/O congestion should be avoided as much as possible. Thus the scheduler selects only a subset of jobs whose aggregate bandwidth is no greater than the maximum storage bandwidth BW_{max} .

- **Cons-FCFS**: The *FCFS* (first-come-first-serve) policy works the way a traditional job scheduler does. It chooses jobs in chronological order based on the start time of concurrent I/O requests. Jobs performing I/O

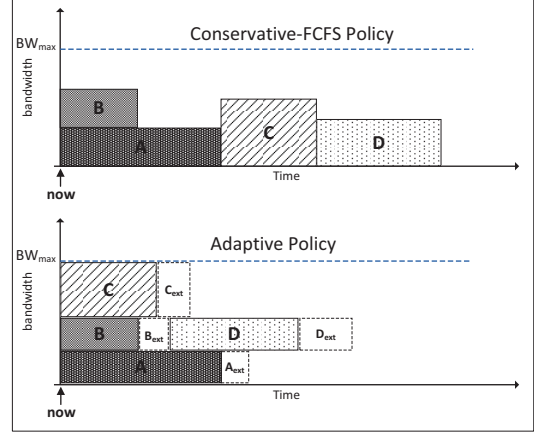


Figure 7. Using the adaptive scheduling policy can more effectively improve job performance, although I/O congestion may incur.

with earlier start time have higher priority and are favored by the scheduler.

- **Cons-MaxUtil**: The *MaxUtil* policy aims at maximizing the system utilization under the storage bandwidth constraint. This can be formalized to a classical optimization problem as follows: *Select a subset of jobs such that their aggregate bandwidth usage doesn't exceed the maximum bandwidth capacity BW_{max} , with the objective of maximizing the number of compute nodes allocated to these jobs.* In our previous work [31] [32], we transformed such problems into a standard 0-1 knapsack model, which can be solved in pseudo-polynomial time by using a dynamic programming method. The objective is system oriented, since it seeks to achieve maximized utilization of compute resource.
- **Cons-MinInstSld**: The *MinInstSld* policy favors jobs with low *InstSld* value. This method is close to *Cons-FCFS*, which orders jobs before making scheduling decision.
- **Cons-MinAggrSld**: Similar to the *MinInstSld*, the *MinAggrSld* favor jobs with low *AggrSld* value.

These four conservative policies are designed for different objectives. *Cons-FCFS* focuses on preserving the user fairness, and *Cons-MaxUtil* focuses on improving system utilization. Both *Cons-MinInstSld* and *Cons-MinAggrSld* aim to minimize the slowdown cause by I/O congestion.

Adaptive: While conservative scheduling seeks to minimize I/O congestion, it may impact job performance. Hence we propose an adaptive scheduling policy that is not restricted to the I/O bandwidth. The adaptive scheduling policy allows more jobs to perform I/O operations, despite the fact that this approach may break the I/O bandwidth bound.

Figure 7 illustrates the difference between the adaptive and conservative policies. In this example, we assume that I/O requests A and B are already scheduled and requests C and D have just arrived. Both requests C and D require more

bandwidth than what is currently available in the system. According to the *Cons-FCFS* policy, requests *C* and *D* will have to wait until enough bandwidth becomes available. This wait time leads to wasted bandwidth and unnecessary performance degradation. Instead, the adaptive scheduling will calculate the overhead of immediately scheduling request *C*, which shares the bandwidth with requests *A* and *B*. The overhead is expressed as the extension of time to finish the I/O request (i.e., square of dashed line labeled with “ X_{ext} ” in Figure 7). If the overhead is lower than suspending I/O request *C*, request *C* is allowed to start. Although this allows the job to perform I/O in contention with other jobs and may slow their performance, it improves the job turnaround time and better utilizes available bandwidth.

Algorithm 1 Adaptive Schedule

Input: S as the set of jobs performing or ready to perform I/O

Output: S_{opt} as the subset of jobs from set S to be allowed to continue or start their I/O

```

1: function ADAPTIVESCHEDULE( $S$ )
2:   Prioritize jobs in  $S$  based their I/O start time in
   FCFS fashion
3:    $BW_{avail} \leftarrow BW_{max}$ 
4:    $S_{opt} \leftarrow \emptyset$ 
5:   for  $J_i \in S$  do
6:      $BW_{Req} \leftarrow b \times N_i$ 
7:     if  $BW_{Req} \leq BW_{avail}$  then
8:       add  $J_i$  to  $S_{opt}$ 
9:        $BW_{avail} \leftarrow BW_{avail} - BW_{Req}$ 
10:    else
11:      Find the earliest time  $T_i$  can start I/O if not
      schedule  $J_i$  now
12:      Calculate  $T_{FCFS}$  as the average time need
      to finish I/O for jobs in  $S_{opt}$  plus  $J_i$ 
13:      Calculate  $T_{Adaptive}$  as the average time on
      I/O if  $J_i$  is allowed to compete for bandwidth with jobs
      in  $S_{opt}$ 
14:      if  $T_{Adaptive} < T_{FCFS}$  then
15:        add  $J_i$  to  $S_{opt}$ 
16:         $BW_{avail} \leftarrow 0$ 
17:      end if
18:    end if
19:  end for
20:  return  $S_{opt}$ 
21: end function

```

The pseudocode of the adaptive policy is shown in Algorithm 1. In the case of I/O congestion, the procedure takes as input the set of jobs performing or ready to start I/O. It first sorts jobs based on the start time of current I/O requests. The earlier it starts, the higher the priority it has (Line 2). Then the batch scheduler selects jobs in descending order of their priorities as long as the remaining

available bandwidth can satisfy its requirement (Lines 5–9). This step works the same as the FCFS policy does. If the remaining available bandwidth is insufficient, the batch scheduler decides whether to still schedule this job. It calculates the average time needed to finish current I/O requests based on the choice of scheduling this job or not (Lines 11–13). If this job is not scheduled, the earliest possible start time is calculated. Otherwise, this job will share the bandwidth with jobs already chosen in the subset S_{opt} . Obviously, other jobs’ I/O will be impacted and thus need more time to finish. If the cost is acceptable (Line 14), this job is marked to for schedule (Line 15). This adaptive policy is a runtime optimization version resembling the conservative backfilling.

IV. EVALUATION

A. Qsim Simulator

Qsim is an event-driven scheduling simulator for Cobalt [26][30]. Taking the historical job trace as input, Qsim quickly replays the job scheduling and resource allocation behavior and generates a new sequence of scheduling events as an output log. Qsim uses the same scheduling and resource allocation code that is used by Cobalt and thus has been proven to provide accurate resource management and scheduling simulation [30][33][34]. Qsim is open source and available along with the Cobalt code releases [26].

B. Job Trace and I/O Trace

We evaluate our design by means of actual workload traces from the production Mira machine. We have collected a three-month job trace from 2014. The job trace provides rich information for our simulation, including submission time, job size, duration, and walltime.

To obtain information of I/O activities, we use the Darshan log collected from Mira [35]. Darshan is a user-level system tool to allow users to characterize the I/O behavior in an efficient and transparent manner at extreme scales. In essence, for each every application running on the systems, Darshan collects the I/O footprint and summarizes it in a compact and statistically reproducible manner. Because it stores only the statistical summary of many I/O operations, the space and I/O overhead when enabling Darshan logs are minimum. By default, Darshan is turned on for all applications on the Blue Gene systems and is completely transparent to the users; no application change is in need. The effectiveness of Darshan logs has been demonstrated at extreme scales—up to 64K nodes.

By combining the job trace and I/O trace via pairing, we have built a comprehensive workload that contains job information (e.g., submission, runtime, size) as well as their I/O characteristics (e.g., number of I/O calls, I/O time, read and write transfer size). We divide the 3-month workload trace into three single-month workload traces and evaluate our design on each of them. This approach enables

us to evaluate our design under workloads with different characteristics.

C. Evaluation Metrics

Three widely used metrics are evaluated:

- *Average job wait time.* This metric denotes the average time elapsed between the moment a job is submitted and the moment it is allocated to run. It is commonly used to reflect the “efficiency” of a scheduling policy.
- *Average response time.* This metric denotes the average time elapsed between the moment a job is submitted and the moment it is completed. Similar to the above metric, it is often used to measure scheduling performance from the user’s perspective.
- *System utilization.* System utilization rate is measured by the ratio of busy node-hours to the total node-hours during a given period of time [36] [37]. The utilization rate at the stabilized system status (excluding warm-up and cool-down phases of a workload) is an important metric of how well a system is utilized.

D. Results

In this work, we compare our design with a default scheduling policy as the baseline. This *BASE_LINE* policy allocates the I/O bandwidth among multiple applications in a fair sharing manner. In the case of no I/O congestion, each application will have the maximal I/O bandwidth it needs; in the case of I/O congestion, the *BASE_LINE* policy will evenly distribute the I/O bandwidth among the concurrent applications. This approach is similar to current I/O schedulers using a round-robin method in HPC systems [15][38].

Figure 8 presents the average wait time of three workloads separately. We first observe that both our designated conservative and adaptive scheduling policies can reduce the wait time compared with the baseline result. The only exception occurs on workload 2 using *MIN_INST_SLD*, which increases the wait time by less than 10%. This proves our motivation that a scheduling policy aware of I/O-congestion is much more efficient than the *BASE_LINE* policy, which never controls I/O-congestion. Second, on all three workloads our *ADAPTIVE* policy has better results in reducing wait time than do the other scheduling policies. On both workloads 1 and 2, the *ADAPTIVE* policy reduces the wait time by at least 30%. These four conservative policies using strict I/O-congestion control lessen the flexibility of bandwidth utilization and thus are less efficient than the *ADAPTIVE* policy. Allowing a certain degree of I/O contention under the *ADAPTIVE* policy can sometimes provide a shorter wait time provided that the overall job runtime extension is smaller than the cost of waiting in the queue. Third, these three scheduling policies (*MIN_INST_SLD*, *MIN_AGGR_SLD*, and *ADAPTIVE*) perform better than the other two policies (*FCFS* and *MAX_UTIL*). The latter two policies are either user fairness or system utilization oriented,

whereas the former three take account of the turnaround time when doing scheduling. Moreover, *MIN_AGGR_SLD* outperforms *MIN_INST_SLD* because the latter focuses only on the local optimum, which is unable to guarantee a global optimum for wait time reduction.

Figure 9 presents the results for the average response time. Similar to Figure 8 on wait time, we see that both *MIN_AGGR_SLD* and *ADAPTIVE* have lower response times than do the other four. The largest reduction on response time is achieved on workload 3 by more than 30% with *ADAPTIVE*; and *MIN_AGGR_SLD* can reduce the response time by more than 20%. On workload 2 we see that *MIN_AGGR_SLD* works even better than *ADAPTIVE*, both of which have improvement up to 15%. We also observe that using *FCFS*, *MAX_UTIL*, and *MIN_INST_SLD* cannot always improve the response time and may even have a negative impact on it. For example, on workload 2 we note that *MAX_UTIL* increases the response time by nearly 10%. In particular, *FCFS* and *MAX_UTIL* are ineffective on lowering the response time, showing nearly the same value as the baseline result. The same phenomena is noticed in Figure 8. According to the definition, the job response time consists of wait time and job runtime. While the decrement of wait time benefits from our I/O-aware scheduling policies, a job may not be assigned with sufficient bandwidth. This situation increases the running time, which is likely to diminish the benefit on wait time.

Figure 10 shows the system utilization on three workloads. For convenience, we normalize the absolute utilization value based on the results of *BASE_LINE*. *MAX_UTIL* is system utilization oriented, and we can clearly see that it improves the utilization on workloads 1 and 3 up to 10%. This indicates that *MAX_UTIL* is more efficient in optimizing the system performance than are the other scheduling policies. The largest utilization drop is seen with *FCFS* and *MIN_INST_SLD* on workload 2, where the utilization decreases by nearly 10%. Other scheduling policies have negligible loss on utilization, which is treated as acceptable cost. Taking all the results into consideration with Figures 8 and 9, the *ADAPTIVE* policy is the best trade-off and achieves a good balance between user satisfaction and system utilization.

E. Sensitivity Study

The workloads used for these simulations are from the real job trace and Darshan log collected from Mira. To assess the scheduling performance under different I/O intensiveness, we also conducted a sensitivity study by tuning job I/O time in the workload. For each simulation, we define an expansion factor (*EF*) as the change of I/O time. For example, *EF*=30% means that each job’s I/O time is compressed to 30%. Similarly, *EF*=150% increases each job’s I/O time by 1.5 times, which implies it needs to transfer 50% more data.

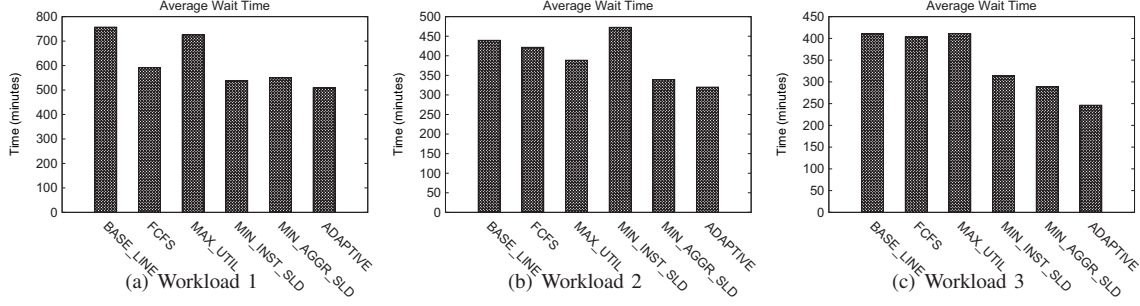


Figure 8. Average wait time

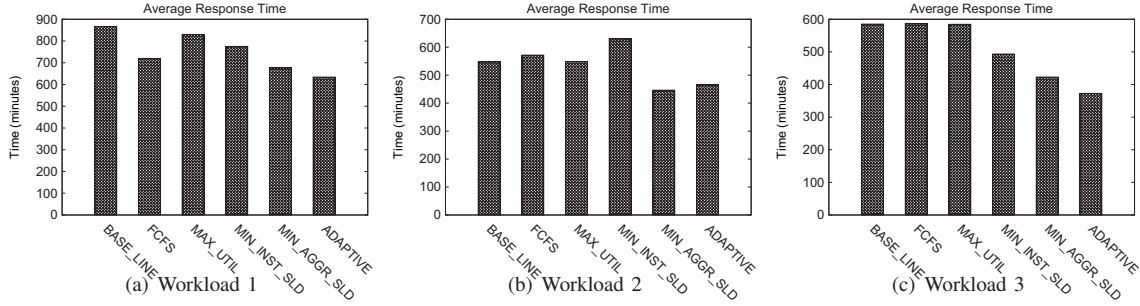


Figure 9. Average response time

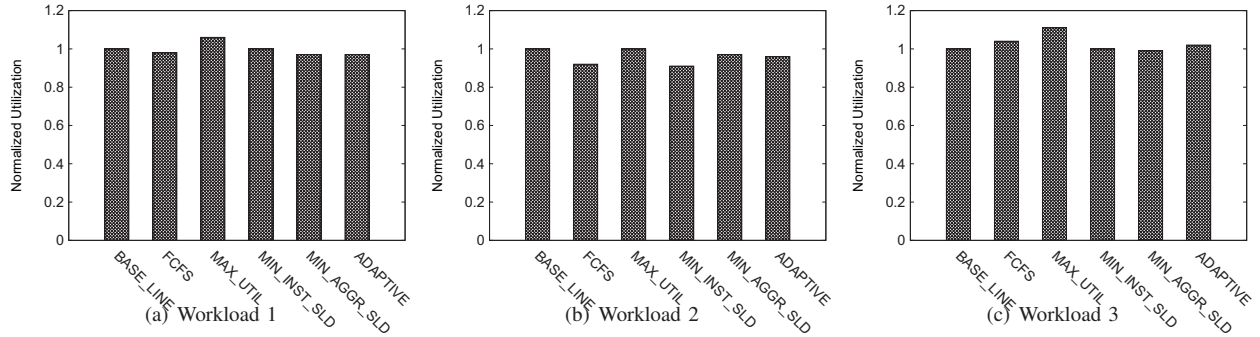


Figure 10. Normalized system utilization

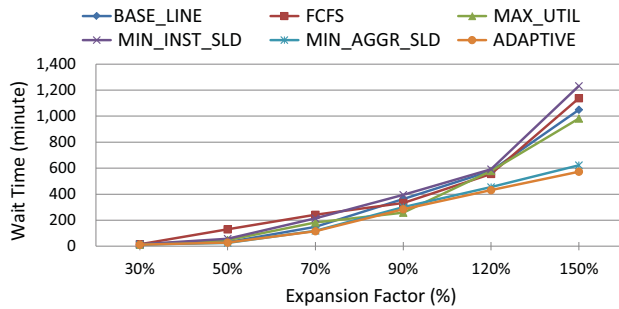


Figure 11. Impact of the I/O intensiveness over average wait time of all policies

In this way, we have built various workloads that can have “heavy” or “light” I/O activities.

Because of space limitations, we present in Figure 11 only the results for average wait time. We set up six expansion factors to compress or enlarge the I/O time. Clearly, the wait time goes up as the job spends more time on I/O. With low expansion factors (30% and 50%), no obvious improvement in wait time occurs, while FCFS even increases the wait time. As the expansion factor becomes larger, we observe that *ADAPTIVE* and *MIN_AGGR_SLD* greatly surpass other policies in reducing the wait time. When the I/O time is expanded by 150%, the largest reduction in wait time is almost 50%. Overall, the *ADAPTIVE* and *MIN_AGGR_SLD* policies should be favored for handling I/O-intensive workload because of their substantial performance improvement.

V. RELATED WORK

One traditional approach to addressing the I/O bottleneck on extreme-scale systems is to employ high-level I/O

libraries, such as HDF5 [16] and NetCDF [17]. These libraries stipulate the data format of the applications as well as their I/O interfaces. In other words, the applications need not assume the POSIX interface but manipulate the data according to the customized programming API. One limitation of this approach is portability: once the application assumes a particular API, considerable effort is needed to migrate it to other platforms.

To improve the portability of the customized API discussed above, researchers proposed several loosely coupled middleware solutions, such as PLFS [12], DataStager [13], and Damaris [14]. These systems work independently of both the underlying storage systems and the upper-level applications and thus greatly generalize the applicability. Nevertheless, an obvious drawback of this approach is the potential overhead introduced by the middleware.

More recently, the trend of solving the I/O imbalance in extreme-scale systems is to move the data closer to the compute resource. Ning et al. [20] propose moving file handling to the I/O nodes in order to ameliorate the I/O pressure from the massive number of compute nodes.

I/O contention in HPC systems draws considerable attention in the community because it is the root cause of parallel applications' performance variability [39][40][41]. Zhang et al. and Lofstead et al. schedule each application's I/O request individually without a global view from system's perspective [41] [40]. Their solutions require support from specific I/O management in the system level for better results. Solutions to the I/O contention between parallel applications have been studied in various works [42][43][44]. Hashimoto and Aida evaluate the performance variability of each job when they run concurrently on the same physical computing server[42]. The researchers show that network I/O sharing introduces most of the performance degradation. Xie et al. analyze the behavior and performance variability of Lustre, a parallel file system on supercomputer Jaguar [7]. They find that the shared filesystem between concurrently running parallel applications cause most of system performance degradation. Lebre et al. propose a new scheduling design for multiple applications with the objective of better aggregating and reordering I/O requests without hurting the fairness across applications [45]. Dorier et al. analyze the I/O interference between two applications [21]. They make a quantified study of performance improvement obtained by interrupting or delaying either one in order to avoid I/O contention. Gainaru et al. analyze the effects of interference on application I/O bandwidth and propose several scheduling techniques that apply to system I/O level to mitigate congestion [38].

VI. CONCLUSION

In this paper, we have presented an I/O-aware batch scheduling framework to alleviate the I/O congestion on petascale HPC systems. In our design, the I/O congestion

scenario is formalized into a classical batch scheduling problem by treating I/O accesses as schedulable subjobs and having the batch scheduler dynamically schedules these I/O accesses. We also have designed two types of I/O-aware scheduling policies, *conservative* and *adaptive*, each focusing on either user-oriented objectives or system performance. Our trace-based simulations clearly demonstrate the performance benefit obtained with these new scheduling policies. In particular the *ADAPTIVE* and *MIN_AGGR_SLD* policies have substantial advantage over other policies regarding user-oriented metrics. *ADAPTIVE* has better scheduling performance than *MIN_AGGR_SLD*, whereas *MIN_AGGR_SLD* has lower time complexity. Furthermore, the *MAX_UTIL* policy should be favored if the optimization for system performance (i.e., system utilization) is given priority. While this study targets Mira, our design is generally applicable to other HPC systems.

To the best of our knowledge, this is the first work on addressing I/O congestion in batch scheduling. Several avenues remain open for future work. One is to build a model to predict an application's I/O behavior based on its past I/O trace. In addition, we plan to expand this work by developing a smart resource management framework for better managing nontraditional resources including I/O and power consumption.

ACKNOWLEDGMENTS

The work at Illinois Institute of Technology is supported in part by US National Science Foundation grants CNS-1320125 and CCF-1422009. This material was based in part upon work supported by the U.S. Department of Energy, Office of Science, under contract DE-AC02-06CH11357.

REFERENCES

- [1] "Top500 supercomputing web site." [Online]. Available: <http://www.top500.org>
- [2] Mira at ANL, "<http://www.alcf.anl.gov/mira/>," Accessed September 5, 2014.
- [3] Titan at ORNL, "<https://www.olcf.ornl.gov/titan/>," Accessed September 5, 2014.
- [4] Exascale computing predicted by 2018, "<http://www.computerworld.com/article/2550451/computer-hardware/scientists-it-community-await-exascale-computers.html>," Accessed September 5, 2014.
- [5] K. Wang, K. Brandstatter, and I. Raicu, "Simmatrix: Simulator for many-task computing execution fabric at exascale," in *Proceedings of the High Performance Computing Symposium*. Society for Computer Simulation International, 2013, p. 9.
- [6] P. A. Freeman, D. L. Crawford, S. Kim, and J. L. Munoz, "Cyber-infrastructure for science and engineering: Promises and challenges," *Proceedings of the IEEE*, vol. 93, no. 3, pp. 682–691, 2005.
- [7] B. Xie, J. Chase, D. Dillow, O. Drokin, S. Klasky, S. Oral, and N. Podhorszki, "Characterizing output bottlenecks in a supercomputer," in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC '12. Los Alamitos, CA, USA: IEEE Computer Society Press, 2012, pp. 8:1–8:11.
- [8] K. Wang, A. Kulkarni, M. Lang, D. Arnold, and I. Raicu, "Exploring the design tradeoffs for extreme-scale high-performance computing system software," 2015.

- [9] —, “Using simulation to explore distributed key-value stores for extreme-scale system services,” in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. ACM, 2013, p. 9.
- [10] “Overview of the IBM Blue Gene/P project,” *IBM J. Res. Dev.*, vol. 52, no. 1/2, pp. 199–220, Jan. 2008.
- [11] Parallel I/O on Mira, “http://www.alcf.anl.gov/files/parallel_io.pdf,” Accessed September 5, 2014.
- [12] J. Bent, G. Gibson, G. Grider, B. McClelland, P. Nowoczynski, J. Nunez, M. Polte, and M. Wingate, “PLFS: A checkpoint filesystem for parallel applications,” in *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, 2009.
- [13] H. Abbasi, M. Wolf, G. Eisenhauer, S. Klasky, K. Schwan, and F. Zheng, “Datastager: Scalable data staging services for petascale applications,” in *Proceedings of the 18th ACM International Symposium on High Performance Distributed Computing*, 2009.
- [14] M. Dorian, G. Antoniu, F. Cappello, M. Snir, and L. Orf, “Damaris: How to efficiently leverage multicore parallelism to achieve scalable, jitter-free I/O,” in *Cluster Computing (CLUSTER), 2012 IEEE International Conference on*, Sept 2012, pp. 155–163.
- [15] X. Zhang, K. Davis, and S. Jiang, “iOrchestrator: Improving the performance of multi-node I/O systems via inter-server coordination,” in *High Performance Computing, Networking, Storage and Analysis (SC), 2010 International Conference for*, Nov 2010, pp. 1–11.
- [16] HDF5, “<http://www.hdfgroup.org/hdf5/doc/index.html>,” Accessed September 5, 2014.
- [17] NetCDF, “<http://www.unidata.ucar.edu/software/netcdf/>,” Accessed September 5, 2014.
- [18] F. Chen, D. A. Koufaty, and X. Zhang, “Hystor: Making the best use of solid state drives in high performance storage systems,” in *Proceedings of the International Conference on Supercomputing*, ser. ICS ’11. New York, NY, USA: ACM, 2011, pp. 22–32.
- [19] X. Zhang, K. Davis, and S. Jiang, “iTransformer: Using SSD to improve disk scheduling for high-performance I/O,” in *Parallel & Distributed Processing Symposium (IPDPS), 2012 IEEE 26th International*. IEEE, 2012, pp. 715–726.
- [20] N. Liu, J. Cope, P. Carns, C. Carothers, R. Ross, G. Grider, A. Crume, and C. Maltzahn, “On the role of burst buffers in leadership-class storage systems,” in *Mass Storage Systems and Technologies (MSST), 2012 IEEE 28th Symposium on*, April 2012, pp. 1–11.
- [21] M. Dorian, G. Antoniu, R. Ross, D. Kimpe, and S. Ibrahim, “CAL-CioM: Mitigating I/O Interference in HPC Systems through Cross-Application Coordination,” in *Parallel and Distributed Processing Symposium, 2014 IEEE 28th International*, May 2014, pp. 155–164.
- [22] K. Wang, X. Zhou, H. Chen, M. Lang, and I. Raicu, “Next generation job management systems for extreme-scale ensemble computing,” in *Proceedings of the 23rd international symposium on High-performance parallel and distributed computing*. ACM, 2014, pp. 111–114.
- [23] K. Wang, M. Lang, X. Zhou, B. McClelland, K. Qiao, and I. Raicu, “Towards scalable distributed workload manager with monitoring-based weakly consistent resource stealing,” 2015.
- [24] M. A. Jette, A. B. Yoo, and M. Grondona, “SLURM: Simple linux utility for resource management,” in *In Lecture Notes in Computer Science: Proceedings of Job Scheduling Strategies for Parallel Processing (JSSPP) 2003*. Springer-Verlag, 2002, pp. 44–60.
- [25] G. Staples, “Torque resource manager,” in *Proceedings of the 2006 ACM/IEEE conference on Supercomputing*. ACM, 2006, p. 8.
- [26] Cobalt Resource Manager, “<https://trac.mcs.anl.gov/projects/cobalt/>,” Accessed September 5, 2014.
- [27] Portable Batch System, “<http://www.pbsworks.com/>,” Accessed September 5, 2014.
- [28] “Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program.” [Online]. Available: <https://www.alcf.anl.gov/incite-program>
- [29] “ASCR Leadership Computing Challenge (ALCC).” [Online]. Available: <http://science.energy.gov/ascr/facilities/alcc/>
- [30] W. Tang, Z. Lan, N. Desai, and D. Buettner, “Fault-aware, utility-based job scheduling on Blue Gene/P systems,” in *IEEE International Conference on Cluster Computing and Workshops, 2009, CLUSTER ’09.*, 2009, pp. 1–10.
- [31] X. Yang, Z. Zhou, S. Wallace, Z. Lan, W. Tang, S. Coghlan, and M. E. Papka, “Integrating dynamic pricing of electricity into energy aware scheduling for HPC systems,” in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC’13. New York, NY, USA: ACM, 2013, pp. 60:1–60:11.
- [32] Z. Zhou, Z. Lan, W. Tang, and N. Desai, “Reducing energy costs for IBM Blue Gene/P via power-aware job scheduling,” in *Job Scheduling Strategies for Parallel Processing*, ser. Lecture Notes in Computer Science, 2014, pp. 96–115.
- [33] W. Tang, N. Desai, D. Buettner, and Z. Lan, “Analyzing and adjusting user runtime estimates to improve job scheduling on the Blue Gene/P,” in *2010 IEEE International Symposium on Parallel Distributed Processing (IPDPS)*, April 2010, pp. 1–11.
- [34] W. Tang, Z. Lan, N. Desai, D. Buettner, and Y. Yu, “Reducing fragmentation on torus-connected supercomputers,” in *Parallel Distributed Processing Symposium (IPDPS), 2011 IEEE International*, 2011, pp. 828–839.
- [35] P. Carns, K. Harms, W. Allcock, C. Bacon, S. Lang, R. Latham, and R. Ross, “Understanding and improving computational science storage access through continuous characterization,” *ACM Transactions on Storage (TOS)*, vol. 7, no. 3, pp. 8:1–8:26, Oct. 2011.
- [36] J. Jones and B. Nitzberg, “Scheduling for parallel supercomputing: A historical perspective of achievable utilization,” in *Job Scheduling Strategies for Parallel Processing*, ser. Lecture Notes in Computer Science, 1999, vol. 1659, pp. 1–16.
- [37] Y. Xu, Z. Zhou, W. Tang, X. Zheng, J. Wang, and Z. Lan, “Balancing job performance with system performance via locality-aware scheduling on torus-connected systems,” in *Cluster Computing (CLUSTER), 2014 IEEE International Conference on*, Sept 2014, pp. 140–148.
- [38] A. Gainaru, G. Aupy, A. Benoit, F. Cappello, Y. Robert, and M. Snir, “Scheduling the I/O of HPC applications under congestion,” LIP, Research Report RR-8519, Oct 2014.
- [39] J. Lofstead and R. Ross, “Insights for exascale IO apis from building a petascale IO api,” in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC’13. New York, NY, USA: ACM, 2013, pp. 87:1–87:12.
- [40] J. Lofstead, F. Zheng, Q. Liu, S. Klasky, R. Oldfield, T. Kordenbrock, K. Schwan, and M. Wolf, “Managing variability in the IO performance of petascale storage systems,” in *High Performance Computing, Networking, Storage and Analysis (SC), 2010 International Conference for*, Nov 2010, pp. 1–12.
- [41] X. Zhang, K. Davis, and S. Jiang, “Opportunistic data-driven execution of parallel programs for efficient I/O services,” in *Parallel Distributed Processing Symposium (IPDPS), 2012 IEEE 26th International*, May 2012, pp. 330–341.
- [42] Y. Hashimoto and K. Aida, “Evaluation of performance degradation in HPC applications with VM consolidation,” in *Proceedings of the 2012 Third International Conference on Networking and Computing*, ser. ICNC ’12. Washington, DC, USA: IEEE Computer Society, 2012, pp. 273–277.
- [43] D. Skinner and W. Kramer, “Understanding the causes of performance variability in HPC workloads,” in *Workload Characterization Symposium, 2005. Proceedings of the IEEE International*, Oct 2005, pp. 137–149.
- [44] A. Uelton, M. Howison, N. Wright, D. Skinner, N. Keen, J. Shalf, K. Karavanic, and L. Oliker, “Parallel I/O performance: From events to ensembles,” in *Parallel Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, April 2010, pp. 1–11.
- [45] A. Lebre, G. Huard, Y. Denneulin, and P. Sowa, “I/o scheduling service for multi-application clusters,” in *Cluster Computing, 2006 IEEE International Conference on*, Sept 2006, pp. 1–10.