



Bridging Storage Semantics Using Data Labels and Asynchronous I/O

ANTHONY KOU GKAS, HARIHARAN DEVARAJAN, and XIAN-HE SUN, Illinois Institute of Technology, Department of Computer Science

In the era of data-intensive computing, large-scale applications, in both scientific and the BigData communities, demonstrate unique I/O requirements leading to a proliferation of different storage devices and software stacks, many of which have conflicting requirements. Further, new hardware technologies and system designs create a hierarchical composition that may be ideal for computational storage operations. In this article, we investigate how to support a wide variety of conflicting I/O workloads under a single storage system. We introduce the idea of a *Label*, a new data representation, and, we present LABIOS: a new, distributed, Label-based I/O system. LABIOS boosts I/O performance by up to 17× via asynchronous I/O, supports heterogeneous storage resources, offers storage elasticity, and promotes *in situ* analytics and software defined storage support via data provisioning. LABIOS demonstrates the effectiveness of storage bridging to support the convergence of HPC and BigData workloads on a single platform.

CCS Concepts: • **Information systems** → **Distributed storage**; **Hierarchical storage management**; *Storage power management*; • **Computer systems organization** → *Distributed architectures*; *Data flow architectures*; **Heterogeneous (hybrid) systems**;

Additional Key Words and Phrases: Label-based I/O, storage bridging, heterogeneous I/O, datalabels, task-based I/O, exascale I/O, energy-aware I/O, elastic storage

ACM Reference format:

Anthony Kougkas, Hariharan Devarajan, and Xian-He Sun. 2020. Bridging Storage Semantics Using Data Labels and Asynchronous I/O. *ACM Trans. Storage* 16, 4, Article 22 (October 2020), 34 pages.
<https://doi.org/10.1145/3415579>

1 INTRODUCTION

Large-scale applications, in both scientific and the BigData communities, demonstrate unique I/O requirements that none of the existing storage solutions can unequivocally address them. This has caused a proliferation of different storage devices, device placements, and software stacks, many of which have conflicting requirements. Each new architecture has been accompanied by new software for extracting performance on the target hardware. Further, to reduce the I/O performance gap, hardware composition of modern storage systems is going through extensive changes

This material is based upon work supported by the National Science Foundation under Grants No. OCI-1835764 and No. CSR-1814872.

Authors' addresses: A. Kougkas, H. Devarajan, and X.-H. Sun, Department of Computer Science, Illinois Institute of Technology, 10 West 35th Street, Chicago, IL 60616; emails: {akougkas, hdevarajan, sun}@iit.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1553-3077/2020/10-ART22 \$15.00

<https://doi.org/10.1145/3415579>

by adding new storage devices. This leads to heterogeneous storage resources where data movement is complex, expensive, and dominating the performance of most applications [40]. For instance, machines with a large amount of RAM allow new computation frameworks, such as Apache Spark [80], to thrive. Supercomputers equipped with node-local fast storage, such as NVMe drives, take scientific simulation to new performance standards [8]. To achieve computational efficiency modern parallel and distributed storage systems must efficiently support a diverse and conflicting set of features.

Data-intensive applications grow more complex as the volume of data increases, creating diverse I/O workloads. Thus, the features a distributed storage system is required to support also increases dramatically in number and are often conflicting. For instance, scientific applications demonstrate a periodic behavior where computations are followed by intense I/O phases. Highly-concurrent write-intensive workloads (e.g., final results, checkpoints), shared file parallel access, frequent in-place data mutations, and complex data structures and formats are the norm in most High-Performance Computing (HPC) workloads [39]. However, iterative write-once, read-many data access, created by the popular MapReduce paradigm, are the defacto patterns in most BigData applications [56]. Another example is the ability of an I/O subsystem to handle data mutations. In HPC, the ability to frequently update data forces storage systems to obey certain standards, such as POSIX, and increase the cost of metadata operations, which is projected to limit the scalability of these systems [64]. In contrast, most cloud storage solutions prefer an immutable representation of data, such as RDDs [79] or key-value pairs. Finally, each application manipulates data in a different data representation (i.e., format) spanning from files, objects, buckets, key-value pairs, and so on, which increases the complexity of the data organization inside a storage system. To navigate this vast and diverse set of contradictory I/O requirements, the software landscape is filled with custom, highly specialized storage solutions varying from high-level I/O libraries to custom data formats, interfaces, and, ultimately, storage systems.

The ability to seamlessly execute different conflicting workloads is a highly desirable feature. However, the tools and cultures of HPC and BigData have diverged, to the detriment of both [63], and unification is essential to address a spectrum of major research domains. This divergence has led organizations to employ separate computing and data analysis clusters. For example, NASA's Goddard Space Flight Center uses one cluster to conduct climate simulation, and another one for the data analysis of the observation data [84]. Due to the data copying between the two clusters, the data analysis is currently conducted off-line, not at runtime. The data transfer between storage systems along with any necessary data transformations are a serious performance bottleneck and cripples the productivity of those systems [41]. Additionally, it increases the wastage of energy and the complexity of the workflow. Integrating analytics into a large scale simulation code has been proven to significantly boost performance and can lead to more accurate and faster solutions. Current storage systems address interoperability (i.e., cross-storage system data access) by adding various connectors, such as IBM's Spectrum Scale HDFS Transparency [34] and Intel's Hadoop Adapter [35], and/or middleware libraries, such as IRIS [41] and Alluxio [45]. Nevertheless, better system support is needed for in-transit, *in situ* analysis, with scheduling being a big challenge and node sharing impossible in existing solutions [57]. However, High-Performance Data Analytics (HPDA) [36], the new generation of Big Data applications, involve sufficient data volumes and algorithmic complexity to require HPC resources. For example, Paypal, an online financial transaction platform, and the US Postal Service are using HPC resources to perform fraud detection in real time on billions of transactions and mail scans. Cycle Computing tested 21 million drug candidate molecules on the Amazon public cloud using a new HPC algorithm [69]. Gaining insights from massive datasets while data is being produced by large-scale simulations can enhance the scalability and flexibility of exascale systems [82].

To address this divergence in storage architectures and workload requirements, we have developed LABIOS, a new, distributed, scalable, and adaptive I/O System. LABIOS, a new class of a storage system, is the first (data) LAbel-based I/O System, is fully decoupled, and is intended to grow in the intersection of HPC and BigData. LABIOS demonstrates the following contributions:

- (1) the effectiveness of **storage malleability**, where resources can automatically grow/shrink based on the workload;
- (2) how to effectively support **synchronous and asynchronous I/O** with configurable heterogeneous storage;
- (3) how to leverage **resource heterogeneity** under a single platform to achieve application and system-admin goals;
- (4) the effectiveness of **data provisioning**, enabling *in situ* data analytics, computational storage, and process-to-process data sharing;
- (5) how to support a diverse set of conflicting I/O workloads, from HPC to BigData analytics, on a single platform, through managed **storage bridging**.

LABIOS achieves these contributions by transforming all I/O requests each into a configurable unit called a *Label*, which is a tuple of an operation and a pointer to data. Labels are pushed from the application to a distributed queue served by a label dispatcher. This queuing system provides the ability to perform asynchronous I/O. LABIOS workers (i.e., storage servers), which can store data in any storage medium (e.g., SSD, HDD, etc.), execute labels independently. LABIOS architecture is fully decoupled and distributed, which allows the worker set to be resized, independently from the clients, based on the I/O traffic within the system. The existence of a temporary intermediate space, the distributed queuing system for labels, and the fact that workers operate independently from the clients allows LABIOS to provide active storage semantics. Using the label structure, that encapsulates an operation and its input data, LABIOS can offer software-defined storage services and QoS guarantees for a variety of workloads on different storage architectures. Finally, the modular design of LABIOS allows its workers to connect to a plethora of external storage services, and, thus, bridging the semantic gap between different storage APIs.

The remainder of the article is organized as follows. Section 2 provides a detailed view of the existing storage solutions and the complexity of managing a diverse set of application requirements. It also presents the motivation behind the design of LABIOS. Section 3 offers the details of LABIOS design, architecture, component analysis, objectives and challenges, the implementation details along with the deployment models. We also provide a discussion on current design considerations and implications as well as some of the limitations of our prototype system. Section 4 details the performance characteristics of our proposed system by experimentally evaluating each individual LABIOS component as well as the entire system as a whole. We present the related work in Section 5, and we conclude this work in Section 6 along with listing some of our future steps.

2 BACKGROUND AND MOTIVATION

2.1 Parallel and Distributed File Systems

Parallel file systems (PFS) are the dominant storage solution in most large-scale machines such as supercomputers and HPC clusters and are therefore well understood in the storage community. As the name implies, a PFS deals with data in the form of files. PFS obey the POSIX standard to offer portable guarantees and strong data consistency. PFS manipulate data in a certain sequence of operations, a paradigm known as *streamlined I/O* (i.e., Unix Standard I/O Streams). Parallel access is achieved by shared file handlers and a complex system of locking mechanisms. Through the years, PFS have been optimized to fit the needs of typical HPC workloads. Application development and

storage system design have grown in harmony with one driving the other, since the HPC ecosystem is relatively closed to external influence. However, PFS face many limitations [33]. Some relevant to this study include:

- (a) *Storage malleability.* Existing high-performance storage solutions are not elastic but static and cannot support power-capped I/O and tunable concurrency control (i.e., QoS guarantees based on job size, priority, input, output, etc.). Sudden workload variations (i.e., I/O demand fluctuations) in distributed systems can be addressed by resource malleability. By dynamically increasing or decreasing the amount of storage resources allocated to an application, the system can reduce its idle resources and therefore achieve lower energy consumption and costs for the end user.
- (b) *Resource utilization.* Storage resources are provisioned for the worst-case scenario where multiple jobs happen to enter their I/O-dominant phases simultaneously leading to over/under-provisioning. This issue is worsened by the growing need to support storage resources sharing across multiple clusters via global mounts. Furthermore, allocation exclusivity and over-provisioning due to ignorance or malicious intent also contribute to erroneous resource utilization.
- (c) *Hardware heterogeneity.* New storage devices (e.g., SSD, NVMe, etc.) are being incorporated into system designs resulting in a diverse heterogeneous storage environment. Existing solutions cannot handle this heterogeneity, since they assume homogeneous servers. Currently, the responsibility for orchestrating data movement, placement, as well as layout within and across nodes falls on both system administrators and users [52].
- (d) *Flexible interface.* Currently, storage is tightly coupled to certain vendor-specific APIs and interfaces. Even though this ensures consistency and reliability of the storage system, it can also lead to reduced productivity; developers either need to learn new APIs, which limits flexibility, or, adopt new storage systems, which leads to environment isolation. Many PFS have introduced various connectors to increase interoperability, but at the cost of performance. Moreover, existing storage systems provide limited facilities for developers to express intent in the form of I/O requirements, semantics, and performance guarantees. Consequently, to achieve good I/O performance, the level of abstraction has been raised. I/O libraries, such as HDF5 [25] and PnetCDF [47] help alleviate this issue but they also add overheads and increase the complexity of use.

In cloud environments the storage scene is different. Innovation is driven by the wide popularity of computing frameworks. As a result, the cloud community has developed a wide variety of storage solutions tailored to serve specific purposes. The most popular storage solution in deployment is the Hadoop Distributed File System (HDFS), which also follows the streamlined I/O paradigm. The architecture and data distribution are somewhat similar to PFS. In HDFS, there are metadata nodes (i.e., namenodes), which are responsible to maintain the namespace, and data nodes that hold the files. However, it has relaxed the POSIX standard to achieve scalability. As the Hadoop ecosystem grows, so are the storage solutions around it: Hive [71] puts a partial SQL interface on front of Hadoop, Pig [58] enables a scripting language in top of MapReduce, HBase [15] applies a partial columnar scheme on top of Hadoop, and HCatalog [27] introduces a metadata layer to simplify access to data stored in Hadoop, and finally, Impala [10] builds a database-like SQL layer on top of Hadoop. This diversity can offer advantages but also undoubtedly increases the complexity of storage. Some of the above solutions suffer from similar limitations as PFS [76], some others are missing critical features [43] and, in general, most of them perform well for the purpose they were designed for.

2.2 Applications' I/O Requirements

Every computing framework expects specific I/O requirements and features from the underlying storage system. Scientific computing, for instance, relies mostly on MPI for its computations-communications, and domain scientists expect POSIX, MPI-IO, and other high-level I/O libraries to cover their I/O needs. The existing collection of storage interfaces, tools, middleware libraries, data formats, and APIs is deeply instilled in the community and has created a certain mindset of what to expect from the storage stack. We present here some storage requirements representative of typical scientific workloads:

- (a) *Strong data consistency*: any data passed to the I/O system must at all times be consistent between operations. A typical example is the read after write access pattern where a set of producers pushes some data to the distributed storage and another different set of consumers reads them back. Moreover, this extends to any data mutation: data must be consistent immediately after an update operation is completed. There is a tradeoff between performance and consistency. In HPC and traditional financial computing workloads, data consistency is not negotiable. However, in many cloud and BigData workloads immutability or eventual consistency might be enough. Some research [73] has proposed tunable data consistency offered by workload, or by file, or even per-operation.
- (b) *Concurrent shared file access*: multiple processes must be able to operate on the same file concurrently. This access pattern is particularly common in HPC workloads where multiple MPI ranks open the same file and read/write data concurrently. Collective I/O optimization [70], concurrent file handlers [26], and complex locking schemes [20, 78, 81] can make this feature possible. However, most storage solutions in the cloud community do not support concurrent data access and prefer a multi-replica scheme to offer a highly available concurrent systems.
- (c) *Hierarchical global namespace*: users must be able to organize data in a hierarchy with directory support and nesting. Moreover, data identifiers (i.e., file names, directory structures, keys, etc.) must be resolved and recognizable in a global namespace that can be accessed from anywhere. This feature is fundamental to most traditional code (i.e., C, C++, Fortran, etc.), where file paths and directory nesting is common in data organization. The hierarchical organization of the global namespace creates a central point of metadata handling, which may lead to performance bottlenecks and limited scalability. Some solutions to this include partitioning of the global namespace and distribution in many servers, client-side caching, and decoupling of the metadata from the data [21, 64, 77, 83]. In contrast, cloud storage solutions offer a flat namespace boosting scalability. However, they cannot support traditional workloads and many efforts [6, 7, 68] are underway to extend cloud storage by introducing file connectors and other abstractions of file namespaces.
- (d) *High bandwidth and parallelism*: scientific discovery is often dependent on the ability of the storage infrastructure to push data into/out of the compute nodes for (near) real-time processing. I/O performance is crucial in achieving computational efficiency and is now the number one priority to most mission-critical workloads. Higher I/O bandwidth is historically the driving force and a desired feature in the development of any parallel storage infrastructure. However, more often than not, I/O bandwidth is the one commodity easily traded for the sake of other features such as fault tolerance, or security. Storage resources are external in most of the supercomputers and as data-intensive computing is now the norm rather the exception, high storage parallelism can boost applications' I/O bandwidth.

However, MapReduce and Hadoop workloads, prominent in cloud environments, discount some of the above and dictate different I/O requirements from the underlying storage that supports them. Please note that the I/O requirements mentioned in this section are not necessarily mutually exclusive but are often contradictory. We highlight some here:

- (a) *Fault tolerance*: hardware faults are expected as most cloud workloads run on data centers consisting of commodity hardware. Hence, errors have been in the front row while designing a distributed cloud storage solution. Applications expect faults to be the norm, and thus, require the storage infrastructure to be able to handle them graciously. Data replication, erasure coding, and data partitioning have been among the most popular defenses against this. However, HPC machines are built with more expensive, sophisticated, and specialized hardware that includes fault tolerant guarantees (i.e., server-graded drives, RAID, etc.), and thus, I/O systems assume a certain level of hardware reliability. This does not mean that fault tolerance is sidestepped in HPC workloads. To handle faults, applications perform frequent checkpoints and/or restart the entire workload.
- (b) *Extreme scalability and multi-tenancy*: the growth of web services gave birth to extreme scale multitenant workloads (e.g., many-task computing [62]), where multiple small applications run on a large datacenter sharing the storage resources. In contrast, HPC machines run few large scale batch jobs. Scalability is desired in both storage camps but the way an I/O system scales differs significantly. BigData workloads involve smaller iterative jobs with an extremely high number of I/O clients. Some cloud storage solutions can expand (i.e., scale-out) adding resources to meet applications' demands. In HPC, multitenancy is addressed by sophisticated job scheduling, synchronization techniques, and I/O optimizations such as data buffering.
- (c) *Data locality*: a major paradigm shift, brought by the widely used MapReduce framework has been data locality (i.e., jobs are spawned where data is). This resulted in system designs with node-local storage devices; fundamentally different than supercomputers where storage is an external remote resource. Distributed storage leverages data locality to minimize possible network bottlenecks but sacrifices global aggregate I/O bandwidth. This further highlights that each storage platform can offer higher performance to certain applications and that the diversity of workloads has led to the divergence of I/O systems.
- (d) *Ease-of-use*: the cloud community strongly advocates for the ease-of-use of their tools and frameworks whereas the scientific community aims to develop high-performance solutions at the expense of user-friendliness and ease of deployment. This is made apparent by the growth of data formats such as key-value pairs in object stores, along with their interfaces and APIs, such as Amazon's S3 [2] or Openstack Swift [55]. HPC storage requires a higher level of expertise from the developer. Ease-of-use is the number one factor of high productivity and can reduce erroneous code.

Table 1 summarizes some I/O requirements and how each storage camp, HPC and Cloud, handles them. It also presents proposed optimizations in the literature. Due to those different I/O requirements, there is no "one storage system for all" approach. This is more evident in large scale computing sites, where distributed storage solutions support multiple concurrent applications with conflicting requirements. We believe that future storage systems need a major re-design to efficiently support the diversity of workloads and the explosion of scale.

2.3 Motivating Examples

Applications perform I/O operations for several reasons: reading initial input, writing final output (i.e., results to persistent layer), temporary I/O for out-of-core computations, defensive I/O

Table 1. Application I/O Requirements

Feature	I/O requirement	HPC	Cloud	Optimizations
Data consistency	Data passed to the I/O system must be consistent between operations.	Strong, POSIX	Eventual, Immutable	Tunable consistency [73]
File access	Multiple processes must be able to operate on the same file concurrently	Shared Concurrent	Multiple replicas	Collective I/O [70], Concurrent file handlers [26], Complex locks [20, 78, 81]
Global namespace	Data identifiers must be resolved and recognizable in a global namespace that can be accessed from anywhere	Hierarchical Directory, Nesting	Flat	Namespace partitioning [77], Client-side caching [21], Decouple data-metadata [64, 83], File connectors [6, 7, 68]
Fault tolerance	Data must be protected against faults and errors	Specialized hardware, Check-pointing	Data replication, Data partitioning	Erasur coding [75]
Scale	Support for extreme scales and multi-tenancy	Few large jobs, Batch processing	Many small jobs, Iterative	Job scheduling, I/O buffering, Scale-out
Locality	Jobs are spawned where data is	Remote storage	Node local	Data aggregations
Ease of use	Interface, user-friendliness and ease of deployment	High-level I/O libraries	Simple data formats	Amazon S3, Openstack Swift

(i.e., checkpointing), and process-to-process data sharing (i.e., inter- and intra-node communications). We provide some examples of workloads that demonstrate the growing need of a storage system that supports diverse workloads on the same single platform.

2.3.1 CM1 (Write-intensive, Final Output). CM1 is a multi-dimensional, non-linear, numerical model designed for idealized studies of atmospheric phenomena [12]. CM1's I/O workload demonstrates a sequential write pattern. The simulation periodically writes collectively its results (e.g., atmospheric points with a set of features) using MPI-IO. Data are written in a binary GrADS format with a shared file access pattern. This workload requires persistence, fault-tolerance, and highly concurrent file access.

2.3.2 HACC (Update-intensive, Check-pointing). HACC stands for Hardware Accelerated Cosmology Code and is a cosmological simulation that studies the formation of structure in collisionless fluids under the influence of gravity in an expanding universe. Each process in HACC periodically saves the state of the simulation along with the dataset using POSIX and a file-per-process pattern. Since HACC runs in time steps, only the last step checkpoint data is needed. Thus, the I/O workload demonstrates an update-heavy pattern. A major performance improvement in HACC workflow is the addition of burst buffers that absorb the checkpointing data faster and perform the last flush of data to the remote PFS.

2.3.3 Montage (Mixed Read/Write, Data Sharing). Montage is a collection of programs comprising an astronomical image mosaic engine. Each phase of building the mosaic takes an input from

the previous phase and outputs intermediate data to the next one. It is an MPI-based engine and therefore Montage's workflow is highly dependent on the data migration between processes. The exchange of data between executables is performed by sharing temporary files in the Flexible Image Transport System (FITS) format via the storage system. At the end a final result is persisted as the final jpeg image. The I/O workload consists of both read and write operations using either POSIX or MPI independent I/O.

2.3.4 K-means Clustering (Read-intensive, Node-local). This application is a typical and widely used BigData kernel that iteratively groups datapoints into disjoint sets. The input datapoints can be numerical, nodes in a graph, or set of objects (e.g., images, tweets, etc.). Implementations using the MapReduce framework [18] remain the most popular clustering algorithm because of the simplicity and performance. The algorithm reads the input dataset in phases and each node computes a set of means, broadcasts them to all machines in the cluster and repeats until convergence. The I/O workload is read-intensive and is performed on data residing on the node locally. K-means clustering is typically I/O bound [59].

3 LABIOS

LABIOS, a new class of a storage system that uses data-labeling to address the issues discussed in Section 2, is a distributed, fully decoupled, and adaptive I/O platform that is intended to grow in the intersection of HPC and BigData. This section presents the design and implementation of LABIOS.

3.1 Design Requirements

As any distributed storage system, LABIOS is designed to be responsible for the organization, storage, retrieval, sharing, and protection of data. LABIOS also contains a representation of the data itself and methods for accessing it (e.g., read/write). However, LABIOS manipulates data in a new paradigm using data labeling and therefore it makes some assumptions and design choices that might be different or complementary to a traditional parallel or distributed storage system. LABIOS' main objective is to *support a wide variety of conflicting I/O workloads under a single platform*.

LABIOS is designed with the following principles in mind:

- (a) *Storage Malleability:* Applications' I/O behavior consists of a collection of I/O bursts. Not all I/O bursts are the same in terms of volume, intensity, and velocity. The storage system should be able to tune the I/O performance by dynamically allocating/deallocating storage resources across and within applications, a feature called data access concurrency control. Storage elasticity enables power-capped I/O, where storage resources can be suspended or shutdown to save energy. Much like modern operating systems shut down the hard drive when not in use, distributed storage solutions should suspend servers when there is no I/O activity.
- (b) *I/O Asynchronicity:* A fully decoupled architecture can offer the desired agility and move I/O operations from the existing streamlined paradigm to a data-labeling one. In data-intensive computing where I/O operations are expected to take a large amount of time, asynchronous I/O and the data-labeling paradigm are a good way to optimize processing efficiency and storage throughput / latency.
- (c) *Resource Heterogeneity:* The hardware composition of the underlying storage should be managed by a single I/O platform. In other words, heterogeneity in hardware (RAM, NVMe, SSD, HDD) but also the presence of multiple layers of storage (i.e., local file systems, shared burst buffers, or remote PFS) should be transparent to the end user. The

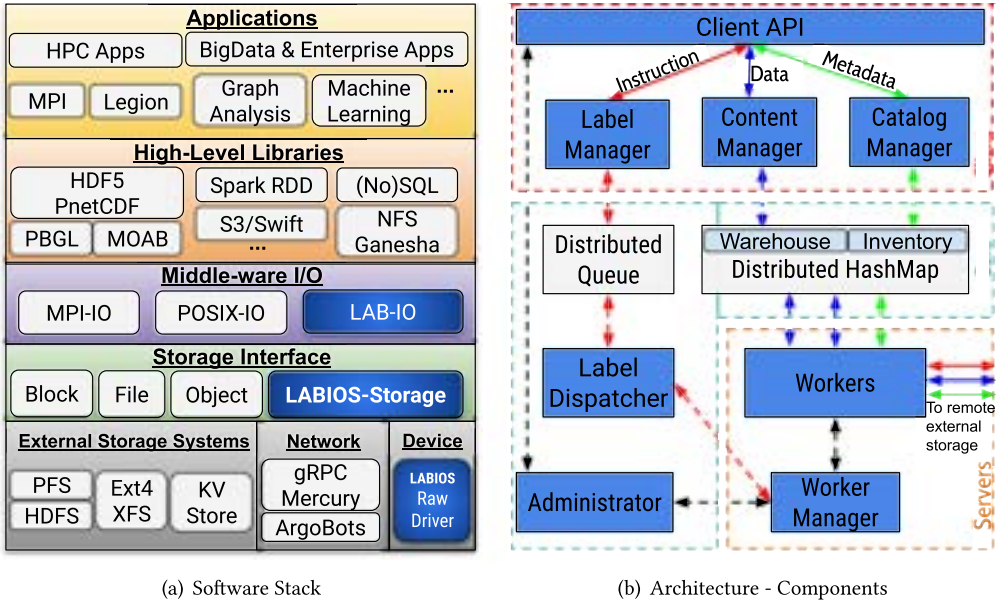
storage infrastructure should be able to dynamically reconfigure itself to meet the I/O demand of running applications and their I/O requirements. Moreover, storage Quality of Service (QoS) guarantees are a highly desired feature that can be achieved by efficiently matching the supply to the I/O demand [51].

- (d) *Data provisioning*: The I/O system should be programmable (i.e., policy-based provisioning and management). Storage must naturally carry out data-centric architectures (e.g., ActiveStorage [66], or ActiveFlash [72]), where data operations can be offloaded to the storage servers relieving the compute nodes of work such as performing data filtering, compression, visualization, deduplication, or calculating statistics (i.e., Software Defined Storage (SDS)). Offloading computation directly to storage and efficient process-to-process data sharing can significantly reduce expensive data movements and is the pinnacle of success for data-centric architectures [60].
- (e) *Storage Bridging*: The I/O system should abstract low-level storage interfaces and support multiple high-level APIs. Modern distributed computing makes use of a variety of storage interfaces ranging from POSIX files to REST objects. Moreover, existing datasets are stored in a universe of storage systems, such as Lustre, HDFS, or Hive. Storage solutions should offer developers the ability to use APIs interchangeably avoiding interface isolation and, thus, boost user productivity while minimizing programmability errors.

3.2 Architecture

3.2.1 Data Model. The core of LABIOS storage is a *Label*, which is effectively a tuple of one or more operations to perform and a pointer to its input data. It resembles a shipping label on top of a Post Office package where information such as source, destination, weight, priority, and so on, clearly describe the contents of the package and what should happen to it. LABIOS' label structure includes: type, uniqueID, source and destination as pointers (i.e., can be a memory pointer, a file path, a server IP, or a network port), operation to be performed as a function pointer (i.e., all functions, user- or pre-defined, are stored in a shared program repository, which servers have access to), and a collection of flags indicating the label's state (i.e., queued, scheduled, pending, cached, invalidated, etc.). In essence, labels encapsulate the instructions to be executed on a piece of data. All I/O operations (e.g., fread() or get(), fwrite() or put(), etc.) are expressed in the form of one or more labels and a scheduling policy to distribute them to the servers. Labels belong to each application exclusively. They are immutable, independent of one another, cannot be reused, and can be executed by any worker anywhere in the cluster. In contrast, labels are not a computation decomposition (i.e., compute task), or a simple data object encapsulation (e.g., RDDs) but rather a storage-independent abstraction that simply expresses the application's intent to operate on certain data. In their most primitive form, labels are of a simple type such as write or read. In other words, a write-label is one where the operation bundled with a piece of data is to read the data from the defined source and persist it at the defined destination. LABIOS organizes information with primitive labels that include: write, read, delete, and move. More complicated operations are expressed with complex labels including: software-defined-storage (SDS) (i.e., where the data operation defined applies one or more transformation on its input data), shared, protected, and priority labels. The LABIOS runtime may be capable to support more label types but only when considering the computation necessary to carry out the label. More investigation on this is required.

3.2.2 Overview. As it can be seen in Figure 1(a), LABIOS can be used either as a middleware I/O library or as a full stack storage solution. Applications can use the LABIOS library to perform I/O using labels and take advantage of the full potential of the system. Each label can carry a set of



(a) Software Stack

(b) Architecture - Components

Fig. 1. LABIOS overview.

functions to be performed by the storage server that executes it. For instance, an application can push write labels and instruct LABIOS to first deduplicate entries, sort the data, compress them, and finally write them to the disk. However, to maintain compatibility with existing systems, legacy applications can keep their I/O stack and issue typical I/O calls (e.g., `fwrite()`). This way, there is no need to modify any code in the user application but only link the LABIOS library (either statically by re-compiling or dynamically using `LD_PRELOAD`) and LABIOS will intercept those I/O calls, transform them into labels automatically, and forward them to the storage servers. In the prototype implementation, we provide UNIX file I/O (i.e., POSIX and `STDIO`) wrappers but more can be added to support other interfaces such as HDF5, netCDF, or MPI-IO. LABIOS can also access data via its own raw driver that handles data on the storage device in the form of labels. By adding more servers, the capacity and performance of them is aggregated in a single namespace. Furthermore, LABIOS can unify multiple namespaces by connecting to external storage systems, a feature that allows LABIOS to offer effective *storage bridging*.

LABIOS offers high speed data access to parallel applications by splitting the data, metadata, and instruction paths and decoupling storage servers from the application, as shown in Figure 1(b). This decoupling of clients and servers is a major architectural choice that enables several key features in LABIOS: the power of the *asynchronous I/O*, the effectiveness of *data provisioning*, and the proliferation of *heterogeneous storage* resources. An incoming application first registers with LABIOS, upon initialization, and passes workload-specific configurations to set up the environment. LABIOS receives the application’s I/O requests via the client API, transforms them, using the label manager, into one or more labels depending mostly on the request size, and then pushes them into a distributed label queue. Users’ data are passed to a distributed data warehouse and a metadata entry is created in an inventory. A *label dispatcher* consumes the label queue and distributes labels using several scheduling policies. Storage servers, called LABIOS *workers*, are organized into a worker pool with a manager being responsible to maintain its state. Workers can be suspended depending on the load of the queue creating an elastic storage system that is able to react to the state of the

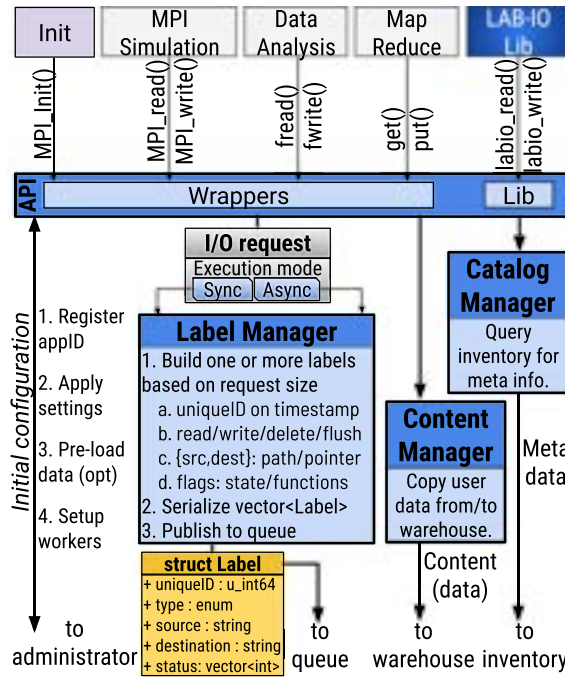


Fig. 2. LABIOS client internal design.

cluster. Last, workers execute their assigned labels independently and read/write data either on their own storage device or through a connection to an external storage system.

3.3 Component Analysis

LABIOS consists of three main components that work in harmony but in a decoupled fashion: the clients, the core lib, and the servers. We present here a detailed analysis of all LABIOS components.

3.3.1 LABIOS Client. This component, shown in Figure 2, interacts with the application and has three main goals: (a) per-application system initialization: register application info (i.e., ID, group name, group credentials and permissions), apply application-specific settings, pre-load data from external sources (if needed), and setup LABIOS workers; (b) accept application’s I/O requests, either by intercepting existing I/O calls using function call wrappers or by exposing LABIOS API, and; (c) build labels based on the incoming I/O request. It consists of the following:

- (a) The *Label Manager* is responsible to create labels based on the incoming I/O requests. The label creation is automatic and does not require any user intervention. The Label Manager builds one or more labels based on the request characteristics (e.g., read/write, size, file path, etc.), serializes and publishes them to the distributed label queue. Each label gets a unique identifier based on the origin of the operation and a timestamp (in nanoseconds), which ensures the order of operations (i.e., this is the constraint in the priority queue). There are two ways labels get created: a) via the native API, or b) transparently via the provided wrappers. In the first case, users create labels programmatically within a configurable maximum data size (much like how key-value pairs are created in any object store). In the case of wrappers, labels are created by a configurable size parameter within a range of min and max values (e.g., min 64 KB–max 4 MB). The data size parameter in each label

is the unit of data distribution in the system. An I/O request larger than the maximum label size will be split into more labels creating a 1-to-N relationship between request and number of labels (e.g., for a 10 MB `fwrite()` and 1 MB `max_label_size`, 10 labels will be created). Any I/O request smaller than the minimum label size will be cached and later aggregated in a special indexed label to create a N-to-1 relationship between number of requests and label (e.g., for ten 100 KB `fwrite()` and 1 MB `min_label_size`, one label will be created). Last, these thresholds can be bypassed for certain operations, mostly for synchronous reads. Setting min and max label size values is dependent on many system parameters such as memory page size, cache size, network type (e.g., TCP buffer size), and type of destination storage (e.g., HDDs, NVMe, SSDs). LABIOS can be configured in a *synchronous* mode, where the application waits for the completion of the label, and in *asynchronous* mode, where the application pushes labels to the system and goes back to computations. A waiting mechanism, much like a barrier, can be used to check the completion of a single or a collection of asynchronously issued labels. The async mode can significantly improve the system's throughput but it also increases the complexity of data consistency and fault tolerance guarantees.

- (b) The *Content Manager* is mainly responsible for handling user data inside a *warehouse*. The warehouse is implemented by a distributed hashmap (i.e., key-value store), it temporarily holds data in-memory effectively serving as the bridge between clients and workers. The warehouse is a collection of system-level structures (i.e., tables in the distributed key-value store), that are application-specific, and has the following requirements: highly available, concurrent data access, fault tolerant, and high throughput. The content manager exposes the warehouse via a simple `get/put/delete` interface to both the clients and the workers. The size and location of the warehouse is configurable based on several parameters such as number of running applications, application's job size, dataset aggregate size, and number of nodes (e.g., one hashtable per node, or per application). Every entry in the warehouse is uniquely identified by a key that is associated with one or more labels. The content manager can also create ephemeral regions of the warehouse (e.g., temporary rooms), which can be used for workflows where data are shared between processes (i.e., data sharing, Section 2.2). Data flows through LABIOS as follows: from application's buffer to the warehouse, and from there to worker storage for persistence or to another application's buffer. Last, the content manager also provides a cache to optimize small size data access; a known issue in distributed storage [13]. I/O requests smaller than a given threshold are kept in a cache and, once aggregated, a special label is created and pushed to the distributed queue to be scheduled to a worker (much like memtables and SSTables in LevelDB). This minimizes network traffic and can boost the performance of the system.
- (c) The *Catalog Manager* is responsible to maintain both user and system metadata information in an inventory, implemented by a distributed hashmap. The catalog manager exposes an interface for each application to query and update the entries within the inventory. Decentralization of the catalog services makes the system scalable and robust. Multiple concurrent processes can query the inventory at the same time. For concurrent updates, LABIOS adopts the semantics of the underlying distributed hashmap with high-availability and concurrent access ensuring the correctness and high throughput of catalog operations. LABIOS also offers the flexibility to place the inventory in memory for high performance, protected by triple replication for fault tolerance. However, this increases the memory footprint of LABIOS and it depends on the availability of resources. The organization of inventory entries depends on the data model (files, objects, etc.) and/or high-level I/O libraries and middleware. For instance, for POSIX files the inventory entries

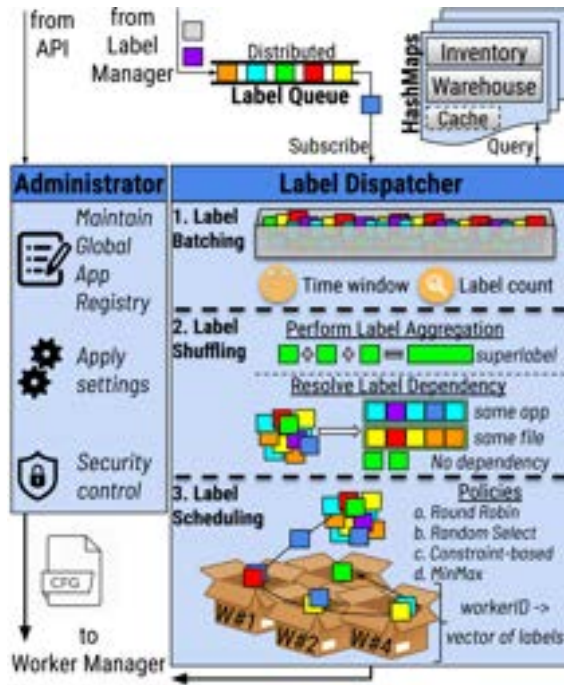


Fig. 3. LABIOS core internal design.

may include: filename to file stat, file handler to filename, file handler to file position in offset, filename to a collection of labels, and others. An HDF5 or a JSON file will have different inventory entries. The distribution of the user’s namespace between machines follows the state-of-the-art metadata services proven within the storage community [64]. LABIOS-specific catalog information include: label status (e.g., in-transit, scheduled, pending), label distribution (e.g., label to workerID), label attributes (e.g., ownership, flags), and location mappings between user’s data and LABIOS internal data structures (e.g., a user’s POSIX file might be stored internally as a collection of objects residing in several workers). Last, when LABIOS is connected to external storage resources, it relies on their metadata service. LABIOS becomes a client to the external storage resources and “pings” their metadata service to acquire needed information. LABIOS does not keep a copy of their respective metadata internally to avoid possible inconsistent states. However, further investigation is needed to optimize this process by avoiding added network latencies from external sources.

3.3.2 *LABIOS Core*. This component, shown in Figure 3, is responsible to manage the instruction, data, and metadata flow separately. It consists of the following:

- (a) The *Administrator* maintains the system’s state by keeping track of all running applications in a global registry, setting up the environment per application (e.g., boot up exclusive workers if needed, pre-load data from external sources, etc.), and performing security control via user authentication and access permission checks.
- (b) The *Label Queue* hosts the labels from the application. LABIOS distributed queuing system has the following requirements: high message throughput, always on and available, at-most-once delivery guarantees, highly concurrent, and fault tolerant. These features

ensure data consistency, since the label dispatcher will consume labels once and in order. The queue concurrency ensures that multiple dispatchers can service the same queue or one dispatcher multiple queues. The number of queues is configurable based on the load (e.g., one queue per application, or one queue per 128 processes, or one queue per node).

- (c) The *Label Dispatcher* subscribes to one or more distributed label queues and dispatches labels to workers using several scheduling policies. The label dispatcher is multi-threaded and can run on one or more nodes depending on the size of the cluster. LABIOS dispatches labels based on either a time window or the number of labels in the queue; both of those parameters are configurable. For example, the dispatcher can be configured to distribute labels one by one or in batches (e.g., every 1000 labels). To avoid stagnation, a timer is also used; if the timer expires, then LABIOS will dispatch all available labels in the queue. Furthermore, the number of label dispatchers is dynamic and depends on the number of deployed queues. There is a fine balance between the volume and velocity of label production stemming from the applications and the rate at which the dispatcher handles them. The relationship between the dispatcher and queuing system increases the flexibility and scalability of the platform and provides an infrastructure to match the rate of incoming I/O. The dispatcher consists of two phases:

- (a) *Label Shuffling*: takes a vector of labels as an input and shuffles them based on type and flags. Two operations are performed by the shuffler. *Data aggregation*: labels that reflect user's requests in consecutive offsets can be combined to one larger label to maintain locality (this feature can be turned on or off). *Label dependencies*: data consistency must be preserved for dependent labels. For instance, a read after write pattern; LABIOS will not schedule a read label before the dependent write label completes. To resolve such dependencies, the shuffler will create a special label, called *supertask*, which embodies a collection of labels that need to be executed in strictly increasing order. After sorting the labels and resolving dependencies, the shuffler sends labels either to the solver to get a scheduling scheme, or directly to the assigner depending on the type (e.g., a read label is preferably assigned to the worker that holds the data to minimize worker-to-worker communication).
- (b) *Label Scheduling*: takes a vector of labels as an input and produces a dispatching plan. For a given set of labels and workers, the scheduler answers three main questions: how many workers are needed, which specific workers, and which labels are assigned to which workers. LABIOS is equipped with several scheduling policies (in detail in Section 3.5). A map of {workerID, vector of labels} is passed to the worker manager to complete the assignment by publishing the labels to each individual worker queue. Labels are published in parallel using a thread pool. The number of threads in the pool depends on the machine the label dispatcher is running on as well as the total number of available workers.

3.3.3 *LABIOS Server*. This component, shown in Figure 4, is responsible for managing the storage servers and has two main entities:

- (a) The *Worker* is essentially the storage server in LABIOS. It is fully decoupled from the applications, is multithreaded, and runs independently. Its main responsibilities are: (a) service its own queue, (b) execute labels, (c) calculate its own worker score and communicate it to the worker manager, (d) auto-suspend itself if there are no labels in its queue for a given time threshold, and (e) connect to external storage sources. The worker score is a new metric, critical to LABIOS operations, that encapsulates several characteristics of the worker into one value, which can then be used by the label dispatcher to assign

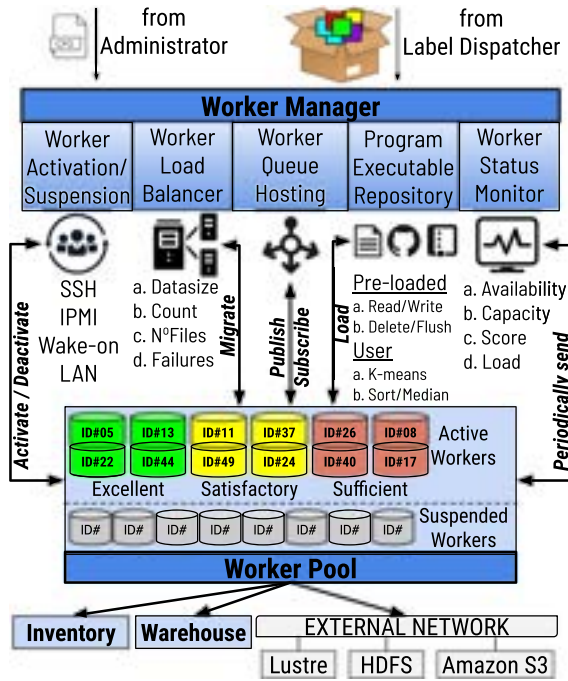


Fig. 4. LABIOS server internal design.

Table 2. LABIOS Worker’s Score—Weighting Examples

Priority	Availability	Capacity	Load	Speed	Energy
Low latency	0.5	0	0.35	0.15	0
Energy savings	0	0.15	0.2	0.15	0.5
High Bandwidth	0	0.15	0.15	0.70	0

any label to any appropriate worker. A higher scored worker is expected to complete the label faster and more efficiently. The score is calculated by every worker independently at an interval or if substantial change of status occurs, and it includes: (i) *availability*: 0, not available (i.e., suspended or busy); 1, available (i.e., active and ready to accept labels). (ii) *capacity*: (double) [0,1] based on the ratio between remaining and total capacity. (iii) *load*: (double) [0,1] based on the ratio between worker’s current queue size and max queue size (the max value is configurable). (iv) *speed*: (integer) [1,5] based on maximum bandwidth of worker’s storage medium and grouped based on ranges (e.g., 1: ≤ 200 MB/s, 2: 200–550 MB/s, ... 5: $\geq 3,500$ MB/s). (v) *energy*: (integer) [1,5] based on worker’s power wattage on full load (e.g., an ARM-based server with flash storage consumes less energy than a Xeon-based server with a spinning HDD). The first three are dynamically changing based on the state of the system whereas speed and energy variables are set during initialization and remain static. Last, each variable is multiplied by a weight. LABIOS’ weighting system is set in place to express the scheduling policy prioritized (examples shown in Table 2). For instance, if energy consumption is the constraint that the label dispatcher aims to optimize, then the energy variable gets a higher weight. The final score is a float in range between 0 and 1 and is calculated as follows: $Score_{worker(i)} = \sum_{j=1}^5 Weight_j * Variable_j$.

- (b) The **Worker Manager** is responsible for managing the worker pool. Its responsibilities are: (a) maintain workers' statuses (e.g., remaining capacity, load, state, and score) in a distributed hashmap (in-memory or on disk), (b) host the workers' queues, (c) perform load balancing between workers, and (d) dynamically commission/decommission workers to the pool. It is connected to the administrator for accepting initial configurations for incoming applications, and to the label dispatcher for publishing labels in each worker's queue. It can be executed independently on its own node by static assignment, or dynamically on one of the worker nodes by election among workers. In a sense, the worker manager partially implements objectives similar to other cluster resource management tools such as Zookeeper, or Google's Borg. One of the most performance-critical goals of the worker manager is to maintain a sorted list of workers based on their score. Workers update their scores constantly, independently, and in a non-deterministic fashion. Therefore, the challenge is to be able to quickly sort the updated scores without decreasing the responsiveness of the worker manager. LABIOS addresses this issue by a custom sorting solution based on buckets. The set of workers are divided on a number of buckets (e.g., high, medium, and low scored workers) and an approximate bin sorting algorithm is applied [29]. A worker score update will only affect a small number of buckets and the complexity time is relevant to the size of the bucket. Last, the worker manager can send activation messages to suspended workers either by using the administrative network, if it exists, (i.e., `ipmitool --power on`), or by a custom solution based on ssh connections and wake-on-lan tools.

3.4 Deployment Models

The power and potential of LABIOS' flexible and decoupled architecture can be seen in the several ways the system can be deployed. Depending on the targeted hardware and the availability of storage resources, LABIOS can: (a) replace an existing parallel or distributed storage solution, or (b) be deployed in conjunction with one or more underlying storage resources as an I/O accelerator (e.g., burst buffer software, I/O forwarding, or software-defined storage in user space). In any of these models, LABIOS will first spawn five separate services (i.e., LABIOS executables) during the application initialization (e.g., `MPI_Init()`). These services are: the warehouse, the queues, the label dispatcher, the worker manager, and the workers. The application, then, can simply use the LABIOS library (i.e., include `labios.hpp`) and start using the system. Leveraging the latest trends in hardware innovation, the machine model we present here as our basis for several deployment schemes is as follows: compute nodes equipped with a large amount of RAM and local NVMe devices, an I/O forwarding layer [37], a shared burst buffer installation based on SSD equipped nodes, and a remote PFS installation based on HDDs (motivated by the recent machines Summit in ORNL or Cori on LBNL). Four equally appropriate deployment examples that can cover different workloads can be seen below:

(a) *LABIOS as I/O accelerator* (Figure 5): Client runs on compute nodes and the distributed queue and hashmaps are placed on each node's memory for lower latency and higher throughput. The label scheduler runs on a separate compute node serving one or more queues per node, and last, one core per node can execute LABIOS worker, who stores data in the local NVMe. This mode can be used as a fast distributed cache for temporary I/O or on top of other external sources. It is also ideal for Hadoop workloads with node-local I/O. However, it must use some compute cores to run its services and I/O traffic will mix with the compute network. **Pros:** fast distributed cache for temporary I/O or on top of external storage resources, good performance for Hadoop workloads. **Cons:** I/O traffic mixed with compute network, overheads by using compute core to run LABIOS services.

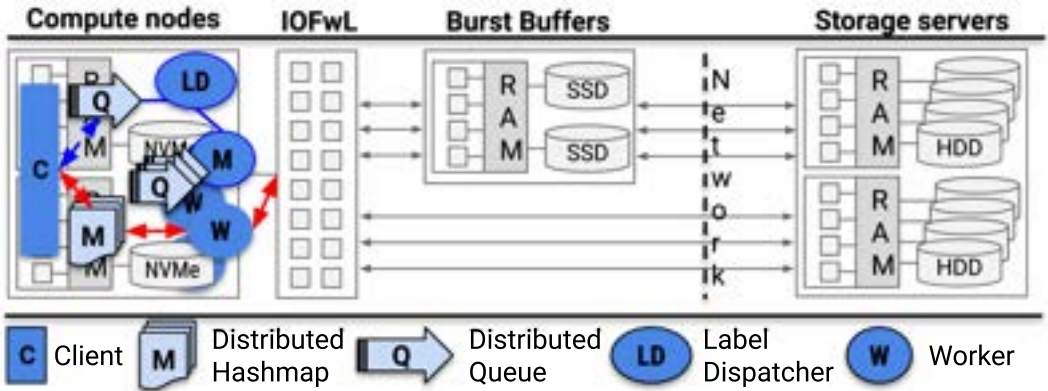


Fig. 5. Deployment example: LABIOS as I/O accelerator in compute nodes.

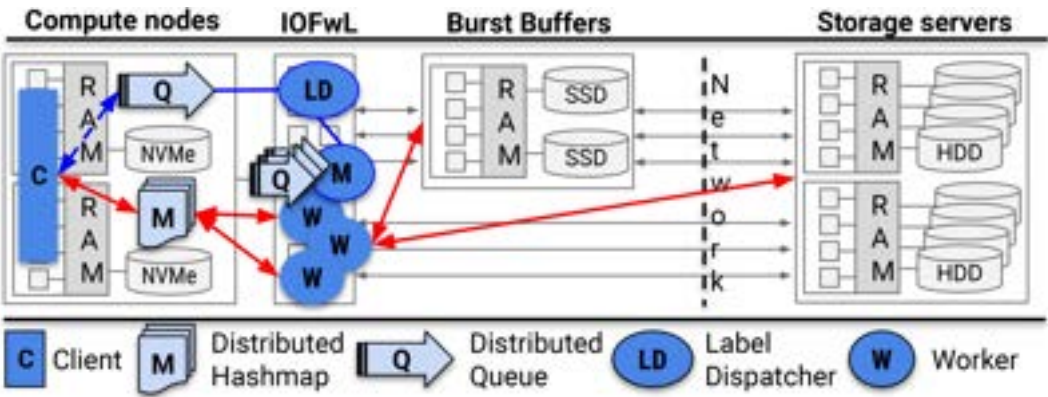


Fig. 6. Deployment example: LABIOS as I/O forwarder solution.

(b) *LABIOS as I/O forwarder* (Figure 6): Client runs on compute nodes and the distributed queue and hashmaps are placed on compute nodes’ memory or NVMe. The label scheduler and workers run on the I/O nodes of the forwarding layer. This mode is ideal for asynchronous I/O calls where applications pass their data to LABIOS, which pushes them in a non-blocking fashion to remote storage, either native to the system or external. However, its scalability is limited by the size of the I/O forwarding layer. **Pros:** good performance for asynchronous non-blocking I/O, great for storage bridging by connecting to many external storage services **Cons:** subject to the existence of an I/O forwarding layer, limited scalability.

(c) *LABIOS as I/O buffering* (Figure 7): Client runs on compute nodes and the distributed queue and hashmaps are placed on compute nodes’ memory or NVMe. The label scheduler can be deployed either in compute or I/O forwarder nodes, serving one or more client queues. Workers are deployed on the burst buffer nodes utilizing the SSD devices to store data. ideal for fast temporary storage, data sharing between applications, and *in situ* visualization and analysis. Requires additional storage and network resources (i.e., burst buffer infrastructure). **Pros:** fast scratch storage space, fast data sharing between applications, enables *in situ* visualization and analysis. **Cons:** requires additional storage and network hardware resources (e.g., burst buffers).

(d) *LABIOS as remote distributed storage* (Figure 8): This can be achieved with various combinations of deploying LABIOS components in different nodes in the cluster. For instance, a natural fit in

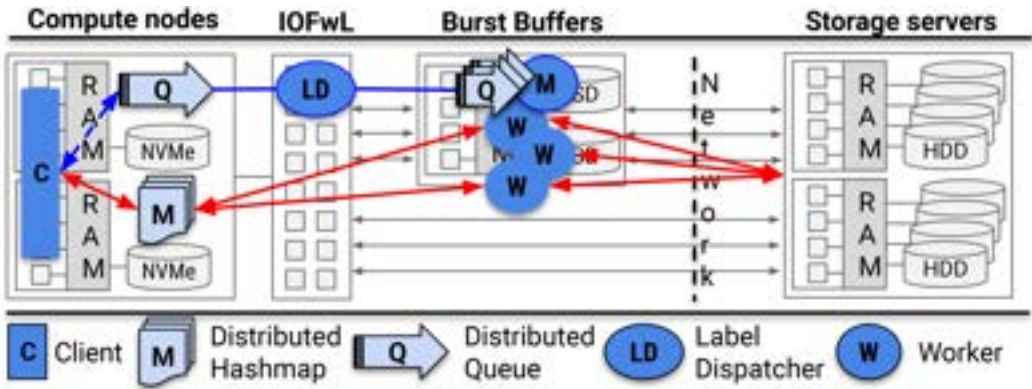


Fig. 7. Deployment example: LABIOS as I/O buffering platform.

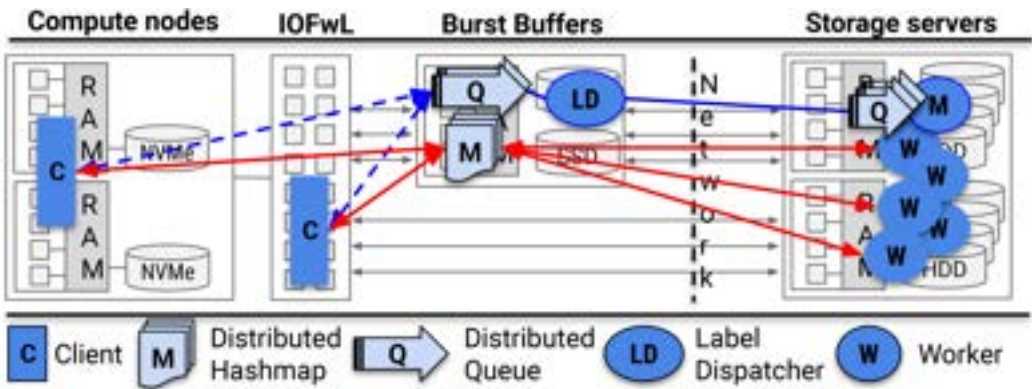


Fig. 8. Deployment example: LABIOS as remote distributed storage.

our machine model is running the client in the I/O forwarding nodes, the distributed queue and hashmaps on the burst buffers, and the workers on the storage servers, effectively replacing a PFS. This mode offers better system scalability by scaling each individual component independently, better resource utilization, and higher flexibility to the system administrator. For instance, one can increase the number of client queues in scenarios when label production is high, or deploy more dispatchers to distribute labels faster. It has, however, higher deployment complexity. LABIOS’ fully decoupled architecture provides greater flexibility and promotes scalability; I/O scales along with the application by simply provisioning additional resources. **Pros:** very scalable, optimal resource utilization, increased flexibility. **Cons:** high deployment complexity, requires system admin involvement.

3.5 Implementation

3.5.1 Usage. LABIOS system is packaged into five separate services that get deployed during the application initialization. When LABIOS is compiled the following executables are created: (a) warehouse, (b) distributed queue(s), (c) label dispatcher(s), (d) worker manager(s), and (e) workers. These executables must be started before the client application can use the system. Depending on the deployment model, these services can start either manually via a collection of scripts or automatically during application’s initialization phase via a bootstrapping method

(e.g., intercept `MPI_Init()` and spawn LABIOS). Clients simply connect to the warehouse and the distributed queues via the LABIOS library. The only user intervention is to include the LABIOS header file. We plan to incorporate the bootstrapping of LABIOS as a plugin to the global system scheduler (e.g., Flux [1]).

3.5.2 Label Scheduling. LABIOS balances the rate of incoming labels, the dispatching cost, and the time to execute the labels and offers an flexible, intent-aware infrastructure. LABIOS provides a custom data distribution by implementing different scheduling policies:

- (1) *Round Robin*: given a set of labels and a list of available workers the dispatcher will distribute labels in a round robin fashion, much like a PFS does. The responsibility of activating workers and compiling a list of available workers for every scheduling window falls under worker manager. This policy demonstrates low scheduling cost but additional load balancing between workers might occur.
- (2) *Random Select*: given a set of labels, the dispatcher will distribute labels to all workers randomly regardless of their state (i.e., active or suspended). This policy ensures the uniform distribution of workload between workers, low scheduling cost, but with no performance guarantees (i.e., possible latency penalty by activating suspended workers, or lack of remaining capacity of worker, etc.).
- (3) *Constraint-based*: in this policy, LABIOS provides the flexibility to express certain priorities on the system. Through the weighting system of worker scores, discussed in Section 3.3, the dispatcher will distribute labels to workers based on the constraint with higher weight value. The constraints used are: (a) *availability*, active workers will have higher score. (b) *worker load*, based on worker's queue size. (c) *worker capacity*, based on worker's remaining capacity. (d) *performance*, workers with higher bandwidth and lower latency get higher score. For a given set of labels, the dispatcher requests a number of workers with the highest score, respective to the prioritized constraint, from the worker manager and distributes the labels evenly among them. The number of workers needed per a set of labels is automatically determined by LABIOS based on the total aggregate I/O size and the selected constraint balancing parallel performance and efficiency. These heuristics can be configured and further optimized based on the workload.
- (4) *MinMax*: given a set of labels and a collection of workers, the dispatcher aims to find a label assignment that maximizes I/O performance while minimizing the system's energy consumption, subject to the remaining capacity and load of the workers; essentially a minmax multidimensional knapsack problem, a well-known NP-hard combinatorial optimization problem [61]. LABIOS solves this problem using an approximate dynamic programming (DP) algorithm [5], which optimizes all constraints from the previous policy. This policy gives a near-optimal matching of labels-workers but with a higher scheduling cost.

3.5.3 Software Defined Storage (SDS) Support. Computational storage, often expressed with terms such as *Active Storage* or *Software Defined Storage*, is the ability of a storage system to carry out some form of computations or data operations right where the data resides. Offloading data-intensive functions to the storage can be beneficial in two meaningful ways: (a) reducing data movement from storage to compute thus reducing network traffic, energy consumption, and ultimately freeing compute performance, and (b) providing functional asynchronicity where applications can exit once they pass the data to the storage (e.g., archival storage with data compression performed on storage). The decoupled architecture of LABIOS along with the label paradigm are a great fit for offering a storage infrastructure for computational storage operations. Further, the decoupled architecture of LABIOS workers allows us to offload computations to both the CPU of

Client		Core		Server		Total LABIOS LOC
Component	LOC	Component	LOC	Component	LOC	
API (native&wrapper)	2588	Adminstrator	669	Worker	1599	
Catalog Manager	340	Label Dispatcher	1530	Worker Manager	2204	
Content Manager	286					
Label Manager	462					
Total Client	3676	Total Core	2199	Total Server	3803	
						9678

Fig. 9. LABIOS prototype implementation lines of code.

the storage node or even to any computational elements of modern storage devices (e.g., ARM processors or FPGAs inside a Flash SSD controller). We have already explored the implications of using the computational power of such storage devices in our previous work [23, 44].

The infrastructure LABIOS offers to enable computational storage capabilities to applications is simple and scalable. First, users submit the source code of their data intensive functions. LABIOS provides a custom compiler (i.e., wrapper over GCC) to build the provided source code into shared libraries. Then, the compiled shared libraries, which reside in a special shared folder in the Worker Manager’s disk (see Figure 4) get dynamically loaded into each worker node’s OS using LD_PRELOAD environment variable. Once this happens, each worker has the definitions required to execute the provided operation. From there, things are very simple for the end user. On label creation, user needs to do the following: define the type of labels as SDS_IN_SITU and include the definition of their data-intensive function (i.e., function name and args) passed as a function pointer. Once a worker picks up a label that includes an instruction like this, it simply performs the operation. More interestingly, those data intensive functions can be stacked. For instance, a user might want to first sort the data, then compress them, and, finally, write them out to storage. This dictates a certain order of operations. Users can define such operation dependencies. In case such dependency is not defined, LABIOS will execute the instructions in the order they were declared. Users need to be aware of such system behavior. A non-resolved dependency can lead to sub-optimal results. For instance, in our previous example of sorting, compressing, and writing, it is widely known that compression is more effective on sorted data (i.e., smaller compressed size).

However, offloading data operations on storage nodes is not a new concept. There must be balance between how much computation one offloads to the often weaker computation elements of storage nodes and the cost of moving data to the more powerful CPUs of the compute nodes. LABIOS provides an efficient infrastructure to perform such offloading but does not currently have a balance check. It is the responsibility of the user to evaluate how much computation should be offloaded to offset the cost of data movement. There are a few insights that can be used as a guide in the literature, most notably Active Disks [65], Active Storage [66], and Active Flash [72].

3.5.4 Prototype Library Implementation Details. LABIOS is written in C++ and has approximately 10K lines of code, excluding external open source projects we used in our LABIOS prototype. Figure 9 shows the summary of LOCs for all LABIOS components in our prototype implementation. For the distributed queuing system, used widely in LABIOS both in client and workers, we used NATS server [17], a simple, high-performance open source messaging system. NATS was selected due to its superiority in throughput (i.e., >200K msg/s, 12× higher than Apache ActiveMQ and 3× than Kafka), low latency, and lightweight nature. For the distributed hashmaps, we used a custom version of Memcached with the *extstore* [54] plugin. Memcached is a simple in-memory key-value store, with easy deployment, development, and great APIs. The *extstore* plugin allows us to store memcached data to a storage device (e.g., NVMe, SSD) instead of RAM. We also modified the default key distribution from randomly hashing keys to servers, to a node local scheme to

```

1  #include <tabios.hpp>
2  ...
3  Client client = InitClient(ip, port, connConfig);
4  std::string path = "pvfs2:/data/integers.dat";
5  LabelSrc src = new LabelSrc(path, src_offset, size);
6  LabelType type = SDS_IN_SITU; //complex type
7  LabelFlag flags = CACHED | MPI_IO; //keep in cache
8  std::function<int(vector<int>)> fn = FindMedian;
9  Label label = client.CreateLabel(type, src, fn, flags);
10 Status status = client.IPublishLabel(label);
11 ... //perform other computations
12 client.WaitLabel(&status);
13 int median = std::static_cast<int>(status.data);

```

Fig. 10. LABIOS API example.

increase throughput and promote locality. Effectively, each node in LABIOS stores its hashtables on its local memcached daemon.

For label serialization, we used Cereal [31], a header-only binary serialization library, which is designed to be fast, light-weight, and easy to extend. For hashing, we used CitiHash algorithms [30], for memory allocations TCMalloc, and last, for metadata indexing an in-memory B-Tree implementation by Google.

3.5.5 API. LABIOS exposes a label API to the application to interact with data. The storage interface expresses I/O operations in the form of labels. The API includes calls to create-delete, publish-subscribe labels, among others. LABIOS' API offers higher flexibility and enables software defined storage capabilities. The active storage operation offloading is achieved via the label type (e.g., SDS) and a set of flags. For instance, the code snippet shown in Figure 10, creates an asynchronous label, which reads a file that includes a collection of integers from an external PFS using the MPI-IO driver, calculates the median value, and passes only the result back to the application via asynchronous I/O.

3.6 Discussion

3.6.1 Considerations. LABIOS' design and architecture promotes its main objective of supporting a diverse variety of conflicting I/O workloads under a single platform. However, additional features could be derived from LABIOS label paradigm. These include the following:

- (a) *Fault tolerance.* In the traditional streamlined I/O paradigm, if an `fwrite()` call fails the entire application fails and it must restart to recover (i.e., using check-pointing mechanisms developed especially in the scientific community). LABIOS' label granularity and decoupled architecture could provide the ability to repeat a failed label and allow the application to continue without restarting.
- (b) *Energy-awareness.* First, LABIOS' ability to dynamically commission/decommission workers to the pool creates an elastic storage solution with tunable performance and concurrency control but also offers a platform that could leverage the energy budget available. One could observe the distinct compute-I/O cycles and redirect energy from compute nodes to activate more LABIOS workers for an incoming I/O burst. Second,

LABIOS' support of heterogeneous workers could lead to energy-aware scheduling where non mission-critical work would be distributed on low-powered storage nodes, effectively trading performance for power consumption.

- (c) *Storage containerization.* Virtualization could be a great fit for LABIOS' decoupled architecture. Workers could execute multiple containers running different storage services. For instance, workers would host one set of containers running Lustre servers and another running MongoDB. The worker manager would act as the container orchestrator and the label dispatcher could manage hybrid workloads by scheduling labels to both services under the same runtime.

3.6.2 Challenges and Limitations. During the design and implementation of LABIOS, we have identified several challenges as follows:

- (a) *Multitenancy and Security.* How does LABIOS handle multiple applications accessing the system? Contention avoiding techniques must be employed. LABIOS could handle this by operating in isolation, where each application would have had its own queues, label dispatchers, and workers; an additional layer of global application orchestrator [42] could deal with the overall system traffic. Further concerns might arise by multitenancy such as user authentications and security. For instance, the LABIOS prototype uses NATS server as the queuing system. As of now, NATS does not support TLS and SSL, which could limit the secure capabilities of the system.
- (b) *Concurrent operations:* There are many components in LABIOS that need to be efficiently accessed at the same time. The distributed queuing system needs to support highly concurrent label insertion. The distributed hashmaps have to support a large number of clients for accessing both data and metadata. For instance, LABIOS' constraint on the priority queue is the timestamp, which means clock skewness across the cluster can affect the order of operations in the system. LABIOS' design relies heavily on the characteristics and performance of those components and any limitations from their side could become limitations of the entire LABIOS system.
- (c) *I/O request decomposition:* The main question LABIOS faces is how to transform I/O requests into labels. What would be an optimal decomposition granularity? LABIOS' label manager addresses this by splitting requests based on an I/O size range. For small requests, LABIOS caches them and aggregate them into a larger label. For large requests, LABIOS splits them into more labels offering a higher degree of parallelism. We plan to leverage the underlying hardware characteristics (network buffers, RAM page size, disk blocks, etc.) to refine LABIOS I/O request decomposition strategy. Another question, relevant to label granularity, is how are label dependencies and session/service management handled? Any task-based system faces these challenges [74]. LABIOS resolves label dependencies based on a configurable policy-driven granularity (i.e., per-application, per-file, per-dataset, etc.). We plan to further LABIOS ability to resolve label dependencies by using dependency graphs.
- (d) *User-defined functions execution safety:* LABIOS allows users to submit their own custom data functions to be used in labels under the software defined storage type (see subsection 3.5.3). We provide a custom compiler, which is a simple wrapper over GCC, that compiles user-submitted code and stores the binary in an executable registry. However, ill-written user code might misbehave and cause adverse effects on LABIOS stability and performance. There are no safety controls in our prototype implementation of LABIOS, but we plan to add some in the future. Our design allows such future extensions where the compiler can reject user code due to safety concerns. LABIOS can, however, currently

protect the integrity of user's data by shielding users from each other. The executable registry is only exposed to the owner of the application and other applications cannot access each others binaries. In other words, user authentication and permission checks, handled by the administrator during registration, provide a level of protection in regards to the integrity of data across multiple users.

4 EVALUATION

4.1 Methodology

4.1.1 Testbed. Experiments were conducted on a bare metal configuration offered by Chameleon systems [14]. The total experimental cluster consists of 64 client nodes, 8 burst buffer nodes, and 32 storage servers. Each node has a dual Intel(R) Xeon(R) CPU E5-2670 v3 @ 2.30 GHz (i.e., a total of 48 cores per node), 128 GB RAM, 10 Gbit Ethernet, and a local Seagate ST2000NX0273 HDD for the OS. Each burst buffer node has the same internal components but, instead of an HDD, it is equipped with Intel SSDSC1BG40 SSDs. LABIOS has been deployed in the cluster as a storage solution (Figure 8). Additional experimental evaluations (tests in subsection 4.2.6 (b) and (c)) have been conducted on the Ares cluster at Illinois Institute of Technology. Each compute node has a dual Intel(R) Xeon Scalable Silver 4114 @ 2.20 GHz (i.e., 40 cores per node), 96 GB RAM, 40 Gbit Ethernet, and a local 512 GB M.2 Samsung 960 Evo NVMe SSD. Each storage node has a dual AMD Opteron 2384 @ 2.7 Ghz, 32 GB RAM, 40 Gbit Ethernet with RoCE, and is also equipped with 2x512 GB Samsung 860 Evo SATA SSD in RAID, and 2TB Seagate 7,200K SATA HDD. The total experimental cluster consists of 2,560 client MPI ranks (i.e., 64 nodes), 4 burst buffer nodes, and 24 storage nodes.

4.1.2 Software. The cluster OS is CentOS 7.2, the PFS we used is OrangeFS 2.9.6. We use MPICH 3.2.1 as the MPI library for the applications. In terms of workloads, we used all four applications from Section 2.3, as well as several application kernels such as word count, integer sorting, numerical statistics, data filtering, and data compression. Finally, several synthetic micro-benchmarks were used to evaluate the internal components of LABIOS. More information on how each application or kernel works are provided with the description on each test and figures presented below.

4.2 Experimental Results

4.2.1 Anatomy of LABIOS Read/Write Operations. Figure 11 shows decomposition of the read and write label execution expressed as time percentage and divided by each LABIOS component. For instance, a *write label* starts with the LABIOS client building a label (at 12 o'clock on the figure), which takes 2% of the total time, it then passes the data to the warehouse (put data 11%), publishes the label to the queue (1%), and finally updates the catalog manager (MDM) about the operation (17%). The total LABIOS client operations take 31% of the total time. The label journey continues in the label dispatcher who picks up the label from the queue (subscribe 5%), schedules it (3%), and pushes it to a specific worker's queue (publish 1%). The most work is done by the LABIOS worker (60% of the total operation time) who first picks up the label from its queue and the data from the warehouse (get data 17%), writes the data down to the disk (29%), and finally updates the catalog manager (1%). Read label time decomposition can also be seen in Figure 11. All results are the average time of executing a 1 MB label 10K times.

4.2.2 Label Dispatching. In this test, we present how LABIOS performs with different scheduling policies and by scaling the number of label dispatcher processes. We report the rate (i.e., labels per second) at which each scheduling policy handles incoming labels. LABIOS client runs on all 64 client machines, the label dispatcher is deployed on its own dedicated node, and LABIOS

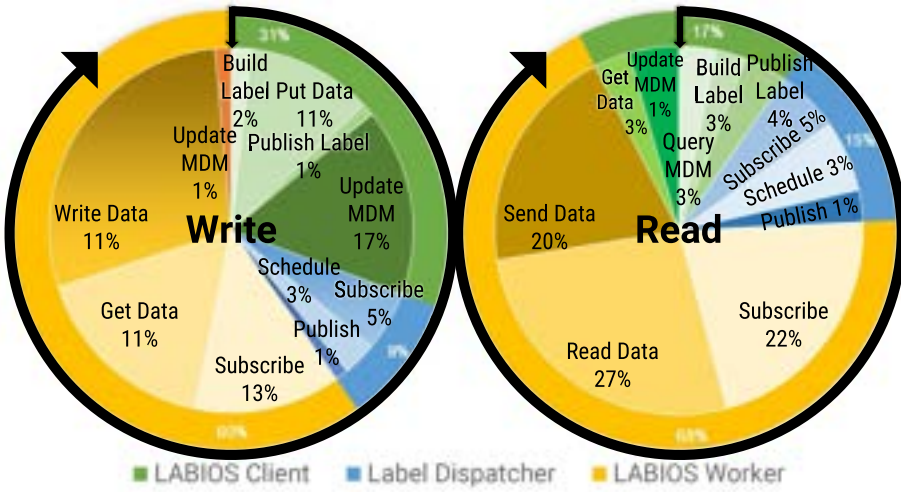


Fig. 11. Anatomy of LABIOS operations.

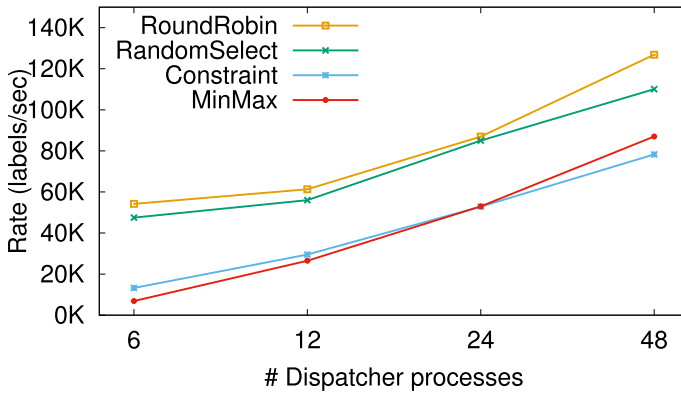


Fig. 12. Label Scheduling.

workers run on the 32 server machines. We measure the time the dispatcher takes to distribute 100K randomly generated labels (i.e., mixed read and write equally sized labels). As it can be seen in Figure 12, all policies scale linearly as we scale the label dispatcher processes from 6 to 48 (i.e., equal to max cores of the node). Round-robin and random-select achieve comparable scheduling rates between 55–125K labels per second. Constraint-based is more communication intensive, since it requires exchanging information about the workers with their manager. MinMax scales better with more resources, since it is more CPU intensive (i.e., DP approach).

4.2.3 Storage Malleability. In this test, we present how LABIOS elastic storage feature affects I/O performance and energy consumption. We issue 4,096 write labels of 1 MB each, and we measure the total I/O time stemming from different ratios between active workers over total workers (e.g., 50% ratio means that 16 workers are active and 16 are suspended). A suspended worker can be activated in about 3 seconds on average (in our testbed between 2.2 and 4.8 s). Figure 13 demonstrates the importance of balancing the added latency to activate more workers and the additional performance we get. We show two worker allocation techniques, the static (S), where labels are

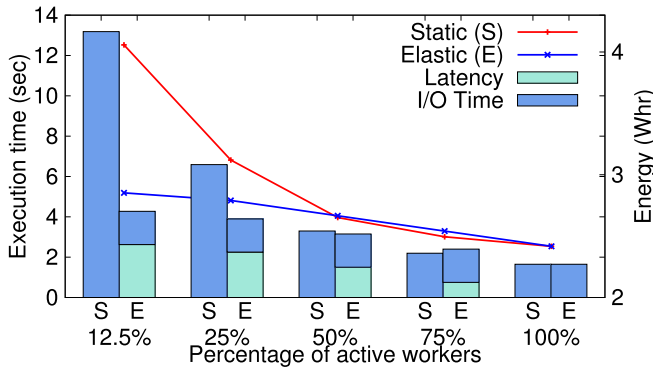


Fig. 13. Storage Malleability.

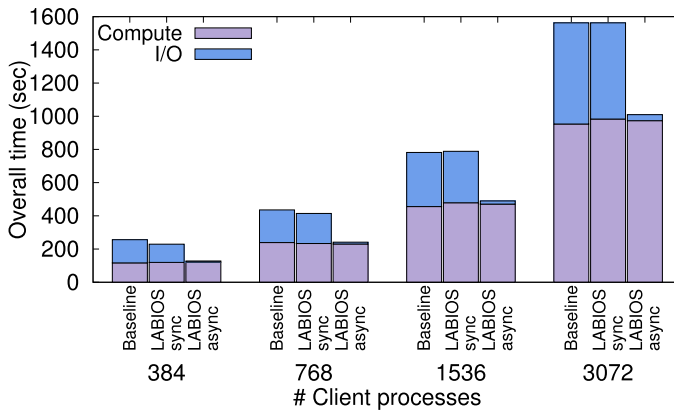


Fig. 14. I/O asynchronicity—CM1 performance.

placed only on the active workers, and the elastic (E), where more workers activate to serve incoming I/O. When LABIOS has a small percentage of active workers, the elastic strategy can boost performance significantly even though we pay the latency penalty to activate more workers. However, when we have a sufficient number of active workers (e.g., 75% or 24 out of 32 total workers), waking up more workers hurts the performance due to the latency penalty. This is further apparent when we see the energy efficiency of the system, expressed in watts per hour (Whr). In our testbed, active workers consume 165 watts, whereas suspended workers only 16 watts. LABIOS elastic worker allocation makes sense until the 75% case where the static allocation is more energy efficient.

4.2.4 I/O Asynchronicity. LABIOS supports both synchronous and asynchronous operations. The potential of a label-based I/O system is more evident by the asynchronous mode where LABIOS can overlap the execution of labels behind other computations. In this test, LABIOS is configured with the round robin scheduling policy, label granularity of 1 MB, and the label dispatcher uses all 48 cores of the node. We scaled the clients from 384 to 3,072 processes (or MPI ranks in this case) to see how LABIOS scales. We run CM1 in 16 iterations (i.e., time steps) with each step first performing computing and then I/O. Each process is performing 32 MB of I/O with the total dataset size reaching 100 GB per step for the largest scale of 3,072. As it can be seen in Figure 14, LABIOS scales well with the synchronous mode, offering competitive performance

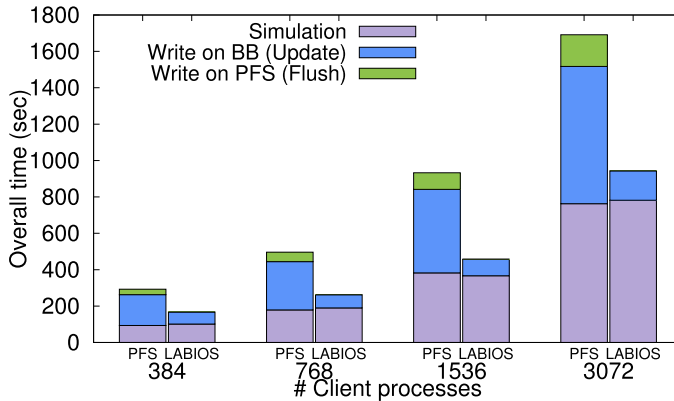


Fig. 15. Resource heterogeneity—HACC performance.

when compared with our baseline, an OrangeFS deployment using the same number of storage servers (i.e., 32 servers). When LABIOS is configured in the async mode, each I/O phase can be executed overlapped with the computation of the next step. This results in a significant 16× I/O performance boost, and a 40% execution time reduction, since the I/O is hidden behind computation. Note that no user code change is required. LABIOS intercepts the I/O calls and builds labels that get executed in a non-blocking fashion.

4.2.5 Resource Heterogeneity. In this test, we run HACC also in 16 time steps. At each step, HACC saves its state on the burst buffers and only at the last step persists the checkpoint data to the remote storage, an OrangeFS deployment. This workload is update-heavy. LABIOS is configured similarly as before but with support of heterogeneous workers, 8 SSD burst buffers and 32 HDD storage servers. LABIOS transparently manages the burst buffers and the servers, and offers 6× I/O performance gains, shown in Figure 15. Moreover, worker to worker flushing is performed in the background.

4.2.6 Data Provisioning. Process-to-process data sharing, *in situ* data analytics, and computational storage capabilities are expressed by provisioning data for each respective purpose. We investigate LABIOS’ ability to support these cases by performing the following tests.

- (a) *Process-to-process data sharing:* in this test, we run Montage, an application that consists of multiple executables that share data between them (i.e., output of one is input to another). LABIOS is configured similarly to the previous set of tests. The baseline uses an OrangeFS deployment of 32 servers. In this test, the simulation produces 50 GB of intermediate data that are written to the PFS and then passed, using temporary files, to the analysis kernel, which produces the final output. As it can be seen in Figure 16, our baseline PFS spends significant time in I/O for this data sharing via the remote storage. This workflow can be significantly boosted by making the data sharing more efficient. LABIOS, instead of sharing intermediate data via the remote storage, passes the labels from the simulation to the analysis via the distributed warehouse. Each intermediate data file creates labels where the destination is not LABIOS workers but the analysis compute nodes. This accelerates the performance in two ways: (a) no temporary files are created in the remote storage servers, and (b) simulation and analysis can now be pipelined (i.e., analysis can start once the first labels are available). As a result, LABIOS offers 65% shorter execution time, boosts I/O performance by 17×, and scales linearly as the number of clients grow.

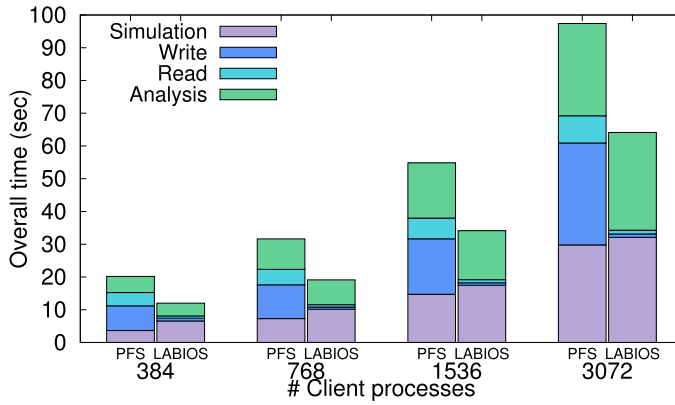


Fig. 16. Data provisioning—Montage performance.

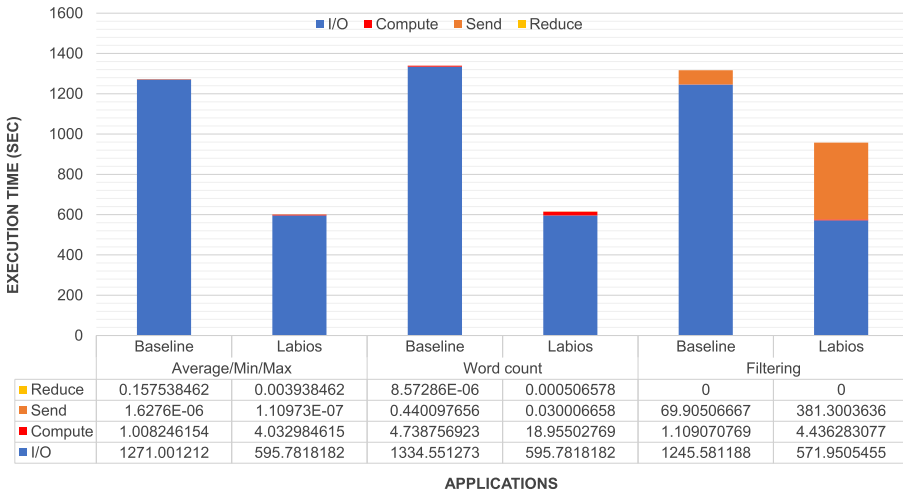


Fig. 17. Data provisioning—*In situ* data analytics performance.

(b) *In situ data analytics*: in this series of tests, we evaluate how effective are LABIOS’ *in situ* analysis capabilities. Through the Label paradigm users can offload certain data intensive operations to the LABIOS workers effectively reducing unnecessary data movements. We run all these tests on Ares cluster, which demonstrates a heterogeneous architecture between compute and storage nodes. Technical details on this hardware composition can be found in subsection 4.1.1. The main idea behind the kernels evaluated in this test is that a large amount of data residing in PFS needs to be analyzed and, thus, moved to compute nodes through the network. This data transfer is expensive and often the bottleneck in such situations. Figure 17 shows results from three data analytics application: a numerical statistics kernel that calculates average, min, and max values over a integer dataset; a word count kernel that counts word frequencies in a textual dataset extracted from Wikipedia articles; and a data filtering kernel that sieves the third quartile of the histogram for a floating point dataset with a gamma distribution. A common theme between these kernels is that the size of input data is significantly larger than the intermediate or output data (i.e., orders of magnitude data reduction). In each of these kernels the amount of input data is

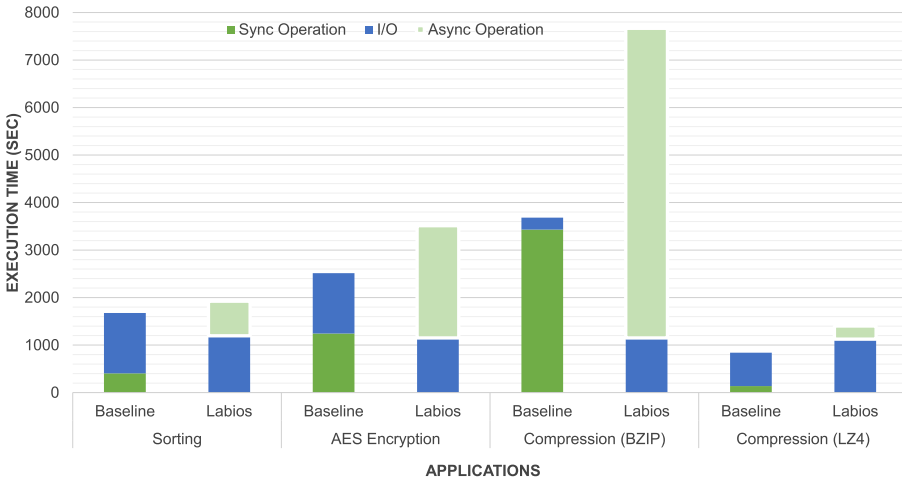


Fig. 18. Data provisioning - Computational storage performance.

4 TB and the number of compute MPI ranks is 2,560. There are 24 storage nodes for the baseline of PFS or LABIOS workers, respectively. Burst buffer nodes are not used. As it can be seen in the figure, for the baseline the I/O time is dominating the overall execution time, since 4 TB of data are transferred from the PFS to the compute nodes. The actual computations are not significantly complex (i.e., in the order of $O(n)$). This is a great fit for *in situ* analysis. Instead of transferring data through the slower network, the functions can be shipped to the storage nodes even if their computational capability is less. In fact, LABIOS cuts the I/O cost down to more than half, since it only performs reads locally. The CPUs present in storage nodes are significantly slower than the compute nodes (i.e., twice as slow) but the reduction in data movement more than justifies the *in situ* analysis. In conclusion, we can observe that the smaller the delta between input and output data size is, the less effective *in situ* analysis becomes. This can be seen in the data filtering kernel where 10% of the data are returned to the compute nodes as the result of the analysis, hence making the LABIOS improvement only 30% over the baseline whereas for statistical and word count kernels is more than 2 \times .

- (c) *Computational Storage*: in this series of test, we evaluate how offloading data operations onto storage can boost applications perceived performance. The main benefit stems from the operation asynchronicity. Applications can exit once they ship their data to storage along with the desired operation to be performed. The following tests were also conducted on the Ares cluster. The main idea behind the applications evaluated in this test is that a large amount of data that reside in compute nodes has to be first processed in some way (i.e., compressed, sorted, etc.) and then written to the remote storage for persistence. The computations applied on the data before they are written out are not necessarily expensive, but they are part of the overall time to exit. Instead, if computational storage is used, then the data transformations can happen in the background by the storage nodes effectively reducing the time to exit. Figure 18 shows results from three application kernels, namely, an integer sorting where data are sorted before written to the remote storage, an AES encryption where data are first encrypted and then written out, and a data compression where all data produced by compute need to first be compressed, reducing their footprint, before written to the archival storage. All tests were conducted with max scale of the cluster (i.e., 2,560 ranks over 64 nodes), and the dataset size is 4 TB. Similarly as

the previous testing, 24 storage nodes were used to run PFS and LABIOS workers, respectively. As it can be seen in the figure, starting from the integer sorting application, for the baseline the sorting algorithm completed in 404 sec and the PFS I/O in 1,271 s for an overall 1,675 s. For LABIOS, the application only needs to write the data (took 1,200 s) and instruct LABIOS to sort them in the background (took 728 s). Note that the computational capabilities of LABIOS workers cannot match that of the compute nodes. However, since the sorting happened in the background (i.e., asynchronously) the application does not experience this slowdown and can exit once I/O has been completed. This leads to a $1.4\times$ speedup. Once the data operation applied to the data before written out to the remote storage becomes more complex the effect of asynchronously execute them on the storage nodes is even more noticeable. Results for the AES data encryption show this clearly. The baseline completed the operation and I/O in 2,525 s, whereas LABIOS allowed the application to exit in only 1,150 s. This is due to the fact that data were written without being encrypted. Data encryption was applied by the LABIOS workers in the background. An interesting observation can be made when we look at the compression case. Compression is a data transformation that spends CPU cycles to perform the compression algorithm for a reduction in data size in return. However, different compression algorithms demonstrate different performance characteristics. For instance, `gzip`, a widely used compression, has a slow compression speed for a high compression ratio. This means that the more time one spends in compressing data the smaller the dataset size becomes. In contrast, when using LABIOS, the application wrote all 4 TB of data to LABIOS workers, which in turn applied compression in the background. However, if we switch the compression algorithm to something more balanced where a more moderate compression ratio is achieved quicker, such as `LZ4`, then we see that the baseline can outperform LABIOS under this configuration where storage nodes' CPU are weaker. The data reduction achieved by lighter compression by `LZ4` was the main reason why the baseline was quicker to finish. In conclusion, computational storage is not a panacea. Its effectiveness depends on few things such as the hardware composition of the storage nodes, the complexity of the computation required, and the application willingness to apply the data operation asynchronously (i.e., in the background).

4.2.7 Storage Bridging. Figure 19 demonstrates the results of running K-means clustering. Our baseline is a 64-node HDFS cluster. LABIOS is configured in two modes: node-local I/O, similar to the HDFS cluster, and remote external storage, similar to an HPC cluster (Section 3.4 (a) and (d)). In the first mode, LABIOS workers run on each of the 64 nodes in the cluster whereas in the second mode, data resides on an external storage running on 32 separate nodes. This application has three distinct phases: (a) *Map*, each mapper reads 32 MB from storage and calculates the squared distance of each point from each initial mean. It then finds the mean that minimizes that distance and writes back to the disk 32 MB of key-value pairs with the mean index as key and the point as value. (b) *Reduce*, each reducer reads 32 MB of key-value pairs written from the mappers and performs a pairwise summation over the values of each key. (c) *Shuffle*, all values across all reducers in the cluster are exchanged via the network (i.e., 32 MB network I/O) to recalculate the new means (i.e., centroids). Finally, it writes the new final centroids back to the disk. An optimized version of this algorithm (i.e., Apache Mahout) avoids writing the key-value pairs back to HDFS during map phase, but instead it emits those values to the reducers avoiding excessive disk I/O (i.e., Hadoop-Memory in Figure 19). This significantly boosts the performance of this algorithm, which is mostly read-intensive, except of the shuffling phase, which is network-heavy [18]. LABIOS supports this workload by having each worker on every node reading the initial dataset

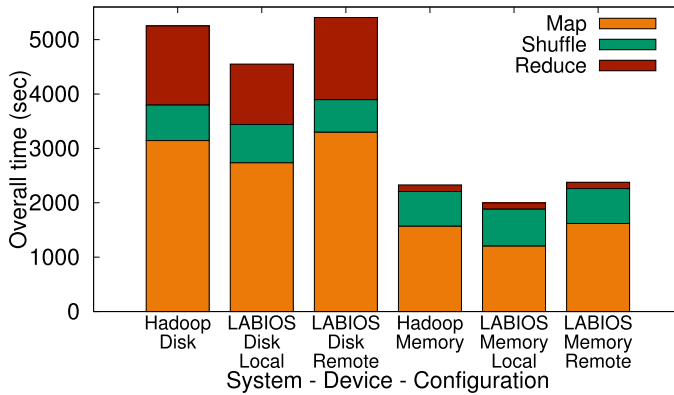


Fig. 19. Storage Bridging—Hadoop K-Means performance.

in an optimized way by performing aggregations, much like MPI collective-I/O where one process reads from storage and distributes the data to all other processes. Further, LABIOS decoupled architecture allows the system to read data from external resources (i.e., LABIOS-Disk-Remote in Figure 19). As it can be seen in the results, reading from external sources is slower than the native node-local I/O mode but it is still a feasible configuration under LABIOS, one that leads to the avoidance of any expensive data movements or data-ingestion approach. In summary, LABIOS supports Hadoop workloads under the same cluster and offers competitive performance with the native HDFS.

5 RELATED WORK

5.0.1 Innovation and New Features in Modern Storage. Fixed reservation with performance guarantees in Ceph [77], in-memory ephemeral storage instances in BeeGFS [32], decoupling of data and metadata path in latest versions of OrangeFS [67], and client-to-client coordination with low server-side coupling in Sirocco [19]. Our work is partially inspired by the above developments. LABIOS is able to offer these features due to its innovative design and architecture.

5.0.2 Active Storage. Comet [28] is an extensible, distributed key-value store that seeks application-specific customization by introducing active storage objects. Comet’s design allows storage operations as a result of executing application specific handlers. ActiveFlash [72] is an *in situ* scientific data analysis approach, wherein data analysis is conducted on where the data already resides. LABIOS workers can independently execute data-intensive operations in a non-blocking fashion, since they are fully decoupled from the clients.

5.0.3 Workflow Interoperability. Running data analysis along with computationally challenging simulations has been explored by Reference [9]. Dataspaces [24] offers a semantically specialized virtual shared space abstraction to support multiple interacting processes and data-intensive application workflows. DAOS [11] integrates a high-performance object store into the HPC storage stack and supports a flexible interface for diverse workloads. However, dedicated analysis resources or expensive data movement between different clusters is still required. LABIOS’ label describes the destination of a certain data operation and can be a memory buffer or a file on another compute node making data sharing easy and efficient.

5.0.4 Task-based Computation Frameworks. Machine independent parallel task-based computing paradigms with new runtime systems such as Charm++ [38] and Legion [3] have been long

advocating for splicing computation to smaller, independent pieces that can be better managed [46], scheduled [48, 50], and executed [53] on heterogeneous environments. LABIOS, in a sense, realizes the same vision of work decomposition but for I/O jobs and not computations.

5.0.5 Storage Malleability. Elasticity is a well explored feature in Cloud storage. Dynamic commission of servers in HDFS [16], transactional database properties with elastic data storage such as ElasTraS [22], and several works exploring energy efficiency in storage systems [4, 49]. LABIOS inherits this feature by its decoupled architecture and the worker pool design and brings storage malleability to HPC as well as BigData.

6 CONCLUSIONS AND FUTURE WORK

Modern large-scale storage systems are required to support a wide range of workflows with different, often conflicting, I/O requirements. Current storage solutions cannot definitively address issues stemming from the scale explosion. In this article, we present the design principles and the architecture of a new, distributed, scalable, elastic, energy-efficient, and fully decoupled label-based I/O system, called LABIOS. We introduce the idea of a *label*, a fundamental piece of LABIOS' architecture, that allows the I/O system to provide storage flexibility, versatility, agility, and malleability. Performance evaluation has shown the potential of LABIOS' architecture by successfully executing multiple conflicting workloads on a single platform. LABIOS can boost I/O performance on certain workloads by up to 17× and reduce overall execution time by 40–60%. Finally, LABIOS provides a platform where users can express intent with software-defined storage abilities and a policy-based execution. As future work, we plan to further develop our system, test it with larger scales, deploy it on more platforms, and extend its functionality with higher fault tolerance semantics, label dependency graphs, and efficient communication protocols.

REFERENCES

- [1] Dong H. Ahn, Ned Bass, Albert Chu, Jim Garlick, Mark Grondona, Stephen Herbein, Helgi I. Ingólfsson, Joseph Koning, Tapasya Patki, Thomas R. W. Scogland, et al. 2020. Flux: Overcoming scheduling challenges for exascale workflows. *Future Gen. Comput. Syst.* 110 (2020), 202–213.
- [2] Amazon Inc. 2018. Amazon S3. Retrieved from <http://docs.aws.amazon.com/AmazonS3/latest/API/Welcome.html>.
- [3] Michael Bauer, Sean Treichler, Elliott Slaughter, and Alex Aiken. 2012. Legion: Expressing locality and independence with logical regions. In *Proceedings of the Conference on High Performance Computing, Networking, Storage and Analysis (SC'12)*. IEEE, 1–11.
- [4] Andreas Berl, Erol Gelenbe, Marco Di Girolamo, Giovanni Giuliani, Hermann De Meer, Minh Quan Dang, and Kostas Pentikousis. 2010. Energy-efficient cloud computing. *Comput. J.* 53, 7 (2010), 1045–1051.
- [5] Dimitris Bertsimas and Ramazan Demir. 2002. An approximate DP approach to multidimensional knapsack problems. *Manage. Sci.* 48, 4 (2002), 550–565.
- [6] Deepavali M. Bhagwat, Marc Eshel, Dean Hildebrand, Manoj P. Naik, Wayne A. Sawdon, Frank B. Schmuck, and Renu Tewari. 2018. Global namespace for a hierarchical set of file systems. U.S. Patent App. 15/397,632.
- [7] Deepavali M. Bhagwat, Marc Eshel, Dean Hildebrand, Manoj P. Naik, Wayne A. Sawdon, Frank B. Schmuck, and Renu Tewari. 2018. Rebuilding the namespace in a hierarchical union mounted file system. U.S. Patent App. 15/397,601.
- [8] Wahid Bhimji, Debbie Bard, Melissa Romanus, David Paul, Andrey Ovsyannikov, Brian Friesen, Matt Bryson, Joaquin Correa, Glenn K. Lockwood, Vakho Tsulaia, et al. 2016. *Accelerating Science with the NERSC Burst Buffer Early User Program*. Technical Report. NERSC.
- [9] John Biddiscombe, Jerome Soumagne, Guillaume Oger, David Guibert, and Jean-Guillaume Piccinali. 2011. Parallel computational steering and analysis for hpc applications using a paraview interface and the hdf5 dsm virtual file driver. In *Proceedings of the Eurographics Symposium on Parallel Graphics and Visualization*. Eurographics Association, 91–100.
- [10] M. K. A. B. V. Bittorf, Taras Bobrovitsky, C. C. A. C. J. Erickson, Martin Grund Daniel Hecht, M. J. I. J. L. Kuff, Dileep Kumar Alex Leblang, N. L. I. P. H. Robinson, David Rorke Silvius Rus, John Russell Dimitris Tsirogiannis Skye Wanderman, and Milne Michael Yoder. 2015. Impala: A modern, open-source SQL engine for Hadoop. In *Proceedings of the 7th Biennial Conference on Innovative Data Systems Research*.

- [11] M. Scot Breitenfeld, Neil Fortner, Jordan Henderson, Jerome Soumagne, Mohamad Chaarawi, Johann Lombardi, and Quincey Koziol. 2017. DAOS for extreme-scale systems in scientific applications. *arXiv* (2017): arXiv-1712.
- [12] George H. Bryan and J. Michael Fritsch. 2002. A benchmark simulation for moist nonhydrostatic numerical models. *Monthly Weather Rev.* 130, 12 (2002), 2917–2928.
- [13] Philip Carns, Sam Lang, Robert Ross, Murali Vilayannur, Julian Kunkel, and Thomas Ludwig. 2009. Small-file access in parallel file systems. In *Proceedings of the IEEE International Symposium on Parallel & Distributed Processing (IPDPS'09)*. IEEE, 1–11.
- [14] Chameleon.org. 2018. Chameleon system. Retrieved from <https://www.chameleoncloud.org/about/chameleon/>.
- [15] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, and Robert E. Gruber. 2008. Bigtable: A distributed storage system for structured data. *ACM Trans. Comput. Syst.* 26, 2 (2008), 4.
- [16] Nathanaël Cherièr, Matthieu Dorier, and Gabriel Antoniu. 2018. *A Lower Bound for the Commission Times in Replication-based Distributed Storage Systems*. Ph.D. Dissertation. Inria Rennes-Bretagne Atlantique.
- [17] Cloud Native Computing Foundation. 2018. NATS Server-C Client. Retrieved from <https://github.com/nats-io/cnats>.
- [18] Xiaoli Cui, Pingfei Zhu, Xin Yang, Keqiu Li, and Changqing Ji. 2014. Optimized big data K-means clustering using MapReduce. *J. Supercomput.* 70, 3 (2014), 1249–1259.
- [19] Matthew L. Curry, H. Lee Ward, and Geoff Danielson. 2015. *Motivation and Design of the Sirocco Storage System Version 1.0*. Technical Report. Sandia National Laboratories. Retrieved from <https://prod-ng.sandia.gov/techlib-noauth/access-control.cgi/2015/156031.pdf>.
- [20] Matthew Curtis-Maury, Vinay Devadas, Vania Fang, and Aditya Kulkarni. 2016. To waffinity and beyond: A scalable architecture for incremental parallelization of file system code. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'16)*. 419–434.
- [21] Matteo D'Ambrosio, Christian Dannewitz, Holger Karl, and Vinicio Vercellone. 2011. MDHT: A hierarchical name resolution service for information-centric networks. In *Proceedings of the ACM Workshop on Information-centric Networking*. ACM, 7–12.
- [22] Sudipto Das, Amr El Abbadi, and Divyakant Agrawal. 2009. ElasTraS: An elastic transactional data store in the cloud. *HotCloud 9* (2009), 131–142.
- [23] Hariharan Devarajan, Anthony Kougkas, X. H. Sun, and H. Chen. 2017. Open ethernet drive: Evolution of energy-efficient storage technology. *Proc. ACM SIGHPC Datacloud 17* (2017).
- [24] Ciprian Docan, Manish Parashar, and Scott Klasky. 2012. Dataspaces: An interaction and coordination framework for coupled simulation workflows. *Cluster Comput.* 15, 2 (2012), 163–181.
- [25] Mike Folk, Albert Cheng, and Kim Yates. 1999. HDF5: A file format and I/O library for high performance computing applications. In *Proceedings of Supercomputing*, Vol. 99. 5–33.
- [26] Kui Gao, Wei-keng Liao, Arifa Nisar, Alok Choudhary, Robert Ross, and Robert Latham. 2009. Using subfilting to improve programming flexibility and performance of parallel shared-file I/O. In *Proceedings of the International Conference on Parallel Processing (ICPP'09)*. IEEE, 470–477.
- [27] Alan Gates. 2012. *HCatalog: An Integration Tool*. Technical Report. Intel.
- [28] Roxana Geambasu, Amit A. Levy, Tadayoshi Kohno, Arvind Krishnamurthy, and Henry M. Levy. 2010. Comet: An active distributed key-value store. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'10)*. 323–336.
- [29] Joachim Giesen, Eva Schuberth, and Miloš Stojaković. 2009. Approximate sorting. *Fundamenta Informaticae* 90, 1–2 (2009), 67–72.
- [30] Google Inc. 2018. CityHash library. Retrieved from <https://github.com/google/cityhash>.
- [31] Grant, W. Shane and Voorhies, Randolph. 2017. Cereal - A C++11 library for serialization by University of Southern California. Retrieved from <http://uscilab.github.io/cereal/>.
- [32] Jan Heichler. 2014. *An Introduction to BeeGFS*. Technical Report.
- [33] Tony Hey, Stewart Tansley, Kristin M. Tolle, et al. 2009. *The Fourth Paradigm: Data-intensive Scientific Discovery*. Vol. 1. Microsoft Research, Redmond, WA.
- [34] IBM. 2018. HDFS Transparency. Retrieved from <https://ibm.co/2Pciyv7>.
- [35] Intel. 2018. Hadoop Adapter for Lustre (HAL). Retrieved from <https://github.com/whamcloud/lustre-connector-for-hadoop>.
- [36] High Performance Data Division Intel Enterprise Edition for Lustre* Software. 2014. *WHITE PAPER Big Data Meets High Performance Computing*. Technical Report. Intel. Retrieved from <https://www.intel.com/content/dam/www/public/us/en/documents/product-briefs/lustre-big-data-white-paper.pdf>.
- [37] Kamil Iskra, John W. Romein, Kazutomo Yoshii, and Pete Beckman. 2008. ZOID: I/O-forwarding infrastructure for petascale architectures. In *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. ACM, 153–162.

- [38] Laxmikant V. Kale and Sanjeev Krishnan. 1996. Charm++: Parallel programming with message-driven objects. *Parallel Programming Using C+* (1996), 175–213.
- [39] Youngjae Kim, Raghul Gunasekaran, Galen M. Shipman, David Dillow, Zhe Zhang, and Bradley W. Settlemyer. 2010. Workload characterization of a leadership class storage cluster. In *Proceedings of the 5th Petascale Data Storage Workshop (PDSW'10)*. IEEE, 1–5.
- [40] Anthony Kougkas, Hariharan Devarajan, and Xian-He Sun. 2018. Hermes: A heterogeneous-aware multi-tiered distributed I/O buffering system. In *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing*. ACM, 219–230.
- [41] Anthony Kougkas, Hariharan Devarajan, and Xian-He Sun. 2018. IRIS: I/O Redirection via Integrated Storage. In *Proceedings of the 32nd ACM International Conference on Supercomputing (ICS'18)*. ACM.
- [42] Anthony Kougkas, Hariharan Devarajan, Xian-He Sun, and Jay Lofstead. 2018. Harmonia: An interference-aware dynamic I/O scheduler for shared non-volatile burst buffers. In *Proceedings of the IEEE Cluster Conference (Cluster'18)*. IEEE.
- [43] Anthony Kougkas, Hassan Eslami, Xian-He Sun, Rajeev Thakur, and William Gropp. 2017. Rethinking key-value store for parallel I/O optimization. *Int. J. High Perform. Comput. Appl.* 31, 4 (2017), 335–356.
- [44] Anthony Kougkas, Anthony Fleck, and Xian-He Sun. 2016. Towards energy efficient data management in hpc: The open ethernet drive approach. In *Proceedings of the 1st Joint International Workshop on Parallel Data Storage and Data Intensive Scalable Computing Systems (PDSW-DISCS'16)*. IEEE, 43–48.
- [45] Haoyuan Li, Ali Ghodsi, Matei Zaharia, Scott Shenker, and Ion Stoica. 2014. Tachyon: Reliable, memory speed storage for cluster computing frameworks. In *Proceedings of the ACM Symposium on Cloud Computing*. ACM, 1–15.
- [46] Jing Li, Jian Jia Chen, Kunal Agrawal, Chenyang Lu, Chris Gill, and Abusayeed Saifullah. 2014. Analysis of federated and global scheduling for parallel real-time tasks. In *Proceedings of the 26th Euromicro Conference on Real-Time Systems (ECRTS'14)*. IEEE, 85–96.
- [47] Jianwei Li, Wei-keng Liao, Alok Choudhary, Robert Ross, Rajeev Thakur, William Gropp, Robert Latham, Andrew Siegel, Brad Gallagher, and Michael Zingale. 2003. Parallel netCDF: A high-performance scientific I/O interface. In *Proceedings of the ACM/IEEE Supercomputing Conference*. ACM/IEEE, 39–39.
- [48] Kenli Li, Xiaoyong Tang, Bharadwaj Veeravalli, and Keqin Li. 2015. Scheduling precedence constrained stochastic tasks on heterogeneous cluster systems. *IEEE Trans. Comput.* 64, 1 (2015), 191–204.
- [49] Harold C. Lim, Shivnath Babu, and Jeffrey S. Chase. 2010. Automated control for elastic storage. In *Proceedings of the 7th International Conference on Autonomic Computing*. ACM, 1–10.
- [50] Juan Liu, Yuyi Mao, Jun Zhang, and Khaled B. Letaief. 2016. Delay-optimal computation task scheduling for mobile-edge computing systems. In *Proceedings of the IEEE International Symposium on Information Theory (ISIT'16)*. IEEE, 1451–1455.
- [51] Yu-Hang Liu and Xian-He Sun. 2015. LPM: Concurrency-driven layered performance matching. In *Proceedings of the 44th International Conference on Parallel Processing (ICPP'15)*. IEEE, 879–888.
- [52] Glenn K. Lockwood, Damian Hazen, Quincey Koziol, R. S. Canon, Katie Antypas, Jan Balewski, Nicholas Balthaser, Wahid Bhimji, James Botts, Jeff Broughton, et al. 2017. *Storage 2020: A Vision for the Future of HPC Storage*. Technical Report. NERSC.
- [53] Yucheng Low, Joseph E. Gonzalez, Aapo Kyrola, Danny Bickson, Carlos E. Guestrin, and Joseph Hellerstein. 2010. Graphlab: A new framework for parallel machine learning. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*. 340–349.
- [54] Memcached. 2018. Extstore plugin. Retrieved from <https://github.com/memcached/memcached/wiki/Extstore>.
- [55] Monty Taylor. 2018. OpenStack Object Storage (Swift). Retrieved from <https://launchpad.net/swift>.
- [56] Wira D. Mulia, Naresh Sehgal, Sohun Sohoni, John M. Acken, C. Lucas Stanberry, and David J. Fritz. 2013. Cloud workload characterization. *IETE Tech. Rev.* 30, 5 (2013), 382–397.
- [57] Ron A. Oldfield, Kenneth Moreland, Nathan Fabian, and David Rogers. 2014. Evaluation of methods to integrate analysis into a large-scale shock physics code. In *Proceedings of the 28th ACM International Conference on Supercomputing*. 83–92. DOI : <https://doi.org/10.1145/2597652.2597668>
- [58] Christopher Olston, Benjamin Reed, Utkarsh Srivastava, Ravi Kumar, and Andrew Tomkins. 2008. Pig latin: A not-so-foreign language for data processing. In *Proceedings of the ACM SIGMOD Conference on Management of Data*. ACM, 1099–1110.
- [59] Fengfeng Pan, Yinliang Yue, Jin Xiong, and Daxiang Hao. 2014. I/O characterization of big data workloads in data centers. In *Proceedings of the Workshop on Big Data Benchmarks, Performance Optimization, and Emerging Hardware*. Springer, 85–97.
- [60] Juan Piernas, Jarek Nieplocha, and Evan J. Felix. 2007. Evaluation of active storage strategies for the lustre parallel file system. In *Proceedings of the ACM/IEEE Conference on Supercomputing*. ACM, 28.
- [61] Jakob Puchinger, Günther R. Raidl, and Ulrich Pferschky. 2010. The multidimensional knapsack problem: Structure and algorithms. *INFORMS J. Comput.* 22, 2 (2010), 250–265.

- [62] Ioan Raicu, Ian Foster, Mike Wilde, Zhao Zhang, Kamil Iskra, Peter Beckman, Yong Zhao, Alex Szalay, Alok Choudhary, Philip Little, et al. 2010. Middleware support for many-task computing. *Cluster Comput.* 13, 3 (2010), 291–314.
- [63] Daniel A. Reed and Jack Dongarra. 2015. Exascale computing and big data. *Commun. ACM* 58, 7 (2015), 56–68.
- [64] Kai Ren, Qing Zheng, Swapnil Patil, and Garth Gibson. 2014. IndexFS: Scaling file system metadata performance with stateless caching and bulk insertion. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC'14)*. IEEE, 237–248.
- [65] Erik Riedel, Christos Faloutsos, Garth A. Gibson, and David Nagle. 2001. Active disks for large-scale data processing. *Computer* 34, 6 (2001), 68–74.
- [66] Erik Riedel, Garth Gibson, and Christos Faloutsos. 1998. Active storage for large-scale data mining and multimedia applications. In *Proceedings of 24th Conference on Very Large Databases*. Citeseer, 62–73.
- [67] Robert B. Ross, Rajeev Thakur, et al. 2000. PVFS: A parallel file system for Linux clusters. In *Proceedings of the 4th Annual Linux Showcase and Conference*.
- [68] Michael W. Shapiro. 2017. Method and system for global namespace with consistent hashing. U.S. Patent 9,787,773.
- [69] Steve Conway. 2015. When Data Needs More Firepower: The HPC, Analytics Convergence. Retrieved from <https://bit.ly/2od68r7>.
- [70] Rajeev Thakur, William Gropp, and Ewing Lusk. 1999. Data sieving and collective I/O in ROMIO. In *Proceedings of the 7th Symposium on the Frontiers of Massively Parallel Computation (Frontiers'99)*. IEEE, 182–189.
- [71] Ashish Thusoo, Joydeep Sen Sarma, Namit Jain, Zheng Shao, Prasad Chakka, Suresh Anthony, Hao Liu, Pete Wyckoff, and Raghotham Murthy. 2009. Hive: A warehousing solution over a map-reduce framework. *Proc. VLDB Endow.* 2, 2 (2009), 1626–1629.
- [72] Devesh Tiwari, Simona Boboila, Sudharshan S. Vazhkudai, Youngjae Kim, Xiaosong Ma, Peter Desnoyers, and Yan Solihin. 2013. Active flash: Toward energy-efficient, *in situ* data analytics on extreme-scale machines. In *Proceedings of the USENIX Conference on File and Storage Technologies (FAST'13)*. 119–132.
- [73] Murali Vilayannur, Partho Nath, and Anand Sivasubramaniam. 2005. Providing tunable consistency for a parallel file store. In *Proceedings of the USENIX Conference on File and Storage Technologies (FAST'05)*, Vol. 5. 2–2.
- [74] Zhenyu Wang and David Garlan. 2000. *Task-driven Computing*. Technical Report. School of Computer Science, Carnegie-Mellon University, Pittsburgh, PA.
- [75] Hakim Weatherspoon and John D. Kubiatowicz. 2002. Erasure coding vs. replication: A quantitative comparison. In *Proceedings of the International Workshop on Peer-to-Peer Systems*. Springer, 328–337.
- [76] Jean-Francois Weets, Manish Kumar Kakhani, and Anil Kumar. 2015. Limitations and challenges of HDFS and MapReduce. In *Proceedings of the International Conference on Green Computing and Internet of Things (ICGCIoT'15)*. IEEE, 545–549.
- [77] Sage A. Weil, Scott A. Brandt, Ethan L. Miller, Darrell D. E. Long, and Carlos Maltzahn. 2006. Ceph: A scalable, high-performance distributed file system. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*. USENIX Association, 307–320.
- [78] Jian Xu and Steven Swanson. 2016. NOVA: A log-structured file system for hybrid volatile/non-volatile main memories. In *Proceedings of the USENIX Conference on File and Storage Technologies (FAST'16)*. 323–338.
- [79] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2012. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation*. USENIX Association, 2–2.
- [80] Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2010. Spark: Cluster computing with working sets. *HotCloud* 10, 10–10 (2010), 95.
- [81] Shuanglong Zhang, Helen Catanese, and An-I. Andy Wang. 2016. The composite-file file system: Decoupling the one-to-one mapping of files and metadata for better performance. In *Proceedings of the USENIX Conference on File and Storage Technologies (FAST'16)*. 15–22.
- [82] Fang Zheng, Hasan Abbasi, Ciprian Docan, Jay Lofstead, Qing Liu, Scott Klasky, Manish Parashar, Norbert Podhorszki, Karsten Schwan, and Matthew Wolf. 2010. PreData—Preparatory data analytics on peta-scale machines. In *Proceedings of the IEEE International Symposium on Parallel & Distributed Processing (IPDPS'10)*. IEEE, 1–12.
- [83] Qing Zheng, Kai Ren, and Garth Gibson. 2014. BatchFS: Scaling the file system control plane with client-funded metadata servers. In *Proceedings of the 9th Parallel Data Storage Workshop*. IEEE, 1–6.
- [84] Shujia Zhou, Bruce H. Van Aartsen, and Thomas L. Clune. 2008. A lightweight scalable I/O utility for optimizing high-end computing applications. In *Proceedings of the IEEE International Symposium on Parallel and Distributed Processing (IPDPS'08)*. IEEE, 1–7.

Received January 2020; revised July 2020; accepted August 2020