

Two-level Throughput and Latency IO Control for Parallel File Systems

Yiqi Xu

Florida International University

Ming Zhao

Florida International University

Abstract

Existing parallel file systems are unable to provide both throughput and response time guarantees for concurrent parallel applications. This limitation prevents different, competing applications from getting their desired performance as high-performance computing (HPC) systems continue to scale up and be used in a shared environment. This paper presents a new two-level scheduler for parallel storage systems, a new solution to address this challenge based on a distributed performance virtualization layer for parallel file systems (vPFS). It provides both bandwidth proportional sharing and response time guarantees by addressing them at different levels of the scheduler in a cooperative manner. The utility and performance of this scheduler are studied on PVFS2, a widely used parallel file system. An experimental evaluation using a typical HPC benchmark (IOR) shows that when the storage is not overloaded, requests complete within 95th percentile response time bound during 90% of the time. The scheduler can further favor more latency-sensitive application under overloaded case.

1 Introduction

High-performance computing (HPC) systems are key to solving challenging problems in many science and engineering domains. In these systems, high-performance I/O is achieved through the use of parallel storage systems. Applications in an HPC system share access to the storage infrastructure through a parallel file system based software layer [3, 10, 2, 13]. The I/O Quality of Service (QoS) each application gets from the storage system determines how fast it can access its data and is critical to its performance guarantee. In a large HPC system, it is common to have multiple applications running at the

same time while sharing and competing for the shared storage. The sharing applications may have distinct I/O characteristics and demands which result in significant performance interference among them.

While the differentiation of applications and their respective throughput needs are addressed in previous work [14], a limitation of existing parallel storage management is the inability to satisfy both I/O bandwidth and latency needs — current work can only guarantee proportional sharing of bandwidth for the entire system with good utilization. This limitation prevents applications from efficiently utilizing the HPC resources while achieving their different desired QoS, especially for applications that are sensitive to delays such as visualization. This problem will become even more serious with the increasing scale of HPC systems and the increasing complexity and number of applications running concurrently on these systems. It presents a hurdle for the further scale-up of HPC systems to support many large, data-intensive applications.

This paper presents a new scheduler upon vPFS [14] to address these challenges through the virtualization of existing parallel file systems, achieving more comprehensive application-QoS-driven storage resource management. The virtualization layer differentiates parallel I/Os received from different application and dispatch them to the underlying storage system according to scheduling algorithms which can be created for different storage management objectives.

Specifically, this paper proposes a two-level scheduling algorithm, in which the upper level controls admission to achieve proportional bandwidth sharing, and the lower level guarantees the deadline only if the client complies with the agreed peak throughput. Proportional sharing algorithms and two-level I/O schedulers have been applied to different storage systems [12, 8, 14, 15], but, to the best of our knowledge, this paper is the first to study the effectiveness of both bandwidth and latency guarantees on typical HPC parallel storage systems.

This work is supported by National Science Foundation grants CCF-0937973 and CCF-0938045.

A prototype of the proposed two-level scheduler was developed upon PVFS2 [2], a widely used parallel file system implementation. It was evaluated with experiments using a typical HPC I/O benchmark IOR [1]. The results demonstrate that the latency guarantees can be met in addition to proportional bandwidth sharing for applications on parallel storage systems. Preliminary experiments show that requests complete within 95th percentile response time bound during 90% of the time when the storage is not overloaded and can favor more latency-sensitive application under overloaded case.

The rest of the paper is organized as follows: Section 2 introduces the motivation and background; Section 3 describes the design and implementations; Section 4 discusses the evaluation; Section 5 examines the related work; and Section 6 concludes our current work and presents future directions.

2 Background and Motivation

In a typical HPC system, applications access their data on a parallel storage system that mainly consists of a parallel file system (PFS) and its associated storage networks and devices. PFS provides the bridge between the compute and storage infrastructures which are typically connected via a high-speed network. To be cost-effective, the storage system always tends to be utilized than idle. That means the applications will contend for limited resources and be backlogged.

QoS delivered to an application in a typical storage system is usually defined as a combination of one or more of the following metrics: throughput *share* in % (or *weight*, proportion of the total bandwidth currently available in the system)), throughput *reservation* (minimum bandwidth/throughput) or throughput *cap* (maximum limit bandwidth/throughput)), and *latency* (response time). The I/O bandwidth that an application gets from the storage system determines how fast it can access its data and is critical to its QoS, while latency of an application is the major performance concern for highly interactive applications.

For example, a large data mining application processes data in bulks. Thus it may need a stable, constant throughput regardless of I/O response time; a hurricane forecast may need far more throughput during a smaller period of time due to emergency; an interactive application such as visualization software may need relatively low throughput but the requires bounded response time.

Figure 1 shows a scientific application benchmark – BTIO’s runtime decrease because of I/O interference from another benchmark IOR. The left data set is 4.7GB with large I/Os (*Large BTIO*) while the right data set is 400MB with very small I/Os (*Small BTIO*). We measure the throughput achieved during two types of BTIOs’ run

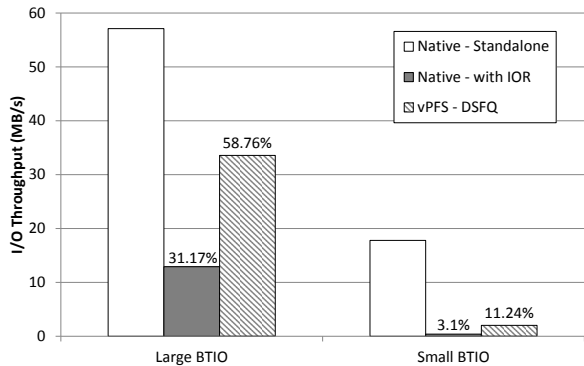


Figure 1: The slowdown of throughput to different kinds of applications due to their I/O size. Left 3 bars represent an application with 1M reads and writes. The right 3 bars represent an application with I/Os of several hundred bytes each.

when 1) it is running alone (*Native-Standalone*); 2) it is running concurrently with IOR (*Native-with IOR*); 3) it is running concurrently with IOR and vPFS framework enabling Distributed SFQ scheduler which favors BTIO (*vPFS-DSFQ*). The figure shows that Large BTIO’s throughput is slowed down by 68.83% while Small BTIO’s throughput is slowed down by 96.9%. This shows that the interactivity of Small BTIO matters more to the final application’s runtime than Large BTIO.

While we can apply SFQ-family algorithm to improve the QoS of BTIO, the increase is much more limited for the Small BTIO than for Large BTIO. The restoration of these two are 58.76% for Large BTIO while merely 11.24% for Small BTIO. The reason behind the Small BTIO’s tendency to under-perform is two fold: First, the vulnerability of performance loss in Small BTIO is because traditional proportional sharing algorithms does not recognize the deadline requirement for small I/Os. Second, the difficulty in restoring the original performance of Small BTIO using SFQ-based algorithms is because the fixed depth does not capture the available storage bandwidth to serve more concurrent I/Os.

For performance loss, the scheduler should respect deadlines of I/Os; for under-utilization, the scheduler should manage both throughput usage and latency usage. Throughput usage can be directly reflected by using request size for cost estimation, while the latency usage should be indicated by how many more I/Os from any application can be dispatched while maintaining all flows’ latency guarantees.

Current parallel storage systems offer minimal QoS for competing applications running in the compute nodes. As the above motivation example concludes, the applications’ diversity in their I/O resource demands the parallel storage system to faithfully provision both

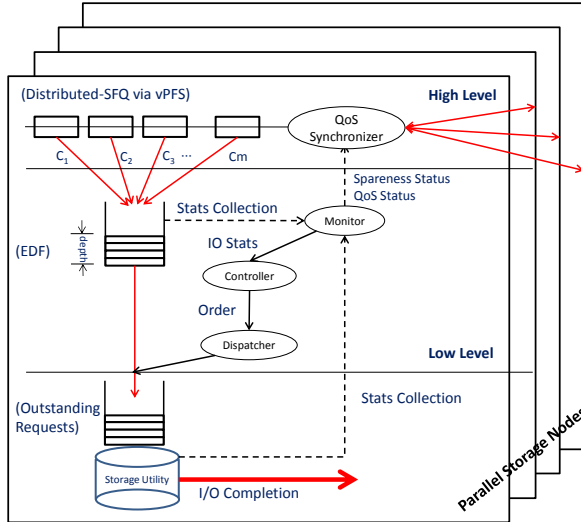


Figure 2: The distributed version of a latency-driven and throughput-driven scheduler for many parallel storage nodes via vPFS framework

throughput and latency needs for those applications. The parallelism nature of parallel storage systems also demand lightweight synchronization in order for parallel nodes to coordinate admission control together.

3 Two-level I/O Control Architecture

We use a tuple of throughput and latency for QoS definition. A (t, l) tuple indicates that the latency target l will be maintained if application's issue rate does not exceed t . Otherwise only t may be provisioned. In this scheme, throughput and latency control are both integral parts of the complete QoS targets for the applications. The architecture takes both application workload and performance requirement into consideration and uses a feedback control loop to improve the queue depth.

Our previous work proposed vPFS [14] to provide QoS-driven management of parallel storage systems through the virtualization of existing parallel file systems. It is based on 1) the capture of parallel file system I/O requests prior to their dispatch to the parallel storage system, 2) queuing of I/O request streams on a per-application basis, 3) scheduling of the queued requests according to application-specific bandwidth and response time allocations, 4) a proxy-based user-level virtualization design which enables the above parallel I/O interposition and scheduling transparently to existing parallel file system implementations and to applications. Such virtual parallel file systems can be dynamically and transparently created upon shared parallel file systems on a per-application basis, where each virtual parallel file system gets a specific *share* of the total bandwidth as a

first priority, and desired *response time* as a second priority.

Based upon the vPFS framework, we propose a two-level scheduler architecture to simultaneously control both I/O throughput and latency for applications on parallel file systems (Figure 2). The design embodies a high-level scheduler for throughput control and low-level scheduler for latency control.

3.1 High-level Scheduler for Throughput

Throughput fairness is realized by an enhanced version of DSFQ [12] on the higher level of the scheduler (Figure 2), which achieves proportional bandwidth sharing of applications' distributed I/Os by synchronizing global service information among all the participating storage nodes. Two global synchronization schemes are studied. The threshold-based broadcasting scheme proves to be effective for parallel storage as long as it uses a reasonable broadcast frequency. The layout-based scheme further takes advantage of parallel I/O application's workload characteristics: large I/Os and uniformly distributed layouts. It only needs to broadcast a layout factor (e.g., the number of servers in use) for each application to the other servers, substantially reducing the broadcast frequency. It can then infer global service using the application's file layout information without any further synchronization.

To ensure fairness, each application is first tagged in the high-level throughput controller with i ($1 \leq i \leq m$), where m is the number of applications. Within each application each request is further tagged with an monotonically increasing start tag which will determine its global order in the single dispatch queue. When the requests are continuously queued, these tags are fixed intervals apart and the interval for each application differentiates their service rate.

We re-designed DSFQ algorithm to use its previous depth D as a credit variable — the depth is not restored whenever a request completes; nor is it restored according to the number of requests completed. Instead, it only fills its depth to the fullest upon the notice of the low-level scheduler explained later. D is replenished whenever idleness is found in the storage and is equal to $\sum_{i=1}^m t_i$, and it is interpreted as throughput instead of slots.

Masking DSFQ as a credit-based rate controller has several benefits: *first*, the proportional sharing of DSFQ on vPFS is effective on parallel storage systems in terms of global service guarantee compared with a single storage node case; *second*, vPFS provides readily efficient global synchronization scheme if needed by lower level; *third*, refilling D is synchronous immediately among all applications while maintaining the fairness.

3.2 Low-level Scheduler for Latency

After receiving requests from the high level scheduler, I/O latency control is realized on the lower level of the two-level scheduler, which combines an EDF real-time scheduler with a feedback-control loop [15] in the controller. The EDF scheduler has a single adaptive queue length for all classes. According to little’s law: $queue\ length = throughput \times latency$, a small queue length helps achieve small latency but may curb the throughput thus under-utilizing the storage. A large queue length could improve utilization by allowing larger throughput but may also undermine latency bound.

3.2.1 Feedback Control Loop

The controller carves continuous service into fixed intervals called *time windows*. It manages the EDF scheduler’s behavior in each time window to achieve 95th percentile response time as well as optimal utilization. The control-feedback loop can be described briefly as follows: First, the monitor senses the actual 95th percentile response time, the number of queued requests in the EDF scheduler, and the number of requests completed by the storage. Second, the controller determines whether in the next time window the storage as a black box is overloaded or underloaded. Third, the controller changes the EDF queue’s depth based on its prediction. Fourth, the storage provides feedback from this time window’s depth change to the next time window.

3.2.2 Detailed Controller Logic

Specifically, the following three criteria are used to determine the queue length change: 1) for each class (or application), the maximum depth ${}_iL_{RT}^O$ allowed without violating the deadline; 2) the minimum depth \underline{L}^O to ensure any request with deadline falling in the time window is completed; 3) the upper bound depth \bar{L}^O for utilization if the latency need is continuously met and utilization should be raised. Define X to be number of requests completed in last time window and L^O to be the current queue depth.

If the current time window is underloaded, to obtain ${}_iL_{RT}^O$, we multiply L^O with an coefficient $E_i = \frac{D_i^O(k)}{T_i^O(k)} = \frac{D_i - MT_i^E(k+1)}{T_i^O(k)}$, where k is the current time window. $D_i^O(k)$ is the class i ’s request deadline, approximated by the difference between their class deadline D_i and the mean waiting time in the EDF queue ($MT_i^E(k+1)$) in the next time window (using last value prediction to estimate, assuming storage capacity does not change abruptly). $T_i^O(k)$ represents the 95th percentile response time of class i in time window k . Multiply \underline{L}^O and \bar{L}^O with $\frac{L^O}{X}$ to

predict capacity in the next time window. If ${}_iL_{RT}^O < \underline{L}^O$, the storage is overloaded next. Otherwise the storage stays underloaded, but it adjusts its queue depth based on the actual need. If ${}_iL_{RT}^O > \bar{L}^O$, \bar{L}^O is the candidate. If ${}_iL_{RT}^O < L^O$ or $L_{max}^O \geq L^O$, ${}_iL_{RT}^O$ is the candidate. If ${}_iL_{RT}^O \geq L^O$ and $L_{max}^O < L^O$, L^O is used, where L_{max}^O is the maximum number of outstanding requests in the current time window.

If the current time window is overloaded, E_i is multiplied by L_{max}^O . If the number of requests in the storage and those in the scheduler and arriving in the next time window is greater than X , the depth stays infinite.

textbfreminderin the previous sentence, what is exactly is the ”number”? the two ”and”s are confusing. If X is larger than number of requests that needs to be satisfied regarding their deadlines until the end of next time window, we choose the maximum between ${}_iL_{RT}^O$ and \underline{L}^O .

The final depth will be obtained from a minimal queue depth among all classes of applications to achieve both utilization and latency requirement in each time window.

3.3 Storage-level Spareness

When there is idleness in the storage, the depth is recovered to serve more requests without violating latency; when the storage is found to be overloaded, the depth is recovered more aggressively (by using *infinite* depth) to prevent cascading the further delay of subsequent requests. The recovery of depth of DSFQ in Section 3.1 is triggered when spareness is found, when request arrives in the high-level scheduler and spareness is checked, or when a *credit time window* expires.

The storage typically has variable capacity, but its idleness is calculated based on a difference between the current workload depth at the storage side and the current time window’s queue threshold (${}_iL_{RT}^O$). In underload case, if the former is less than 0.9 times the latter (to provide margin for inaccuracies), the low level monitor reports its idleness to order upper level D refilling.

4 Preliminary Evaluation

4.1 Setup

The hardware setup is as follows. We use one physical node to run two MPI applications. This node has two 2.4GHz six-core Intel Xeon CPUs, 24GB of RAM, one 500GB 7.2K SAS disk and is used to generate up to 64 MPI processes. We use another physical node to run the parallel file system PVFS2 server, which has two six-core 2.4GHz AMD Opteron CPUs, 32GB of RAM, one 500GB 7.2K RPM SAS disk. PVFS 2.8.2 is configured with the default settings. The underlying local file system is EXT3 with the default settings.

The proposed two-level scheduler was evaluated using IOR, a highly configurable I/O workload generator which simulates large-scale scientific applications performing checkpointing. The response time is measured inside IOR clients by counting the difference between the I/O completion time and I/O issue time. By tracking this latency, we can have an understanding of the feedback-control algorithm used in the scheduler.

4.2 Results

Two IOR instances run concurrently, each with 32 processes continuously writing 1MB requests. The respective service-level objectives (SLOs) for App1 and App2 are (100ms, 40MB/s) and (300ms, 20MB/s). App1 issues I/Os with an on-off pattern every 50 seconds with 50MB/s issue rate; App2 issues I/O with gradually increasing rate from 20MB/s to 50 MB/s which violates the SLO from the start. The achieved issue rates (in MB/s) of both application are shown in Figure 3. The storage’s service capacity is roughly at 50MB/s.

4.2.1 Adaptive Queue Depth

Figure 4 shows the scheduler can transit between overload and steady cases. *First*, when App1 is suddenly on, storage receives far more I/O requests than can be completed in the last time window (the storage’s capacity is roughly at 50MB/s). The scheduler is able to transit from small depths (below 20) which guarantees latency, to infinite depth (100) to gain throughput regardless of the risk of missing a small number of deadlines. *Second*, the large queue length is transitioned to small length because both throughput and latency are enforced on App2 when App1 is not on during the first 600s. This can be proved by the three drops between the first four “on” periods of App1 where depth reaches on the top twice, each followed by one drop. From the 5th “on” period after 600s, there is only one drop, because the combined issue rate from the two applications exceeds their SLO agreements and scheduler chose to favor throughput and the queue length stays high more frequently.

Figure 4 also shows that the scheduler can determine a good depth value in steady cases. *First*, at the end of first 4 “on” periods, the depth can always drop to around 15, a different value than 32 (the concurrency of App2), meaning that the scheduler can choose a reasonable storage concurrency to honor adequate throughput, regardless of the application concurrency. *Second*, at the 5th “on” period with both App1 and App2 in the system in Figure 3, the storage system is saturated. Thus, the issue rate is not able to increase anymore because each IOR issues a new request after the previous one finishes. The correspond-

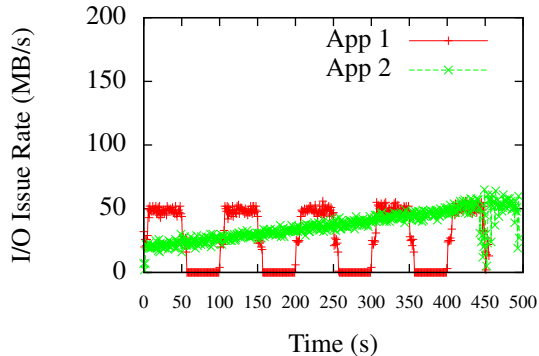


Figure 3: Issue rate for two IORs.

ing queue length in Figure 4 remains high from the 850th to the 1000th time window.

4.2.2 Latency

Figure 5 records 95th response time each second, observed from IOR clients. It shows the effectiveness of the response time control for App1 and App2 under constant overload case because the storage’s capacity is roughly at 50MB/s. Although neither App1 or App2’s issue rate exceeds storage capacity, both break their respective throughput agreement from the start. In the first 200 seconds, 90% of 95th percentile response time is within 300ms for both applications, while 100ms is far more difficult to meet for either application. However, App1 appears 10 times more than App2 within periods that are compliant to 100ms. After 200 seconds, an increasing number of periods’ 95th percentile response time go beyond 300ms because of the increased total issue rate.

To ensure work-conserving and utilization, the bandwidth proportional sharing is reflected by application’s respective agreement on throughput if they conform. In our setup, however, both applications violate their agreements. By design, whenever App1 is on, combined issue rate is more than the storage can handle well. Both applications get throughput equal to the issue rates (instead of being proportional to their agreed rates between the two applications), and do not necessarily get compliant latency, because the scheduler tries to maximize throughput to end the overloaded case as soon as possible.

Whenever App1 is off, the system can handle the issue rate of App2 until starting from the 350th second when App2’s issue rate also exceeds the storage capacity. This is confirmed by the increase in the outliers in Figure 5 after the 500s, where more I/Os’ latencies are out of bound (many of them are far away above 1000ms).

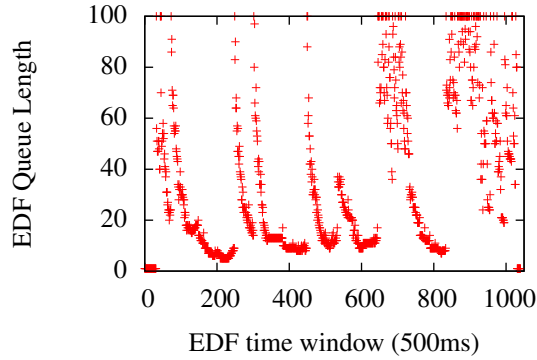


Figure 4: EDF queue length on a server.

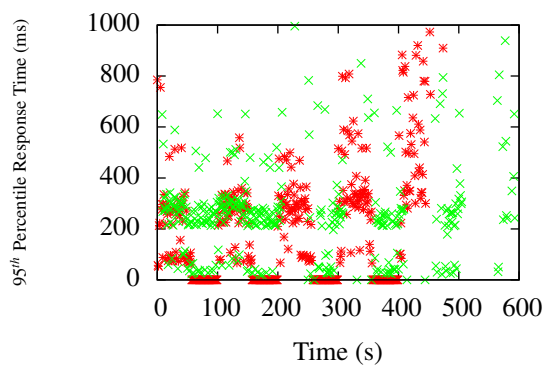


Figure 5: 95th percentile response time from two IORs.

5 Related Work

Storage resource management has been studied in related work in order to service competing I/O workloads and meet their desired throughput and latency goals. Such management can be embedded in the shared storage resources’ internal scheduler (e.g., disk schedulers) (Cello [11], Stonehenge [7], YFQ [4], PVFS [2]), which has direct control over the resource but requires the internal scheduler to be accessible and modifiable. The management can also be implemented via virtualization by interposing a layer between clients and their shared storage resources (Façade [9], SLEDS [5] [8], SFQ(D) [8], GVFS [16]). This approach does not need any knowledge of the storage resource’ internals or any changes to its implementation. It is transparent to the existing storage deployments and supports different types of storage systems. vPFS [14] takes this approach to virtualize parallel file systems and achieve proportional bandwidth sharing. This paper further proposes a new two-level scheduler upon vPFS that controls both I/O throughput and latency for applications on parallel file systems.

Various scheduling algorithms have been investigat-

ed in related storage management solutions. They employ techniques such as virtual clocks, leaky buckets, and credits for proportional sharing, earliest-deadline first (EDF) scheduling to guarantee latency bounds, feedback-control with request rate throttling, adaptive control of request queue lengths based on latency measurements, and scheduling of multi-layer storage resources based on online modeling. The effectiveness of these scheduling algorithms is unknown for a HPC parallel storage system. In particular, our two-level scheduling algorithm is inspired by the related work [15] which achieves throughput and latency control for a centralized storage system. This paper differentiates in that 1) it adopts an enhanced DSFQ algorithm as for the upper layer which provides sound, theoretically provable global fairness guarantees, 2) and it accomplishes throughput and latency control for parallel file system based storage.

The majority of the storage resource schedulers in the literature focuses on the allocation of a single storage resource (e.g., a storage server, device, or a cluster of interchangeable storage resources) and addresses the local throughput or latency objectives. LexAS [6] was proposed for fair bandwidth scheduling on a storage system with parallel disks, but the I/Os are not stripped and the scheduling is done with a centralized controller. DSFQ [12] is a distributed algorithm that can realize total service proportional sharing across all the storage resources that satisfy workload requests. However, it faces challenges of efficient global scheduling when applied to a HPC parallel storage system, which are addressed by the vPFS and the enhanced algorithms enabled upon it.

6 Conclusions and Future Work

Modern HPC systems are shared by an increasing number of types of data-intensive applications with diverse I/O workloads. Currently there is no complete solution to the faithful provisioning of those application’s I/O needs for both throughput and latency guarantees, resulting in unpredictable performance. This paper presents a two-level I/O scheduler for both throughput and latency control on a parallel file system. The scheduler was implemented upon vPFS, our previous virtualization framework on parallel storage systems. It combines DSFQ and EDF algorithms and coordinates distributed scheduling among multiple storage nodes. The experiments demonstrate that it can choose optimal queue depth for both throughput and latency. The feedback-control algorithm in the scheduler can effectively respond to workload changes along with storage capacity change.

In the future we will optimize the scheduling of I/Os of different sizes. We will also manage I/O latency aggregated from multiple storage nodes via a distributed version of EDF algorithm.

References

- [1] IOR parallel file system benchmark.
- [2] PVFS2, parallel virtualized file system.
- [3] Lustre file system: High-performance storage architecture and scalable cluster file system. White Paper, October 2009.
- [4] BRUNO, J., BRUSTOLONI, J., GABBER, E., OZDEN, B., AND SILBERSCHATZ, A. Disk scheduling with quality of service guarantees. In *Proceedings of the IEEE International Conference on Multimedia Computing and Systems - Volume 2* (Washington, DC, USA, 1999), ICMCS '99, IEEE Computer Society, pp. 400–.
- [5] CHAMBLISS, D., ALVAREZ, G., PANDEY, P., JADAV, D., XU, J., MENON, R., AND LEE, T. Performance virtualization for large-scale storage systems. In *Reliable Distributed Systems, 2003. Proceedings. 22nd International Symposium on* (oct. 2003), pp. 109 – 118.
- [6] GULATI, A., AND VARMAN, P. Lexicographic qos scheduling for parallel i/o. In *Proceedings of the seventeenth annual ACM symposium on Parallelism in algorithms and architectures* (New York, NY, USA, 2005), SPAA '05, ACM, pp. 29–38.
- [7] HUANG, L., PENG, G., AND CHIUEH, T.-C. Multi-dimensional storage virtualization. In *Proceedings of the joint international conference on Measurement and modeling of computer systems* (New York, NY, USA, 2004), SIGMETRICS '04/Performance '04, ACM, pp. 14–24.
- [8] JIN, W., CHASE, J. S., AND KAUR, J. Interposed proportional sharing for a storage service utility. *SIGMETRICS Perform. Eval. Rev.* 32 (June 2004), 37–48.
- [9] LUMB, C. R., MERCHANT, A., AND ALVAREZ, G. A. Façade: Virtual storage devices with performance guarantees. In *Proceedings of the 2nd USENIX Conference on File and Storage Technologies* (Berkeley, CA, USA, 2003), FAST '03, USENIX Association, pp. 131–144.
- [10] SCHMUCK, F., AND HASKIN, R. GPFS: A shared-disk file system for large computing clusters. In *Proceedings of the 1st USENIX Conference on File and Storage Technologies* (Berkeley, CA, USA, 2002), FAST '02, USENIX Association.
- [11] SHENOY, P. J., AND VIN, H. M. Cello: a disk scheduling framework for next generation operating systems. In *Proceedings of the 1998 ACM SIGMETRICS joint international conference on Measurement and modeling of computer systems* (New York, NY, USA, 1998), SIGMETRICS '98/PERFORMANCE '98, ACM, pp. 44–55.
- [12] WANG, Y., AND MERCHANT, A. Proportional-share scheduling for distributed storage systems. In *Proceedings of the 5th USENIX conference on File and Storage Technologies* (Berkeley, CA, USA, 2007), USENIX Association, pp. 4–4.
- [13] WELCH, B., UNANGST, M., ABBASI, Z., GIBSON, G., MUELLER, B., SMALL, J., ZELENKA, J., AND ZHOU, B. Scalable performance of the Panasas parallel file system. In *Proceedings of the 6th USENIX Conference on File and Storage Technologies* (Berkeley, CA, USA, 2008), FAST'08, USENIX Association, pp. 2:1–2:17.
- [14] XU, Y., ARTEAGA, D., ZHAO, M., LIU, Y., FIGUEIREDO, R., AND SEELAM, S. vpfs: Bandwidth virtualization of parallel storage systems. In *Mass Storage Systems and Technologies (MSST), 2012 IEEE 28th Symposium on* (2012), pp. 1–12.
- [15] ZHANG, J., SIVASUBRAMANIAM, A., WANG, Q., RISK, A., AND RIEDEL, E. Storage performance virtualization via throughput and latency control. *Trans. Storage 2* (August 2006), 283–308.
- [16] ZHAO, M., ZHANG, J., AND FIGUEIREDO, R. J. Distributed file system virtualization techniques supporting on-demand virtual machine environments for grid computing. *Cluster Computing* 9, 1 (Jan. 2006), 45–56.