



File Systems Unfit as Distributed Storage Backends: Lessons from 10 Years of Ceph Evolution

Abutalib Aghayev
Carnegie Mellon University

Sage Weil
Red Hat, Inc.

Michael Kuchnik
Carnegie Mellon University

Mark Nelson
Red Hat, Inc.

Gregory R. Ganger
Carnegie Mellon University

George Amvrosiadis
Carnegie Mellon University

Abstract

For a decade, the Ceph distributed file system followed the conventional wisdom of building its storage backend on top of local file systems. This is a preferred choice for most distributed file systems today because it allows them to benefit from the convenience and maturity of battle-tested code. Ceph's experience, however, shows that this comes at a high price. First, developing a zero-overhead transaction mechanism is challenging. Second, metadata performance at the local level can significantly affect performance at the distributed level. Third, supporting emerging storage hardware is painstakingly slow.

Ceph addressed these issues with BlueStore, a new backend designed to run directly on raw storage devices. In only two years since its inception, BlueStore outperformed previous established backends and is adopted by 70% of users in production. By running in user space and fully controlling the I/O stack, it has enabled space-efficient metadata and data checksums, fast overwrites of erasure-coded data, inline compression, decreased performance variability, and avoided a series of performance pitfalls of local file systems. Finally, it makes the adoption of backwards-incompatible storage hardware possible, an important trait in a changing storage landscape that is learning to embrace hardware diversity.

CCS Concepts • Information systems → Distributed storage; • Software and its engineering → File systems management; Software performance.

Keywords Ceph, object storage, distributed file system, storage backend, file system

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.
SOSP '19, October 27–30, 2019, Huntsville, ON, Canada

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-6873-5/19/10...\$15.00

<https://doi.org/10.1145/3341301.3359656>

1 Introduction

Distributed file systems operate on a cluster of machines, each assigned one or more roles such as cluster state monitor, metadata server, and storage server. Storage servers, which form the bulk of the machines in the cluster, receive I/O requests over the network and serve them from locally attached storage devices using *storage backend* software. Sitting in the I/O path, the storage backend plays a key role in the performance of the overall system.

Traditionally distributed file systems have used local file systems, such as ext4 or XFS, directly or through middleware, as the storage backend [29, 34, 37, 41, 74, 84, 93, 98, 101, 102]. This approach has delivered reasonable performance, precluding questions on the suitability of file systems as a distributed storage backend. Several reasons have contributed to the success of file systems as the storage backend. First, they allow delegating the hard problems of data persistence and block allocation to a well-tested and highly performant code. Second, they offer a familiar interface (POSIX) and abstractions (files, directories). Third, they enable the use of standard tools (`ls`, `find`) to explore disk contents.

Ceph [98] is a widely-used, open-source distributed file system that followed this convention for a decade. Hard lessons that the Ceph team learned using several popular file systems led them to question the fitness of file systems as storage backends. This is not surprising in hindsight. Stonebraker, after building the INGRES database for a decade, noted that “operating systems offer all things to all people at much higher overhead” [90]. Similarly, exokernels demonstrated that customizing abstractions to applications results in significantly better performance [31, 50]. In addition to the performance penalty, adopting increasingly diverse storage hardware is becoming a challenge for local file systems, which were originally designed for a single storage medium.

The first contribution of this experience paper is to *outline the main reasons behind Ceph's decision to develop BlueStore*, a new storage backend deployed directly on raw storage devices. First, it is hard to implement efficient transactions on top of existing file systems. A significant body of work aims to introduce transactions into file systems [39, 63, 65, 72, 77, 80, 87, 106], but none of these approaches have been adopted due to their high performance overhead, limited functionality, interface complexity, or implementation complexity.

The experience of the Ceph team shows that the alternative options, such as leveraging the limited internal transaction mechanism of file systems, implementing Write-Ahead Logging in user space, or using a transactional key-value store, also deliver subpar performance.

Second, the local file system's metadata performance can significantly affect the performance of the distributed file system as a whole. More specifically, a key challenge that the Ceph team faced was enumerating directories with millions of entries fast, and the lack of ordering in the returned result. Both Btrfs and XFS-based backends suffered from this problem, and directory splitting operations meant to distribute the metadata load were found to clash with file system policies, crippling overall system performance.

At the same time, the rigidity of mature file systems prevents them from adopting emerging storage hardware that abandon the venerable block interface. The history of production file systems shows that on average they take a decade to mature [30, 56, 103, 104]. Once file systems mature, their maintainers tend to be conservative when it comes to making fundamental changes due to the consequences of mistakes. On the other hand, novel storage hardware aimed for data centers introduce backward-incompatible interfaces that require drastic changes. For example, to increase capacity, hard disk drive (HDD) vendors are moving to Shingled Magnetic Recording (SMR) technology [38, 62, 82] that works best with a backward-incompatible *zone* interface [46]. Similarly, to eliminate the long I/O tail latency in solid state drives (SSDs) caused by the Flash Translation Layer (FTL) [35, 54, 108], vendors are introducing Zoned Namespace (ZNS) SSDs that eliminate the FTL, again, exposing the zone interface [9, 26]. Cloud storage providers [58, 73, 109] and storage server vendors [17, 57] are already adapting their private software stacks to use the zoned devices. Distributed file systems, however, are stalled by delays in the adoption of zoned devices in local file systems.

In 2015, the Ceph project started designing and implementing BlueStore, a user space storage backend that stores data directly on raw storage devices, and metadata in a key-value store. By taking full control of the I/O path, BlueStore has been able to efficiently implement full data checksums, inline compression, and fast overwrites of erasure-coded data, while also improving performance on common customer workloads. In 2017, after just two years of development, BlueStore became the default production storage backend in Ceph. A 2018 survey among Ceph users shows that 70% use BlueStore in production with hundreds of petabytes in deployed capacity [60]. As a second contribution, this paper *introduces the design of BlueStore, the challenges its design overcomes, and opportunities for future improvements*. Novelties of BlueStore include (1) storing low-level file system metadata, such as extent bitmaps, in a key-value store, thereby avoiding on-disk format changes and reducing implementation complexity; (2) optimizing clone operations and minimizing the

overhead of the resulting extent reference-counting through careful interface design; (3) BlueFS—a user space file system that enables RocksDB to run faster on raw storage devices; and (4) a space allocator with a fixed 35 MiB memory usage per terabyte of disk space.

In addition to the above contributions, *we perform several experiments that evaluate the improvement of design changes from Ceph's previous production backend, FileStore, to BlueStore*. We experimentally measure the performance effect of issues such as the overhead of journaling file systems, double writes to the journal, inefficient directory splitting, and update-in-place mechanisms (as opposed to copy-on-write).

2 Background

This section aims to highlight the role of distributed storage backends and the features that are essential for building an efficient distributed file system (§ 2.1). We provide a brief overview of Ceph's architecture (§ 2.2) and the evolution of Ceph's storage backend over the last decade (§ 2.3), introducing terms that will be used throughout the paper.

2.1 Essentials of Distributed Storage Backends

Distributed file systems aggregate storage space from multiple physical machines into a single unified data store that offers high-bandwidth and parallel I/O, horizontal scalability, fault tolerance, and strong consistency. While distributed file systems may be designed differently and use unique terms to refer to the machines managing data placement on physical media, the storage backend is usually defined as the software module directly managing the storage device attached to physical machines. For example, Lustre's Object Storage Servers (OSSs) store data on Object Storage Targets [102] (OSTs), GlusterFS's Nodes store data on Bricks [74], and Ceph's Nodes store data on Object Storage Devices (OSDs) [98]. In these, and other systems, the storage backend is the software module that manages space on disks (OSTs, Bricks, OSDs) attached to physical machines (OSSs, Nodes).

Widely-used distributed file systems such as Lustre [102], GlusterFS [74], OrangeFS [29], BeeGFS [93], XtremFS [41], and (until recently) Ceph [98] rely on general local file systems, such as ext4 and XFS, to implement their storage backends. While different systems require different features from a storage backend, two of these features, (1) **efficient transactions** and (2) **fast metadata operations** appear to be common; another emerging requirement is (3) **support for novel, backward-incompatible storage hardware**.

Transaction support in the storage backend simplifies implementing strong consistency that many distributed file systems provide [41, 74, 98, 102]. A storage backend can seamlessly provide transactions if the backing file system already supports them [55, 77]. Yet, most file systems implement the POSIX standard, which lacks a transaction concept.

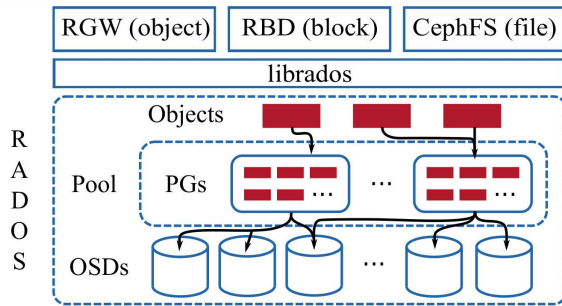


Figure 1. High-level depiction of Ceph’s architecture. A single pool with 3× replication is shown. Therefore, each placement group (PG) is replicated on three OSDs.

Therefore, distributed file system developers typically resort to using inefficient or complex mechanisms, such as implementing a Write-Ahead Log (WAL) on top of a file system [74], or leveraging a file system’s internal transaction mechanism [102].

Metadata management is another recurring pain point in distributed file systems [69]. Inability to efficiently enumerate large directory contents or handle small files at scale in local file systems can cripple performance for both centralized [101, 102] and distributed [74, 98] metadata management designs. To address this problem, distributed file system developers use metadata caching [74], deep directory hierarchies arranged by data hashes [98], custom databases [89], or patches to local file systems [12, 13, 112].

An emerging requirement for storage backends is support for novel storage hardware that operates using backward-incompatible interfaces. For example, SMR can boost HDD capacity by more than 25% and hardware vendors claim that by 2023, over half of data center HDDs will use SMR [83]. Another example is ZNS SSDs that eliminate FTL and do not suffer from uncontrollable garbage collection delays [9], allowing better tail-latency control. Both of these new classes of hardware storage present backward-incompatible interfaces that are challenging for local, block-based file systems to adopt.

2.2 Ceph Distributed Storage System Architecture

Figure 1 shows the high-level architecture of Ceph. At the core of Ceph is the Reliable Autonomic Distributed Object Store (RADOS) service [100]. RADOS scales to thousands of Object Storage Devices (OSDs), providing self-healing, self-managing, replicated object storage with strong consistency. Ceph’s librados library provides a transactional interface for manipulating objects and object collections in RADOS. Out of the box, Ceph provides three services implemented using librados: the RADOS Gateway (RGW), an object storage similar to Amazon S3 [5]; the RADOS Block Device (RBD), a virtual block device similar to Amazon EBS [4]; and CephFS, a distributed file system with POSIX semantics.

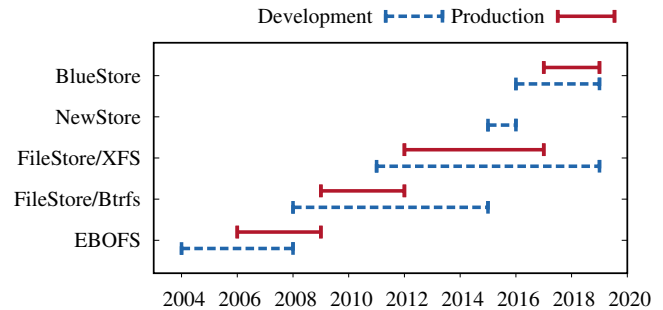


Figure 2. Timeline of storage backend evolution in Ceph. For each backend, the period of development, and the period of being the default production backend is shown.

Objects in RADOS are stored in logical partitions called *pools*. Pools can be configured to provide redundancy for the contained objects either through replication or erasure coding. Within a pool, the objects are sharded among aggregation units called *placement groups* (PGs). Depending on the replication factor, PGs are mapped to multiple OSDs using CRUSH, a pseudo-random data distribution algorithm [99]. Clients also use CRUSH to determine the OSD that should contain a given object, obviating the need for a centralized metadata service. PGs and CRUSH form an indirection layer between clients and OSDs that allows the migration of objects between OSDs to adapt to cluster or workload changes.

In every node of a RADOS cluster, there is a separate *Ceph OSD* daemon per local storage device. Each OSD processes client I/O requests from librados clients and cooperates with peer OSDs to replicate or erasure code updates, migrate data, or recover from failures. Data is persisted to the local device via the internal *ObjectStore* interface, which provides abstractions for objects, object collections, a set of primitives to inspect data, and transactions to update data. A transaction combines an arbitrary number of primitives operating on objects and object collections into an atomic operation. In principle, each OSD may make use of a different backend implementation of the ObjectStore interface, although clusters tend to be uniform in practice.

2.3 Evolution of Ceph’s Storage Backend

The first implementation of the ObjectStore interface was in fact a user space file system called Extent and B-Tree-based Object File System (EBOFS). In 2008, Btrfs was emerging with attractive features such as transactions, deduplication, checksums, and transparent compression, which were lacking in EBOFS. Therefore, as shown in Figure 2, EBOFS was replaced by FileStore, an ObjectStore implementation on top of Btrfs.

In FileStore, an object collection is mapped to a directory and object data is stored in a file. Initially, object attributes were stored in POSIX extended file attributes (*xattrs*), but were later moved to LevelDB when object attributes exceeded

size or count limitations of `xattr`s. FileStore on Btrfs was the production backend for several years, throughout which Btrfs remained unstable and suffered from severe data and metadata fragmentation. In the meantime, the ObjectStore interface evolved significantly, making it impractical to switch back to EBOFS. Instead, FileStore was ported to run on top of XFS, `ext4`, and later ZFS. Of these, FileStore on XFS became the de facto backend because it scaled better and had faster metadata performance [36].

While FileStore on XFS was stable, it still suffered from metadata fragmentation and did not exploit the full potential of the hardware. Lack of native transactions led to a user space WAL implementation that performed full data journaling and capped the speed of read-modify-write workloads, a typical Ceph workload, to the WAL's write speed. In addition, since XFS was not a copy-on-write file system, clone operations used heavily by snapshots were significantly slower.

NewStore was the first attempt at solving the metadata problems of file-system-based backends. Instead of using directories to represent object collections, NewStore stored object metadata in RocksDB, an ordered key-value store, while object data was kept in files. RocksDB was also used to implement the WAL, making read-modify-write workloads efficient due to a combined data and metadata log. Storing object data as files and running RocksDB on top of a journaling file system, however, introduced high consistency overhead. This led to the implementation of BlueStore, which used raw disks. The following section describes the challenges BlueStore aimed to resolve. A complete description of BlueStore is given in § 4.

3 Building Storage Backends on Local File Systems is Hard

This section describes the challenges faced by the Ceph team while trying to build a distributed storage backend on top of local file systems.

3.1 Challenge 1: Efficient Transactions

Transactions simplify application development by encapsulating a sequence of operations into a single atomic unit of work. Thus, a significant body of work aims to introduce transactions into file systems [39, 63, 65, 72, 77, 80, 87, 106]. None of these works have been adopted by production file systems, however, due to their high performance overhead, limited functionality, interface complexity, or implementation complexity.

Hence, there are three tangible options for providing transactions in a storage backend running on top of a file system: (1) hooking into a file system's internal (but limited) transaction mechanism; (2) implementing a WAL in user space; and (3) using a key-value database with transactions as a WAL. Next, we describe why each of these options results in significant performance or complexity overhead.

3.1.1 Leveraging File System Internal Transactions

Many file systems implement an in-kernel transaction framework that enables performing compound internal operations atomically [18, 23, 86, 94]. Since the purpose of this framework is to ensure internal file system consistency, its functionality is generally limited, and thus, unavailable to users. For example, a rollback mechanism is not available in file system transaction frameworks because it is unnecessary for ensuring internal consistency of a file system.

Until recently, Btrfs was making its internal transaction mechanism available to users through a pair of system calls that atomically applied operations between them to the file system [23]. The first version of FileStore that ran on Btrfs relied on these system calls, and suffered from the lack of a rollback mechanism. More specifically, if a Ceph OSD encountered a fatal event in the middle of a transaction, such as a software crash or a KILL signal, Btrfs would commit a partial transaction and leave the storage backend in an inconsistent state.

Solutions attempted by the Ceph and Btrfs teams included introducing a single system call for specifying the entire transaction [96], and implementing rollback through snapshots [95], both of which proved costly. Btrfs authors recently deprecated transaction system calls [15]. This outcome is similar to Microsoft's attempt to leverage NTFS's in-kernel transaction framework for providing an atomic file transaction API, which was deprecated due to its high barrier to entry [53].

These experiences strongly suggest that it is hard to leverage the internal transaction mechanism of a file system in a storage backend implemented in user space.

3.1.2 Implementing the WAL in User Space

An alternative to utilizing the file system's in-kernel transaction framework was to implement a logical WAL in user space. While this approach worked, it suffered from three major problems.

Slow Read-Modify-Write. Typical Ceph workloads perform many read-modify-write operations on objects, where preparing the next transaction requires reading the effect of the previous transaction. A user space WAL implementation, on the other hand, performs three steps for every transaction. First, the transaction is serialized and written to the log. Second, `fsync` is called to commit the transaction to disk. Third, the operations specified in the transaction are applied to the file system. The effect of a transaction cannot be read by upcoming transactions until the third step completes, which is dependent on the second step. As a result, every read-modify-write operation incurred the full latency of the WAL commit, preventing efficient pipelining.

Non-Idempotent Operations. In FileStore, objects are represented by files and collections are mapped to directories.

With this data model, replaying a logical WAL after a crash is challenging due to non-idempotent operations. While the WAL is trimmed periodically, there is always a window of time when a committed transaction that is still in the WAL has already been applied to the file system. For example, consider a transaction consisting of three operations: ① clone $a \rightarrow b$; ② update a ; ③ update c . If a crash happens after the second operation, replaying the WAL corrupts object b . As another example, consider a transaction: ① update b ; ② rename $b \rightarrow c$; ③ rename $a \rightarrow b$; ④ update d . If a crash happens after the third operation, replaying the WAL corrupts object a , which is now named b , and then fails because object a does not exist anymore.

FileStore on Btrfs solved this problem by periodically taking persistent snapshots of the file system and marking the WAL position at the time of snapshot. Then on recovery the latest snapshot was restored, and the WAL was replayed from the position marked at the time of the snapshot.

When FileStore abandoned Btrfs in favor of XFS (§ 2.3), the lack of efficient snapshots caused two problems. First, on XFS the sync system call is the only option for synchronizing file system state to storage. However, in typical deployments with multiple drives per node, sync is too expensive because it synchronizes all file systems on all drives. This problem was resolved by adding the `syncfs` system call [97] to the Linux kernel, which synchronizes only a given file system.

The second problem was that with XFS, there is no option to restore a file system to a specific state after which the WAL can be replayed without worrying about non-idempotent operations. Guards (sequence numbers) were added to avoid replaying non-idempotent operations, however, verifying correctness of guards for complex operations was hard due to the large problem space. Tooling was written to generate random permutations of complex operation sequences, and it was combined with failure injection to semi-comprehensively verify that all failure cases were correctly handled. However, the FileStore code ended up fragile and hard-to-maintain.

Double Writes. The final problem with the WAL in FileStore is that data is written twice: first to the WAL and then to the file system, halving the disk bandwidth. This is a known problem that leads most file systems to only log metadata changes, allowing data loss after a crash. It is possible to avoid the penalty of double writes for new data, by first writing it to disk and then logging only the respective metadata. However, FileStore's approach of using the state of the file system to infer the namespace of objects and their states makes this method hard to use due to corner cases, such as partially written files. While FileStore's approach turned out to be problematic, it was chosen for a technical reason: the alternative required implementing an in-memory cache for data and metadata to any updates waiting on the WAL, despite the kernel having a page and inode cache of its own.

3.1.3 Using a Key-Value Store as the WAL

With NewStore, the metadata was stored in RocksDB, an ordered key-value store, while the object data were still represented as files in a file system. Hence, metadata operations could be performed atomically; data overwrites, however, were logged into RocksDB and executed later. We first describe how this design addresses the three problems of a logical WAL, and then show that it introduces high consistency overhead that stems from running atop a journaling file system.

First, slow read-modify-write operations are avoided because the key-value interface allows reading the new state of an object without waiting for the transaction to commit.

Second, the problem of non-idempotent operation replay is avoided because the read side of such operations is resolved at the time when the transaction is prepared. For example, for clone $a \rightarrow b$, if object a is small, it is copied and inserted into the transaction; if object a is large, a copy-on-write mechanism is used, which changes both a and b to point to the same data and marks the data read-only.

Finally, the problem of double writes is avoided for new objects because the object namespace is now decoupled from the file system state. Therefore, data for a new object is first written to the file system and then a reference to it is atomically added to the database.

Despite these favorable properties, the combination of RocksDB and a journaling file system introduces high consistency overhead, similar to the *journaling of journal* problem [48, 81]. Creating an object in NewStore entails two steps: (1) writing to a file and calling `fsync`, and (2) writing the object metadata to RocksDB synchronously [44], which also calls `fsync`. Ideally, the `fsync` in each step should issue one expensive FLUSH CACHE command [105] to disk. With a journaling file system, however, each `fsync` issues two flush commands: after writing the data, and after committing the corresponding metadata changes to the file system journal. Hence, creating an object in NewStore results in four expensive flush commands to disk.

We demonstrate the overhead of journaling using a benchmark that emulates a storage backend creating many objects. The benchmark has a loop where each iteration first writes 0.5 MiB of data and then inserts a 500-byte metadata to RocksDB. We run the benchmark on two setups. The first setup emulates NewStore, issuing four flush operations for every object creation: data is written as a file to XFS, and the metadata is inserted to stock RocksDB running on XFS. The second setup emulates object creation on raw disk, which issues two flush operations for every object creation: data is written to the raw disk and the metadata is inserted to a modified RocksDB that runs on a raw disk with a preallocated pool of WAL files.

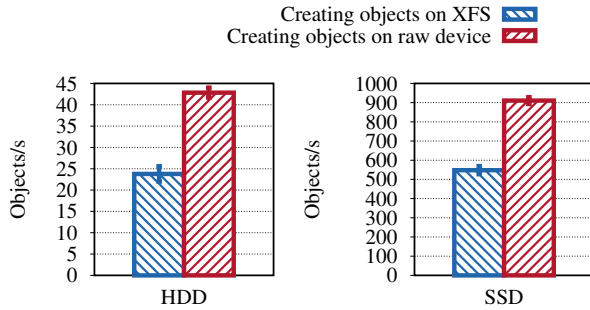


Figure 3. The overhead of running an object store workload on a journaling file system. Object creation throughput is 80% higher on a raw HDD (4 TB Seagate ST4000NM0023) and 70% higher on a raw NVMe SSD (400 GB Intel P3600).

Figure 3 shows that the object creation throughput is 80% higher on raw disk than on XFS when running on a HDD and 70% when running on an NVMe SSD.

3.2 Challenge 2: Fast Metadata Operations

Inefficiency of metadata operations in local file systems is a source of constant struggle for distributed file systems [69, 71, 112]. One of the key metadata challenges in Ceph with the FileStore backend stems from the slow directory enumeration (`readdir`) operations on large directories, and the lack of ordering in the returned result [85].

Objects in RADOS are mapped to a PG based on a hash of their name, and enumerated by hash order. Enumeration is necessary for operations like scrubbing [78], recovery, or for serving `librados` calls that list objects. For objects with long names—as is often the case with RGW—FileStore works around the file name length limitation in local file systems using extended attributes, which may require a `stat` call to determine the object name. FileStore follows a commonly-adopted solution to the slow enumeration problem: a directory hierarchy with large fan-out is created, objects are distributed among directories, and then selected directories’ contents are sorted after being read.

To sort them quickly and to limit the overhead of potential `stat` calls, directories are kept small (a few hundred entries) by splitting them when the number of entries in them grows. This is a costly process at scale, for two primary reasons. First, processing millions of inodes at once reduces the effectiveness of dentry cache, resulting in many small I/Os to disk. And second, XFS places subdirectories in different *allocation groups* [45] to ensure there is space for future directory entries to be located close together [61]; therefore, as the number of objects grows, directory contents spread out, and split operations take longer due to seeks. As a result, when all Ceph OSDs start splitting in unison the performance suffers. This is a well-known problem that has been affecting many Ceph users over the years [16, 27, 88].

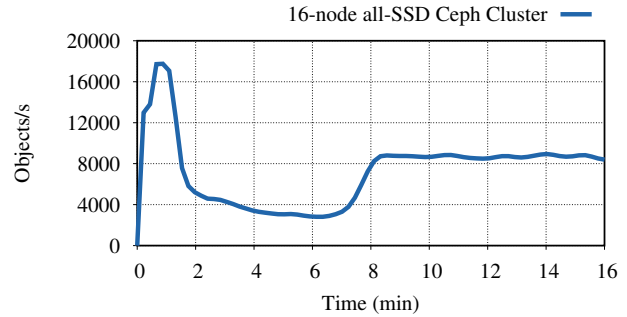


Figure 4. The effect of directory splitting on throughput with FileStore backend. The workload inserts 4 KiB objects using 128 parallel threads at the RADOS layer to a 16-node Ceph cluster (setup explained in § 6). Directory splitting brings down the throughput for 7 minutes on an all-SSD cluster. Once the splitting is complete, the throughput recovers but does not return to peak, due to combination of deeper nesting of files, increased size of the underlying file system, and an imperfect implementation of the directory hashing code in FileStore.

To demonstrate this effect, we configure a 16-node Ceph cluster (§ 6) with roughly half the recommended number of PGs to increase load per PG and accelerate splitting, and insert millions of 4 KiB objects with queue depth of 128 at the RADOS layer (§ 2.2). Figure 4 shows the effect of the splitting on FileStore for an all-SSD cluster. While the first split is not noticeable in the graph, the second split causes a precipitous drop that kills the throughput for 7 minutes on an all-SSD and 120 minutes on an all-HDD cluster (not shown), during which a large and deep directory hierarchy with millions of entries is scanned and even a deeper hierarchy is created. The recovery takes an order of magnitude longer on an all-HDD cluster due to high cost of seeks.

3.3 Challenge 3: Support For New Storage Hardware

The changing storage hardware landscape presents a new challenge for distributed file systems that depend on local file systems. To increase capacity, hard disk drive vendors are shifting to SMR that works best when using a backward-incompatible interface. While the vendors have produced *drive-managed* SMR drives that are backward compatible, these drives have unpredictable performance [1]. For leveraging the extra capacity and achieving predictable performance at the same time, *host-managed* SMR drives with a backward-incompatible zone interface should be used [46]. The zone interface, on the other hand, manages the disk as a sequence of 256 MiB regions that must be written sequentially, encouraging a log-structured, copy-on-write design [76]. This design is in direct opposition to in-place overwrite design followed by most mature file systems.

Data center SSDs are going through a similar change. OpenChannel SSDs eliminate the FTL, leaving the management of raw flash to the host. Lacking an official standard, several vendors have introduced different methods of interfacing OpenChannel SSDs, resulting in fragmented implementations [11, 21, 33]. To prevent this, major vendors have joined forces to introduce a new NVMe standard called Zoned Namespaces (ZNS) that defines an interface for managing SSDs without an FTL [10]. Eliminating the FTL results in many advantages, such as reducing the write amplification, improving latency outliers and throughput, reducing overprovisioning by an order of magnitude, and cutting the cost by reducing DRAM—the highest costing component in SSD after the NAND flash.

Both of these technologies—host-managed SMR drives and ZNS SSDs—are becoming increasingly important for distributed file systems, yet, both have a backward incompatible zone interface that requires radical changes to local file systems [9, 26]. It is not surprising that attempts to modify production file systems, such as XFS and ext4, to work with the zone interface have so far been unsuccessful [19, 68], primarily because these are overwrite file systems, whereas the zone interface requires a copy-on-write approach to data management.

3.4 Other Challenges

Many public and private clouds rely on distributed storage systems like Ceph for providing storage services [67]. Without the complete control of the I/O stack, it is hard for distributed file systems to enforce storage latency SLOs. One cause of high-variance request latencies in file-system-based storage backends is the OS page cache. To improve user experience, most OSs implement the page cache using write-back policy, in which a write operation completes once the data is buffered in memory and the corresponding pages are marked as *dirty*. On a system with little I/O activity, the dirty pages are written back to disk at regular intervals, synchronizing the on-disk and in-memory copies of data. On a busy system, on the other hand, the write-back behavior is governed by a complex set of policies that can trigger writes at arbitrary times [8, 24, 107].

Hence, while the write-back policy results in a responsive system for users with lightly loaded systems, it complicates achieving predictable latency on busy storage backends. Even with a periodic use of `fsync`, FileStore has been unable to bound the amount of deferred inode metadata write-back, leading to inconsistent performance.

Another challenge for file-system-based backends is implementing operations that work better with copy-on-write support, such as snapshots. If the backing file system is copy-on-write, these operations can be implemented efficiently. However, even if the copy-on-write is supported, a file system may have other drawbacks, like fragmentation in FileStore on Btrfs (§ 2.3). If the backing file system is not

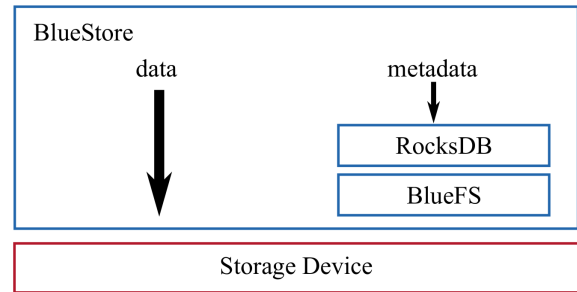


Figure 5. The high-level architecture of BlueStore. Data is written to the raw storage device using direct I/O. Metadata is written to RocksDB running on top of BlueFS. BlueFS is a user space library file system designed for RocksDB, and it also runs on top of the raw storage device.

copy-on-write, then these operations require performing expensive full copies of objects, which makes snapshots and overwriting of erasure-coded data prohibitively expensive in FileStore (§ 5.2).

4 BlueStore: A Clean-Slate Approach

BlueStore is a storage backend designed from scratch to solve the challenges (§ 3) faced by backends using local file systems. Some of the main goals of BlueStore were:

1. Fast metadata operations (§ 4.1)
2. No consistency overhead for object writes (§ 4.1)
3. Copy-on-write clone operation (§ 4.2)
4. No journaling double-writes (§ 4.2)
5. Optimized I/O patterns for HDD and SSD (§ 4.2)

BlueStore achieved all of these goals within just two years and became the default storage backend in Ceph. Two factors played a key role in why BlueStore matured so quickly compared to general-purpose POSIX file systems that take a decade to mature [30, 56, 103, 104]. First, BlueStore implements a small, special-purpose interface, and not a complete POSIX I/O specification. Second, BlueStore is implemented in user space, which allows it to leverage well-tested and high-performance third-party libraries. Finally, BlueStore’s control of the I/O stack enables additional features whose discussion we defer to § 5.

The high-level architecture of BlueStore is shown in Figure 5. BlueStore runs directly on raw disks. A space allocator within BlueStore determines the location of new data, which is asynchronously written to disk using direct I/O. Internal metadata and user object metadata is stored in RocksDB, which runs on BlueFS, a minimal user space file system tailored to RocksDB. The BlueStore space allocator and BlueFS share the disk and periodically communicate to balance free space. The remainder of this section describes metadata and data management in BlueStore.

4.1 BlueFS and RocksDB

BlueStore achieves its first goal, **fast metadata operations**, by storing metadata in RocksDB. BlueStore achieves its second goal of **no consistency overhead** with two changes. First, it writes data directly to raw disk, resulting in one cache flush for data write. Second, it changes RocksDB to reuse WAL files as a circular buffer, resulting in one cache flush for metadata write—a feature that was upstreamed to the mainline RocksDB.

RocksDB itself runs on BlueFS, a minimal file system designed specifically for RocksDB that runs on a raw storage device. RocksDB abstracts out its requirements from the underlying file system in the *Env* interface. BlueFS is an implementation of this interface in the form of a user space, extent-based, and journaling file system. It implements basic system calls required by RocksDB, such as `open`, `mkdir`, and `write`. A possible on-disk layout of BlueFS is shown in Figure 6. BlueFS maintains an inode for each file that includes the list of extents allocated to the file. The superblock is stored at a fixed offset and contains an inode for the journal. The journal has the only copy of all file system metadata, which is loaded into memory at mount time. On every metadata operation, such as directory creation, file creation, and extent allocation, the journal and in-memory metadata are updated. The journal is not stored at a fixed location; its extents are interleaved with other file extents. The journal is compacted and written to a new location when it reaches a preconfigured size, and the new location is recorded in the superblock. These design decisions work because large files and periodic compactions limit the volume of metadata at any point in time.

Metadata Organization. BlueStore keeps multiple namespaces in RocksDB, each storing a different type of metadata. For example, object information is stored in the *O* namespace (that is, RocksDB keys start with *O* and their values represent object metadata), block allocation metadata is stored in the *B* namespace, and collection metadata is stored in the *C* namespace. Each collection maps to a PG and represents a shard of a pool’s namespace. The collection name includes the pool identifier and a prefix shared by the collection’s object names. For example, a key-value pair `C12.e4-6` identifies a collection in pool 12 with objects that have hash values starting with the 6 significant bits of `e4`. Hence, the object `O12.e532` is a member, whereas the object `O12.e832` is not. Such organization of metadata allows a collection of millions of objects to be split into multiple collections merely by changing the number of significant bits. This *collection splitting* operation is necessary to rebalance data across OSDs when, for example, a new OSD is added to the cluster to increase the aggregate capacity or an existing OSD is removed from the cluster due to a malfunction. With FileStore, collection splitting, which is different than directory splitting (§ 3.2), was an expensive operation that was done by renaming directories.

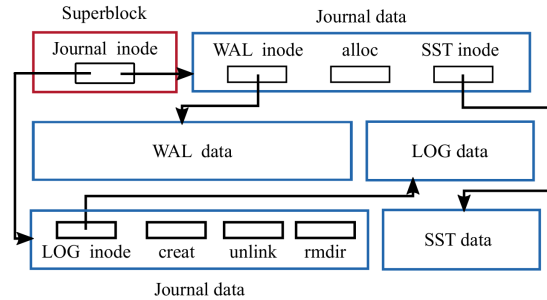


Figure 6. A possible on-disk data layout of BlueFS. The metadata in BlueFS lives only in the journal. The journal does not have a fixed location—its extents are interleaved with file data. The WAL, LOG, and SST files are write-ahead log file, debug log file, and a sorted-string table files, respectively, generated by RocksDB.

4.2 Data Path and Space Allocation

BlueStore is a copy-on-write backend. For incoming writes larger than a *minimum allocation size* (64 KiB for HDDs, 16 KiB for SSDs) the data is written to a newly allocated extent. Once the data is persisted, the corresponding metadata is inserted to RocksDB. This allows BlueStore to provide an **efficient clone** operation. A clone operation simply increments the reference count of dependent extents, and writes are directed to new extents. It also allows BlueStore to **avoid journal double-writes** for object writes and partial overwrites that are larger than the minimum allocation size.

For writes smaller than the minimum allocation size, both data and metadata are first inserted to RocksDB as promises of future I/O, and then asynchronously written to disk after the transaction commits. This deferred write mechanism has two purposes. First, it batches small writes to increase efficiency, because new data writes require two I/O operations whereas an insert to RocksDB requires one. Second, it **optimizes I/O based on the device type**. 64 KiB (or smaller) overwrites of a large object on an HDD are performed asynchronously in place to avoid seeks during reads, whereas in-place overwrites only happen for I/O sizes less than 16 KiB on SSDs.

Space Allocation. BlueStore allocates space using two modules: the FreeList manager and the Allocator. The FreeList manager is responsible for a *persistent* representation of the parts of the disk currently in use. Like all metadata in BlueStore, this free list is also stored in RocksDB. The first implementation of the FreeList manager represented in-use regions as key-value pairs with offset and length. The disadvantage of this approach was that the transactions had to be serialized: the old key had to be deleted first before inserting a new key to avoid an inconsistent free list. The second implementation is bitmap-based. Allocation and deallocation operations use RocksDB’s merge operator to flip

bits corresponding to the affected blocks, eliminating the ordering constraint. The merge operator in RocksDB performs a deferred atomic read-modify-write operation that does not change the semantics and avoids the cost of point queries [43].

The Allocator is responsible for allocating space for the new data. It keeps a copy of the free list in memory and informs the FreeList manager as allocations are made. The first implementation of Allocator was extent-based, dividing the free extents into power-of-two-sized bins. This design was susceptible to fragmentation as disk usage increased. The second implementation uses a hierarchy of indexes layered on top of a single-bit-per-block representation to track whole regions of blocks. Large and small extents can be efficiently found by querying the higher and lower indexes, respectively. This implementation has a fixed memory usage of 35 MiB per terabyte of capacity.

Cache. Since BlueStore is implemented in user space and accesses the disk using direct I/O, it cannot leverage the OS page cache. As a result, BlueStore implements its own write-through cache in user space, using the scan resistant 2Q algorithm [49]. The cache implementation is sharded for parallelism. It uses an identical sharding scheme to Ceph OSDs, which shard requests to collections across multiple cores. This avoids false sharing, so that the same CPU context processing a given client request touches the corresponding 2Q data structures.

5 Features Enabled by BlueStore

In this section we describe new features implemented in BlueStore. These features were previously lacking because implementing them efficiently requires full control of the I/O stack.

5.1 Space-Efficient Checksums

Ceph scrubs metadata every day and data every week. Even with scrubbing, however, if the data is inconsistent across replicas it is hard to be sure which copy is corrupt. Therefore, checksums are indispensable for distributed storage systems that regularly deal with petabytes of data, where bit flips are almost certain to occur.

Most local file systems do not support checksums. When they do, like Btrfs, the checksum is computed over 4 KiB blocks to make block overwrites possible. For 10 TiB of data, storing 32-bit checksums of 4 KiB blocks results in 10 GiB of checksum metadata, which makes it difficult to cache checksums in memory for fast verification.

On the other hand, most of the data stored in distributed file systems is read-only and can be checksummed at a larger granularity. BlueStore computes a checksum for every write and verifies the checksum on every read. While multiple checksum algorithms are supported, `crc32c` is used by default because it is well-optimized on both x86 and ARM

architectures, and it is sufficient for detecting random bit errors. With full control of the I/O stack, BlueStore can choose the checksum block size based on the I/O hints. For example, if the hints indicate that writes are from the S3-compatible RGW service, then the objects are read-only and the checksum can be computed over 128 KiB blocks, and if the hints indicate that objects are to be compressed, then a checksum can be computed after the compression, significantly reducing the total size of checksum metadata.

5.2 Overwrite of Erasure-Coded Data

Ceph has supported erasure-coded (EC) pools (§ 2.2) through the FileStore backend since 2014. However, until BlueStore, EC pools only supported object appends and deletions—overwrites were slow enough to make the system unusable. As a result, the use of EC pools were limited to RGW; for RBD and CephFS only replicated pools were used.

To avoid the “RAID write hole” problem [92], where crashing during a multi-step data update can leave the system in an inconsistent state, Ceph performs overwrites in EC pools using two-phase commit. First, all OSDs that store a chunk of the EC object make a copy of the chunk so that they can roll back in case of failure. After all of the OSDs receive the new content and overwrite their chunks, the old copies are discarded. With FileStore on XFS, the first phase is expensive because each OSD performs a physical copy of its chunk. BlueStore, however, makes overwrites practical because its copy-on-write mechanism avoids full physical copies.

5.3 Transparent Compression

Transparent compression is crucial for scale-out distributed file systems because 3× replication increases storage costs [32, 40]. BlueStore implements transparent compression where written data is automatically compressed before being stored.

Getting the full benefit of compression requires compressing over large 128 KiB chunks, and compression works well when objects are written in their entirety. For partial overwrites of a compressed object, BlueStore places the new data in a separate location and updates metadata to point to it. When the compressed object gets too fragmented due to multiple overwrites, BlueStore compacts the object by reading and rewriting. In practice, however, BlueStore uses hints and simple heuristics to compress only those objects that are unlikely to experience many overwrites.

5.4 Exploring New Interfaces

Despite multiple attempts [19, 68], local file systems are unable to leverage the capacity benefits of SMR drives due to their backward-incompatible interface, and it is unlikely that they will ever do so efficiently [28, 30]. Supporting these denser drives, however, is important for scale-out distributed file systems because it lowers storage costs [59].

Unconstrained by the block-based designs of local file systems, BlueStore has the freedom of exploring novel interfaces and data layouts. This has recently enabled porting RocksDB and BlueFS (§ 4.1) to run on host-managed SMR drives, and an effort is underway to store object data on such drives next [3]. In addition, the Ceph community is exploring a new backend that targets a combination of persistent memory and emerging NVMe devices with novel interfaces, such as ZNS SSDs [9, 26] and KV SSDs [51].

6 Evaluation

This section compares the performance of a Ceph cluster using FileStore, a backend built on a local file system, and BlueStore, a backend using the storage device directly. We first compare the throughput of object writes to the RADOS distributed object storage (§ 6.1). Then, we compare the end-to-end throughput of random writes, sequential writes, and sequential reads to RBD, the Ceph virtual block device built on RADOS (§ 6.2). Finally, we compare the throughput of random writes to an RBD device allocated on an erasure-coded pool (§ 6.3).

We run all experiments on a 16-node Ceph cluster connected with a Cisco Nexus 3264-Q 64-port QSFP+ 40GbE switch. Each node has a 16-core Intel E5-2698Bv3 Xeon 2GHz CPU, 64GiB RAM, 400GB Intel P3600 NVMe SSD, 4TB 7200RPM Seagate ST4000NM0023 HDD, and a Mellanox MCX314A-BCCT 40GbE NIC. All nodes run Linux kernel 4.15 on Ubuntu 18.04, and the Luminous release (v12.2.11) of Ceph. We use the default Ceph configuration parameters.

6.1 Bare RADOS Benchmarks

We start by comparing the performance of object writes to RADOS when using the FileStore and BlueStore backends. We focus on write performance improvements because most BlueStore optimizations affect writes.

Figure 7 shows the throughput for different object sizes written with a queue depth of 128. At the steady state, the throughput on BlueStore is 50-100% greater than FileStore. The throughput improvement on BlueStore stems from avoiding double writes (§ 3.1.2) and consistency overhead (§ 3.1.3).

Figure 8 shows the 95th and above percentile latencies of object writes to RADOS. BlueStore has an order of magnitude lower tail latency than FileStore. In addition, with BlueStore the tail latency increases with the object size, as expected, whereas with FileStore even small-sized object writes may have high tail latency, stemming from the lack of control over writes (§ 3.4).

The read performance on BlueStore (not shown) is similar or better than on FileStore for I/O sizes larger than 128 KiB; for smaller I/O sizes FileStore is better because of the kernel read-ahead [6]. BlueStore does not implement read-ahead on purpose. It is expected that the applications implemented on top of RADOS will perform their own read-ahead.

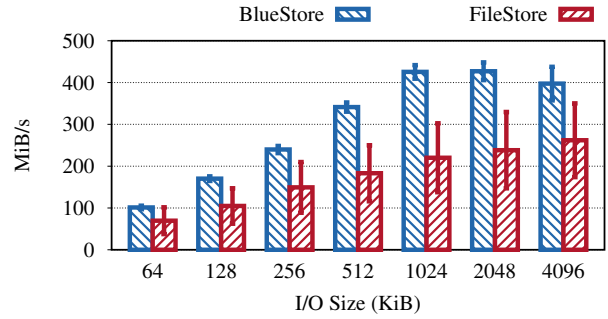


Figure 7. Throughput of steady state object writes to RADOS on a 16-node all-HDD cluster with different sizes using 128 threads. Compared to FileStore, the throughput is 50-100% greater on BlueStore and has a significantly lower variance.

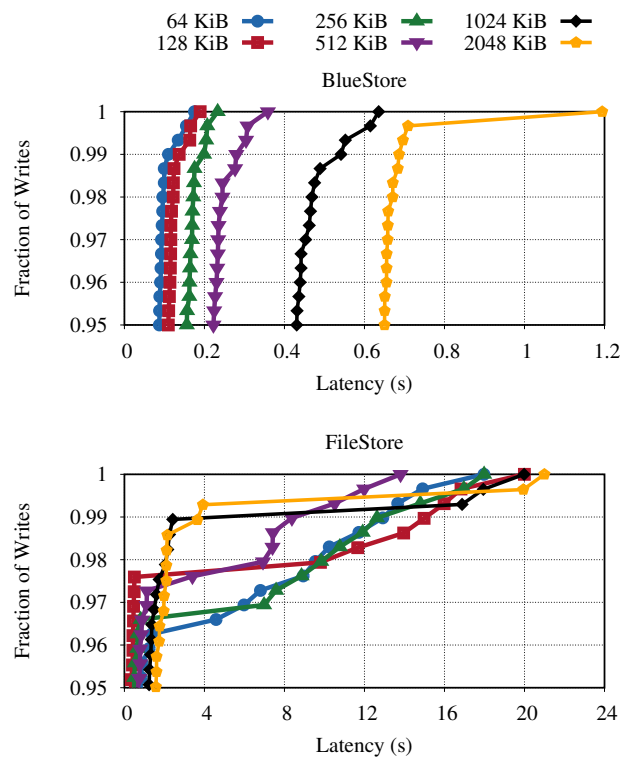


Figure 8. 95th and above percentile latencies of object writes to RADOS on a 16-node all-HDD cluster with different sizes using 128 threads. BlueStore (top graph) has an order of magnitude lower tail latency than FileStore (bottom graph).

BlueStore eliminates the directory splitting effect of FileStore by storing metadata in an ordered key-value store. To demonstrate this, we repeat the experiment that showed the splitting problem in FileStore (§ 3.2) on an identically configured Ceph cluster using a BlueStore backend. Figure 9 shows that the throughput on BlueStore does not suffer the precipitous drop, and in the steady state it is 2× higher than FileStore throughput on SSD (and 3× higher than FileStore

throughput on HDD—not shown). Still, the throughput on BlueStore drops significantly before reaching a steady state due to RocksDB compaction whose cost grows with the object corpus.

6.2 RADOS Block Device Benchmarks

Next, we compare the performance of RADOS Block Device (RBD), a virtual block device service implemented on top of RADOS, when using the BlueStore and FileStore backends. RBD is implemented as a kernel module that exports a block device to the user, which can be formatted and mounted like a regular block device. Data written to the device is striped into 4 MiB RADOS objects and written in parallel to multiple OSDs over the network.

For RBD benchmarks we create a 1 TB virtual block device, format it with XFS, and mount it on the client. We use `fiio` [7] to perform sequential and random I/O with queue depth of 256 and I/O sizes ranging from 4 KiB to 4 MiB. For each test, we write about 30 GiB of data. Before starting every experiment, we drop the OS page cache for FileStore, and we restart OSDs for BlueStore to eliminate caching effects in read experiments. We first run all the experiments on a Ceph cluster installed with FileStore backend. We then tear down the cluster, reinstall it with BlueStore backend, and repeat all the experiments.

Figure 10 shows the results for sequential writes, random writes, and sequential reads. For I/O sizes larger than 512 KiB, sequential and random write throughput is on average 1.7× and 2× higher with BlueStore, respectively, again mainly due to avoiding double-writes. BlueStore also displays a significantly lower throughput variance because it can deterministically push data to disk. In FileStore, on the other hand, arbitrarily-triggered writeback (§ 3.4) conflicts with the foreground writes to the WAL and introduces long request latencies.

For medium I/O sizes (128–512 KiB) the throughput difference decreases for sequential writes because XFS masks out part of the cost of double writes in FileStore. With medium I/O sizes the writes to WAL do not fully utilize the disk. This leaves enough bandwidth for another write stream to go through and not have a large impact on the foreground writes to WAL. After writing the data synchronously to the WAL, FileStore then asynchronously writes it to the file system. XFS buffers these asynchronous writes and turns them into one large sequential write before issuing to disk. XFS cannot do the same for random writes, which is why the high throughput difference continues even for medium-sized random writes.

Finally, for I/O sizes smaller than 64 KiB (not shown) the throughput of BlueStore is 20% higher than that of FileStore. For these I/O sizes BlueStore performs deferred writes by inserting data to RocksDB first, and then asynchronously overwriting the object data to avoid fragmentation (§ 4.2).

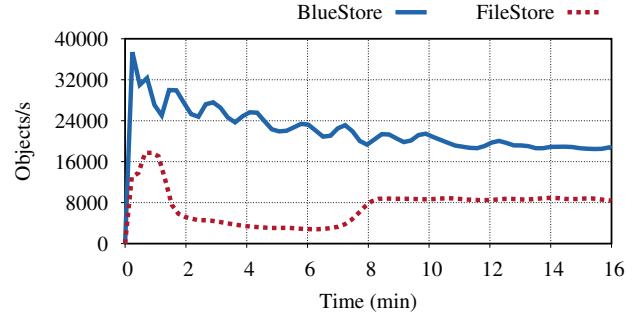


Figure 9. Throughput of 4 KiB RADOS object writes with queue depth of 128 on a 16-node all-SSD cluster. At steady state, BlueStore is 2× faster than FileStore on SSD. BlueStore does not suffer from directory splitting; however, its throughput is gradually brought down by the RocksDB compaction overhead.

The throughput of read operations in BlueStore is similar or slightly better than that of FileStore for I/O sizes larger than 32 KiB. For smaller I/O sizes, as the rightmost graph in Figure 10 shows, FileStore throughput is better because of the kernel readahead. While RBD does implement a readahead, it is not as well-tuned as the kernel readahead.

6.3 Overwriting Erasure-Coded Data

One of the features enabled by BlueStore is the efficient overwrite of EC data. We have measured the throughput of random overwrites for both BlueStore and FileStore. Our benchmark creates 1 TB RBD using one client. The client mounts the block device and performs 5 GiB of random 4 KiB writes with queue depth of 256. Since the RBD is striped in 4 MiB RADOS objects, every write results in an object overwrite. We repeat the experiment on a virtual block device allocated on a replicated pool, on an EC pool with parameters $k = 4$ and $m = 2$ (EC4-2), and $k = 5$ and $m = 1$ (EC5-1).

Figure 11 compares the throughput of replicated and EC pools when using BlueStore and FileStore backends. BlueStore EC pools achieve 6× more IOPS on EC4-2, and 8× more IOPS on EC5-1 than FileStore. This is due to BlueStore avoiding full physical copies during the first phase of the two-phase commit required for overwriting EC objects (§ 5.2). As a result, it is practical to use EC pools with applications that require data overwrite, such as RBD and CephFS, with the BlueStore backend.

7 Challenges of Building Efficient Storage Backends on Raw Storage

This section describes some of the challenges that the Ceph team faced when building a storage backend on raw storage devices from scratch.

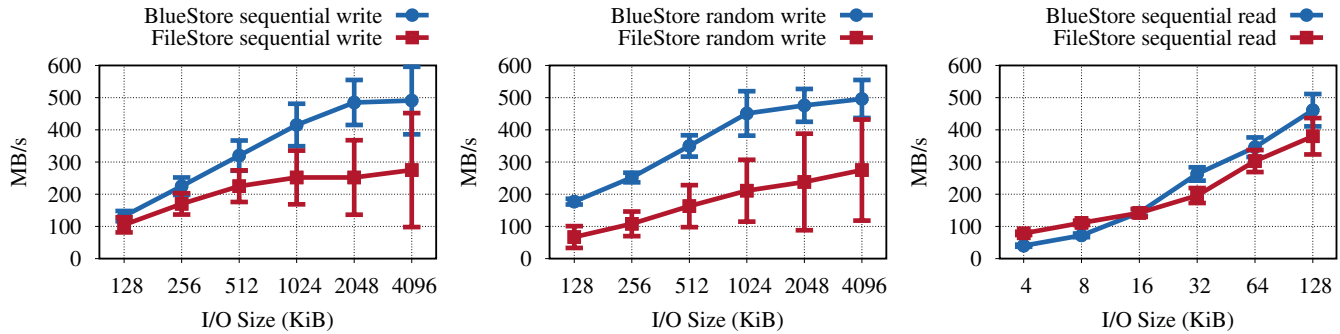


Figure 10. Sequential write, random write, and sequential read throughput with different I/O sizes and queue depth of 256 on a 1 TB Ceph virtual block device (RBD) allocated on a 16-node all-HDD cluster. Results for an all-SSD cluster were similar but not shown for brevity.

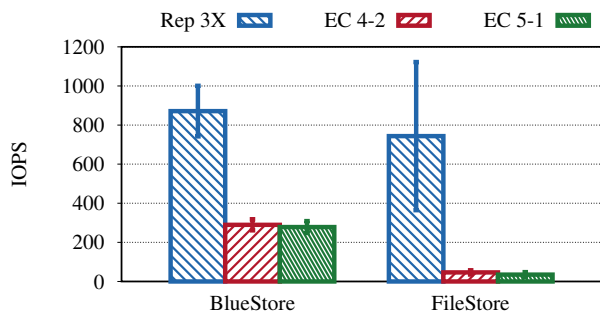


Figure 11. IOPS observed from a client performing random 4 KiB writes with queue depth of 256 to a Ceph virtual block device (RBD). The device is allocated on a 16-node all-HDD cluster.

7.1 Cache Sizing and Writeback

The OS fully utilizes the machine memory by dynamically growing or shrinking the size of the page cache based on the application’s memory usage. It writes back the dirty pages to disk in the background trying not to adversely affect foreground I/O, so that memory can be quickly reused when applications ask for it.

A storage backend based on a local file system automatically inherits the benefits of the OS page cache. A storage backend that bypasses the local file system, however, has to implement a similar mechanism from scratch (§ 4.2). In BlueStore, for example, the cache size is a fixed configuration parameter that requires manual tuning. Building an efficient user space cache with the dynamic resizing functionality of the OS page cache is an open problem shared by other projects, like PostgreSQL [25] and RocksDB [42]. With the arrival of fast NVMe SSDs, such a cache needs to be efficient enough that it does not incur overhead for write-intensive workloads—a deficiency that current page cache suffers from [20].

7.2 Key-value Store Efficiency

The experience of the Ceph team demonstrates that moving all of the metadata to an ordered key-value store, like RocksDB, significantly improves the efficiency of metadata operations. However, the Ceph team has also found that embedding RocksDB in BlueStore is problematic in multiple ways: (1) RocksDB’s compaction and high write amplification have been the primary performance limiters when using NVMe SSDs in OSDs; (2) since RocksDB is treated as a black box, data is serialized and copied in and out of it, consuming CPU time; and (3) RocksDB has its own threading model, which limits the ability to do custom sharding. These and other problems with RocksDB and similar key-value stores keeps the Ceph team researching better solutions.

7.3 CPU and Memory Efficiency

Modern compilers align and pad basic datatypes in memory so that CPU can fetch data efficiently, thereby increasing performance. For applications with complex structs, the default layout can waste a significant amount of memory [22, 64]. Many applications are rightly not concerned with this problem because they allocate short-lived data structures. A storage backend that bypasses the OS page cache, on the other hand, runs continuously and controls almost all of a machine’s memory. Therefore, the Ceph team spent a lot of time packing structures stored in RocksDB to reduce the total metadata size and also compaction overhead. The main tricks used were delta and variable-integer encoding.

Another observation with BlueStore is that on high-end NVMe SSDs the workloads are becoming increasingly CPU-bound. For its next-generation backend, the Ceph community is exploring techniques that reduce CPU consumption, such as minimizing data serialization-deserialization, and using the SeaStar framework [79] with shared-nothing model that avoids context switches due to locking.

8 Related Work

The primary motivator for BlueStore is the lack of transactions and unscalable metadata operations in local file systems. In this section we compare BlueStore to previous research that aims to address these problems.

Transaction Support. Previous works have generally followed three approaches when introducing transactional interface to file system users.

The first approach is to leverage the in-kernel transaction mechanism present in the file systems. Examples of this are Btrfs' export of transaction system calls to userspace [23], Transactional NTFS [52], Valor [87], and TxFS [39]. The drawbacks of this approach are the complexity and incompleteness of the interface, and the a significant implementation complexity. For example, Btrfs and NTFS both recently deprecated their transaction interface [15, 53] citing difficulty guaranteeing correct or safe usage, which corroborates FileStore's experience (§ 3.1.1). Valor [87], while not tied to a specific file system, also has a nuanced interface that requires correct use of a combo of seven system calls, and a complex in-kernel implementation. TxFS is a recent work that introduces a simple interface built on ext4's journaling layer; however, its implementation requires non-trivial amount of change to the Linux kernel. BlueStore, informed by FileStore's experience, avoids using file systems' in-kernel transaction infrastructure.

The second approach builds a user space file system atop a database, utilizing existing transactional semantics. For example, Amino [106] relies on Berkeley DB [66] as the backing store, and Inversion [65] stores files in a POSTGRES database [91]. While these file systems provide seamless transactional operations, they generally suffer from high performance overhead because they accrue the overhead of the layers below. BlueStore similarly leverages a transactional database, but incurs zero overhead because it eliminates the local file system and runs the database on a raw disk.

The third approach provides transactions as a first-class abstraction in the OS and implements all services, including the file system, using transactions. QuickSilver [77] is an example of such system that uses built-in transactions for implementing a storage backend for a distributed file system. Similarly, TxOS [72] adds transactions to the Linux kernel and converts ext3 into a transactional file system. This approach, however, is too heavyweight for achieving file system transactions, and such a kernel is tricky to maintain [39].

Metadata Optimizations. A large body of work has produced a plethora of approaches to metadata optimizations in local file systems. BetrFS [47] introduces B^e-Tree as an indexing structure for efficient large scans. DualFS [70], hFS [110], and ext4-lazy [2] abandon traditional FFS [61] cylinder group design and aggregate all metadata in one place to achieve

significantly faster metadata operations. TableFS [75] and DeltaFS [111] store metadata in LevelDB running atop a file system and achieve orders of magnitude faster metadata operations than local file systems.

While BlueStore also stores metadata in RocksDB—a LevelDB derivative—to achieve similar speedup, it differs from the above in two important ways: (1) in BlueStore, RocksDB runs atop a raw disk incurring zero overhead, and (2) BlueStore keeps all metadata, including the internal metadata, in RocksDB as key-value pairs. Storing internal metadata as variable-sized key-value pairs, as opposed to fixed-sized records on disk, scales more easily. For example, the Lustre distributed file system that uses an ext4-derivate called LD-ISKFS for the storage backend, has changed on-disk format twice in a short period to accommodate for increasing disk sizes [12, 13].

9 Conclusion

Distributed file system developers conventionally adopt local file systems as their storage backend. They then try to fit the general-purpose file system abstractions to their needs, incurring significant accidental complexity [14]. At the core of this convention lies the belief that developing a storage backend from scratch is an arduous process, akin to developing a new file system that takes a decade to mature.

Our paper, relying on the Ceph team's experience, shows this belief to be inaccurate. Furthermore, we find that developing a *special-purpose*, user space storage backend from scratch (1) reclaims the significant performance left on the table when building a backend on a general-purpose file system, (2) makes it possible to adopt novel, backward incompatible storage hardware, and (3) enables new features by gaining complete control of the I/O stack. We hope that this experience paper will initiate discussions among storage practitioners and researchers on fresh approaches to designing distributed file systems and their storage backends.

Acknowledgments

We thank Robert Morris (our shepherd), Matias Bjørling, and the anonymous reviewers for their feedback. We would like to acknowledge the BlueStore authors, which include Igor Fedotov, Xie Xingguo, Ma Jianpeng, Allen Samuels, Varada Kari, and Haomai Wang. We also thank the members and companies of the PDL Consortium: Alibaba Group, Amazon, Datrium, Facebook, Google, Hewlett Packard Enterprise, Hitachi Ltd., Intel Corporation, IBM, Micron, Microsoft Research, NetApp, Inc., Oracle Corporation, Salesforce, Samsung Semiconductor Inc., Seagate Technology, and Two Sigma for their interest, insights, feedback, and support.

Abutalib Aghayev is supported by an SOSP 2019 student scholarship from the National Science Foundation. Michael Kuchnik is supported by an SOSP 2019 student scholarship from the ACM Special Interest Group in Operating Systems and by an NDSEG Fellowship.

References

- [1] Abutalib Aghayev and Peter Desnoyers. 2015. Skylight—A Window on Shingled Disk Operation. In *13th USENIX Conference on File and Storage Technologies (FAST 15)*. USENIX Association, Santa Clara, CA, USA, 135–149. <https://www.usenix.org/conference/fast15/technical-sessions/presentation/aghayev>
- [2] Abutalib Aghayev, Theodore Ts'o, Garth Gibson, and Peter Desnoyers. 2017. Evolving Ext4 for Shingled Disks. In *15th USENIX Conference on File and Storage Technologies (FAST 17)*. USENIX Association, Santa Clara, CA, 105–120. <https://www.usenix.org/conference/fast17/technical-sessions/presentation/aghayev>
- [3] Abutalib Aghayev, Sage Weil, Greg Ganger, and George Amvrosiadis. 2019. *Reconciling LSM-Trees with Modern Hard Drives using BlueFS*. Technical Report CMU-PDL-19-102. CMU Parallel Data Laboratory. http://www.pdl.cmu.edu/PDL-FTP/FS/CMU-PDL-19-102_abs.shtml
- [4] Amazon.com, Inc. 2019. Amazon Elastic Block Store. <https://aws.amazon.com/ebs/>.
- [5] Amazon.com, Inc. 2019. Amazon S3. <https://aws.amazon.com/s3/>.
- [6] Jens Axboe. 2009. Queue sysfs files. <https://www.kernel.org/doc/Documentation/block/queue-sysfs.txt>.
- [7] Jens Axboe. 2016. Flexible I/O Tester. <git://git.kernel.dk/fio.git>.
- [8] Jens Axboe. 2016. Throttled Background Buffered Writeback. <https://lwn.net/Articles/698815/>.
- [9] Matias Björling. 2019. From Open-Channel SSDs to Zoned Namespaces. In *Linux Storage and Filesystems Conference (Vault 19)*. USENIX Association, Boston, MA.
- [10] Matias Björling. 2019. New NVMe Specification Defines Zoned Namespaces (ZNS) as Go-To Industry Technology. <https://nvmexpress.org/new-nvmetm-specification-defines-zoned-namespaces-zns-as-go-to-industry-technology/>.
- [11] Matias Björling, Javier Gonzalez, and Philippe Bonnet. 2017. Light-NVM: The Linux Open-Channel SSD Subsystem. In *15th USENIX Conference on File and Storage Technologies (FAST 17)*. USENIX Association, Santa Clara, CA, 359–374. <https://www.usenix.org/conference/fast17/technical-sessions/presentation/bjorling>
- [12] Artem Blagodarenko. 2016. Scaling LDISKFS for the future. <https://www.youtube.com/watch?v=ubbZGpxV6zk>.
- [13] Artem Blagodarenko. 2017. Scaling LDISKFS for the future. Again. https://www.youtube.com/watch?v=HLfEd0_Dq0U.
- [14] Frederick P Brooks Jr. 1986. No Silver Bullet—Essence and Accident in Software Engineering.
- [15] Btrfs. 2019. Btrfs Changelog. <https://btrfs.wiki.kernel.org/index.php/Changelog>.
- [16] David C. 2018. [ceph-users] Luminous | PG split causing slow requests. <http://lists.ceph.com/pipermail/ceph-users-ceph.com/2018-February/024984.html>.
- [17] Luoqing Chao and Thunder Zhang. 2015. Implement Object Storage with SMR based key-value store. https://www.snia.org/sites/default/files/SDC15_presentations/smr/QingchaoLuo_Implement_Object_Storage_SMR_Key-Value_Store.pdf.
- [18] Dave Chinner. 2010. XFS Delayed Logging Design. <https://www.kernel.org/doc/Documentation/filesystems/xfs-delayed-logging-design.txt>.
- [19] Dave Chinner. 2015. SMR Layout Optimization for XFS. <http://xfs.org/images/f/f6/Xfs-smr-structure-0.2.pdf>.
- [20] Dave Chinner. 2019. Re: pagecache locking (was: bcachefs status update) merged. <https://lkml.org/lkml/2019/6/13/1794>.
- [21] Alibaba Cloud. 2018. Alibaba Deploys Alibaba Open Channel SSD for Next Generation Data Centers. https://www.alibabacloud.com/blog/alibaba-deploys-alibaba-open-channel-ssd-for-next-generation-data-centers_593802.
- [22] William Cohen. 2016. How to avoid wasting megabytes of memory a few bytes at a time. <https://developers.redhat.com/blog/2016/06/01/how-to-avoid-wasting-megabytes-of-memory-a-few-bytes-at-a-time/>.
- [23] Jonathan Corbet. 2009. Supporting transactions in Btrfs. <https://lwn.net/Articles/361457/>.
- [24] Jonathan Corbet. 2011. No-I/O dirty throttling. <https://lwn.net/Articles/456904/>.
- [25] Jonathan Corbet. 2018. PostgreSQL's fsync() surprise. <https://lwn.net/Articles/752063/>.
- [26] Western Digital. 2019. Zoned Storage. <http://zonedstorage.io>.
- [27] Anton Dmitriev. 2017. [ceph-users] All OSD fails after few requests to RGW. <http://lists.ceph.com/pipermail/ceph-users-ceph.com/2017-May/017950.html>.
- [28] Jake Edge. 2015. Filesystem support for SMR devices. <https://lwn.net/Articles/637035/>.
- [29] Jake Edge. 2015. The OrangeFS distributed filesystem. <https://lwn.net/Articles/643165/>.
- [30] Jake Edge. 2015. XFS: There and back ... and there again? <https://lwn.net/Articles/638546/>.
- [31] D. R. Engler, M. F. Kaashoek, and J. O'Toole, Jr. 1995. Exokernel: An Operating System Architecture for Application-level Resource Management. In *Proceedings of the Fifteenth ACM Symposium on Operating Systems Principles (SOSP '95)*. ACM, New York, NY, USA, 251–266. <https://doi.org/10.1145/224056.224076>
- [32] Andrew Fikes. 2010. Storage Architecture and Challenges. https://cloud.google.com/files/storage_architecture_and_challenges.pdf.
- [33] Mary Jo Foley. 2018. Microsoft readies new cloud SSD storage spec for the Open Compute Project. <https://www.zdnet.com/article/microsoft-readies-new-cloud-ssd-storage-spec-for-the-open-compute-project/>.
- [34] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. 2003. The Google File System. In *Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles (SOSP '03)*. ACM, New York, NY, USA, 29–43. <https://doi.org/10.1145/945445.945450>
- [35] Mingzhe Hao, Gokul Soundararajan, Deepak Kenchamma-Hosekote, Andrew A. Chien, and Haryadi S. Gunawi. 2016. The Tail at Store: A Revelation from Millions of Hours of Disk and SSD Deployments. In *14th USENIX Conference on File and Storage Technologies (FAST 16)*. USENIX Association, Santa Clara, CA, 263–276. <https://www.usenix.org/conference/fast16/technical-sessions/presentation/hao>
- [36] Christoph Hellwig. 2009. XFS: The Big Storage File System for Linux. *USENIX ;login issue* 34, 5 (2009).
- [37] J. Howard, M. Kazar, S. Menees, D. Nichols, M. Satyanarayanan, Robert N. Sidebotham, and M. West. 1987. Scale and Performance in a Distributed File System. In *Proceedings of the Eleventh ACM Symposium on Operating Systems Principles (SOSP '87)*. ACM, New York, NY, USA, 1–2. <https://doi.org/10.1145/41457.37500>
- [38] Joel Hruska. 2019. Western Digital to Demo Dual Actuator HDD, Will Use SMR to Hit 18TB Capacity. <https://www.extremetech.com/computing/287319-western-digital-to-demo-dual-actuator-hdd-will-use-smr-to-hit-18tb-capacity>.
- [39] Yige Hu, Zhiting Zhu, Ian Neal, Youngjin Kwon, Tianyu Cheng, Vijay Chidambaram, and Emmett Witchel. 2018. TxFS: Leveraging File-System Crash Consistency to Provide ACID Transactions. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. USENIX Association, Boston, MA, 879–891. <https://www.usenix.org/conference/atc18/presentation/hu>
- [40] Cheng Huang, Huseyin Simitci, Yikang Xu, Aaron Ogus, Brad Calder, Parikshit Gopalan, Jin Li, and Sergey Yekhanin. 2012. Erasure Coding in Windows Azure Storage. In *Presented as part of the 2012 USENIX Annual Technical Conference (USENIX ATC 12)*. USENIX, Boston, MA, 15–26. <https://www.usenix.org/conference/atc12/technical-sessions/presentation/huang>

- [41] Felix Hupfeld, Toni Cortes, Björn Kolbeck, Jan Stender, Erich Focht, Matthias Hess, Jesus Malo, Jonathan Marti, and Eugenio Cesario. 2008. The XtremFS Architecture – a Case for Object-based File Systems in Grids. *Concurrency and Computation: Practice and Experience* 20, 17 (Dec. 2008), 2049–2060. <https://doi.org/10.1002/cpe.v20:17>
- [42] Facebook Inc. 2019. RocksDB Direct IO. <https://github.com/facebook/rocksdb/wiki/Direct-IO>.
- [43] Facebook Inc. 2019. RocksDB Merge Operator. <https://github.com/facebook/rocksdb/wiki/Merge-Operator>.
- [44] Facebook Inc. 2019. RocksDB Synchronous Writes. <https://github.com/facebook/rocksdb/wiki/Basic-Operations#synchronous-writes>.
- [45] Silicon Graphics Inc. 2006. XFS Allocation Groups. http://xfs.org/docs/xfsdocs-xml-dev/XFS_FileSystem_Structure/tmp/en-US/html/Allocation_Groups.html.
- [46] INCITS T10 Technical Committee. 2014. *Information technology - Zoned Block Commands (ZBC)*. Draft Standard T10/BSR INCITS 536. American National Standards Institute, Inc. Available from <http://www.t10.org/drafts.htm>.
- [47] William Jannen, Jun Yuan, Yang Zhan, Amogh Akshintala, John Esmet, Yizheng Jiao, Ankur Mittal, Prashant Pandey, Phaneendra Reddy, Leif Walsh, Michael A. Bender, Martin Farach-Colton, Rob Johnson, Bradley C. Kuzmaul, and Donald E. Porter. 2015. BetrFS: Write-Optimization in a Kernel File System. *Trans. Storage* 11, 4, Article 18 (Nov. 2015), 29 pages. <https://doi.org/10.1145/2798729>
- [48] Sooman Jeong, Kisung Lee, Seongjin Lee, Seoungbum Son, and Youjip Won. 2013. I/O Stack Optimization for Smartphones. In *Presented as part of the 2013 USENIX Annual Technical Conference (USENIX ATC 13)*. USENIX, San Jose, CA, 309–320. <https://www.usenix.org/conference/atc13/technical-sessions/presentation/jeong>
- [49] Theodore Johnson and Dennis Shasha. 1994. 2Q: A Low Overhead High Performance Buffer Management Replacement Algorithm. In *Proceedings of the 20th International Conference on Very Large Data Bases (VLDB '94)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 439–450. <http://dl.acm.org/citation.cfm?id=645920.672996>
- [50] M. Frans Kaashoek, Dawson R. Engler, Gregory R. Ganger, Hector M. Briceño, Russell Hunt, David Mazières, Thomas Pinckney, Robert Grimm, John Jannotti, and Kenneth Mackenzie. 1997. Application Performance and Flexibility on Exokernel Systems. In *Proceedings of the Sixteenth ACM Symposium on Operating Systems Principles (SOSP '97)*. ACM, New York, NY, USA, 52–65. <https://doi.org/10.1145/268998.266644>
- [51] Yangwook Kang, Rekha Pitchumani, Pratik Mishra, Yang-suk Kee, Francisco Londono, Sangyoon Oh, Jongyeol Lee, and Daniel D. G. Lee. 2019. Towards Building a High-performance, Scale-in Key-value Storage System. In *Proceedings of the 12th ACM International Conference on Systems and Storage (SYSTOR '19)*. ACM, New York, NY, USA, 144–154. <https://doi.org/10.1145/3319647.3325831>
- [52] John Kennedy and Michael Satran. 2018. About Transactional NTFS. <https://docs.microsoft.com/en-us/windows/desktop/fileio/about-transactional-ntfs>.
- [53] John Kennedy and Michael Satran. 2018. Alternatives to using Transactional NTFS. <https://docs.microsoft.com/en-us/windows/desktop/fileio/deprecation-of-txf>.
- [54] Jaeho Kim, Donghee Lee, and Sam H. Noh. 2015. Towards SLO Complying SSDs Through OPS Isolation. In *13th USENIX Conference on File and Storage Technologies (FAST 15)*. USENIX Association, Santa Clara, CA, 183–189. https://www.usenix.org/conference/fast15/technical-sessions/presentation/kim_jaeho
- [55] Butler Lampson and Howard E Sturgis. 1979. Crash recovery in a distributed data storage system. (1979).
- [56] Adam Leventhal. 2016. APFS in Detail: Overview. <http://dtrace.org/blogs/ahl/2016/06/19/apfs-part1/>.
- [57] Peter Macko, Xiongzi Ge, John Haskins Jr., James Kelley, David Slik, Keith A. Smith, and Maxim G. Smith. 2017. SMORE: A Cold Data Object Store for SMR Drives (Extended Version). *CoRR* abs/1705.09701 (2017). <http://arxiv.org/abs/1705.09701>
- [58] Magic Pocket & Hardware Engineering Teams. 2018. Extending Magic Pocket Innovation with the first petabyte scale SMR drive deployment. <https://blogs.dropbox.com/tech/2018/06/extending-magic-pocket-innovation-with-the-first-petabyte-scale-smr-drive-deployment/>.
- [59] Magic Pocket & Hardware Engineering Teams. 2019. SMR: What we learned in our first year. <https://blogs.dropbox.com/tech/2019/07/smr-what-we-learned-in-our-first-year/>.
- [60] Lars Marowsky-Brée. 2018. Ceph User Survey 2018 results. <https://ceph.com/ceph-blog/ceph-user-survey-2018-results/>.
- [61] Marshall K McKusick, William N Joy, Samuel J Leffler, and Robert S Fabry. 1984. A Fast File System for UNIX. *ACM Transactions on Computer Systems (TOCS)* 2, 3 (1984), 181–197.
- [62] Chris Mellor. 2019. Toshiba embraces shingling for next-gen MAMR HDDs. <https://blocksandfiles.com/2019/03/11/toshiba-mamr-statements-have-shingling-absence/>.
- [63] Changwoo Min, Woon-Hak Kang, Taesoo Kim, Sang-Won Lee, and Young Ik Eom. 2015. Lightweight Application-Level Crash Consistency on Transactional Flash Storage. In *2015 USENIX Annual Technical Conference (USENIX ATC 15)*. USENIX Association, Santa Clara, CA, 221–234. <https://www.usenix.org/conference/atc15/technical-session/presentation/min>
- [64] Sumedh N. 2013. Coding for Performance: Data alignment and structures. <https://software.intel.com/en-us/articles/coding-for-performance-data-alignment-and-structures>.
- [65] Michael A. Olson. 1993. The Design and Implementation of the Inversion File System. In *USENIX Winter*.
- [66] Michael A. Olson, Keith Bostic, and Margo Seltzer. 1999. Berkeley DB. In *Proceedings of the Annual Conference on USENIX Annual Technical Conference (ATEC '99)*. USENIX Association, Berkeley, CA, USA, 43–43. <http://dl.acm.org/citation.cfm?id=1268708.1268751>
- [67] OpenStack Foundation. 2017. 2017 Annual Report. <https://www.openstack.org/assets/reports/OpenStack-AnnualReport2017.pdf>.
- [68] Adrian Palmer. 2015. SMRFFS-EXT4—SMR Friendly File System. https://github.com/Seagate/SMR_FS-EXT4.
- [69] Swapnil Patil and Garth Gibson. 2011. Scale and Concurrency of GIGA+: File System Directories with Millions of Files. In *Proceedings of the 9th USENIX Conference on File and Storage Technologies (FAST'11)*. USENIX Association, Berkeley, CA, USA, 13–13. <http://dl.acm.org/citation.cfm?id=1960475.1960488>
- [70] Juan Piernas. 2002. DualFS: A New Journaling File System without Meta-data Duplication. In *In Proceedings of the 16th International Conference on Supercomputing*, 137–146.
- [71] Poornima G and Rajesh Joseph. 2016. Metadata Performance Bottlenecks in Gluster. <https://www.slideshare.net/GlusterCommunity/performance-bottlenecks-for-metadata-workload-in-gluster-with-poornima-gurusiddaiah-rajesh-joseph>.
- [72] Donald E. Porter, Owen S. Hofmann, Christopher J. Rossbach, Alexander Benn, and Emmett Witchel. 2009. Operating System Transactions. In *Proceedings of the ACM SIGOPS 22Nd Symposium on Operating Systems Principles (SOSP '09)*. ACM, New York, NY, USA, 161–176. <https://doi.org/10.1145/1629575.1629591>
- [73] Lee Prewitt. 2019. SMR and ZNS – Two Sides of the Same Coin. <https://www.youtube.com/watch?v=jBxzO6YyMxU>.
- [74] Red Hat Inc. 2019. GlusterFS Architecture. <https://docs.gluster.org/en/latest/Quick-Start-Guide/Architecture/>.
- [75] Kai Ren and Garth Gibson. 2013. TABLEFS: Enhancing Metadata Efficiency in the Local File System. In *Presented as part of the 2013 USENIX Annual Technical Conference (USENIX ATC 13)*. USENIX, San Jose, CA, USA, 145–156. <https://www.usenix.org/conference/atc13/technical-sessions/presentation/ren>

- [76] Mendel Rosenblum and John K. Ousterhout. 1991. The Design and Implementation of a Log-structured File System. In *Proceedings of the Thirteenth ACM Symposium on Operating Systems Principles (SOSP '91)*. ACM, New York, NY, USA, 1–15. <https://doi.org/10.1145/121132.121137>
- [77] Frank Schmuck and Jim Wylie. 1991. Experience with Transactions in QuickSilver. In *Proceedings of the Thirteenth ACM Symposium on Operating Systems Principles (SOSP '91)*. ACM, New York, NY, USA, 239–253. <https://doi.org/10.1145/121132.121171>
- [78] Thomas J. E. Schwarz, Qin Xin, Ethan L. Miller, Darrell D. E. Long, Andy Hospodor, and Spencer Ng. 2004. Disk Scrubbing in Large Archival Storage Systems. In *Proceedings of the The IEEE Computer Society's 12th Annual International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunications Systems (MAS-COTS '04)*. IEEE Computer Society, Washington, DC, USA, 409–418. <http://dl.acm.org/citation.cfm?id=1032659.1034226>
- [79] Seastar. 2019. Shared-nothing Design. <http://seastar.io/shared-nothing/>.
- [80] Margo I. Seltzer. 1993. Transaction Support in a Log-Structured File System. In *Proceedings of the Ninth International Conference on Data Engineering*. IEEE Computer Society, Washington, DC, USA, 503–510. <http://dl.acm.org/citation.cfm?id=645478.654970>
- [81] Kai Shen, Stan Park, and Men Zhu. 2014. Journaling of Journal Is (Almost) Free. In *Proceedings of the 12th USENIX Conference on File and Storage Technologies (FAST 14)*. USENIX, Santa Clara, CA, 287–293. <https://www.usenix.org/conference/fast14/technical-sessions/presentation/shen>
- [82] Anton Shilov. 2017. Seagate Ships 35th Millionth SMR HDD, Confirms HAMR-Based Drives in Late 2018. <https://www.anandtech.com/show/11315/seagate-ships-35th-millionth-smr-hdd-confirms-hamr-based-hard-drives-in-late-2018>.
- [83] A. Shilov. 2019. Western Digital: Over Half of Data Center HDDs Will Use SMR by 2023. <https://www.anandtech.com/show/14099/western-digital-over-half-of-dc-hdds-will-use-smr-by-2023>.
- [84] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. 2010. The Hadoop Distributed File System. In *Proceedings of the 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST) (MSST '10)*. IEEE Computer Society, Washington, DC, USA, 1–10. <https://doi.org/10.1109/MSST.2010.5496972>
- [85] Chris Siebenmann. 2011. About the order that readdir() returns entries in. <https://utcc.utoronto.ca/~cks/space/blog/unix/ReaddirOrder>.
- [86] Chris Siebenmann. 2013. ZFS transaction groups and the ZFS Intent Log. <https://utcc.utoronto.ca/~cks/space/blog/solaris/ZFSTXGsAndZILs>.
- [87] Richard P. Spillane, Sachin Gaikwad, Manjunath Chinni, Erez Zadok, and Charles P. Wright. 2009. Enabling Transactional File Access via Lightweight Kernel Extensions. In *7th USENIX Conference on File and Storage Technologies (FAST 09)*. USENIX Association, San Francisco, CA. <https://www.usenix.org/conference/fast-09/enabling-transactional-file-access-lightweight-kernel-extensions>
- [88] Stas Starikov. 2016. [ceph-users] RadosGW performance degradation on the 18 millions objects stored. <http://lists.ceph.com/pipermail/ceph-users-ceph.com/2016-September/012983.html>.
- [89] Jan Stender, Björn Kolbeck, Mikael Höggqvist, and Felix Hupfeld. 2010. BabuDB: Fast and Efficient File System Metadata Storage. In *Proceedings of the 2010 International Workshop on Storage Network Architecture and Parallel I/Os (SNAPI '10)*. IEEE Computer Society, Washington, DC, USA, 51–58. <https://doi.org/10.1109/SNAPI.2010.14>
- [90] Michael Stonebraker. 1981. Operating System Support for Database Management. *Communications of the ACM* 24, 7 (July 1981), 412–418. <https://doi.org/10.1145/358699.358703>
- [91] Michael Stonebraker and Lawrence A. Rowe. 1986. The Design of POSTGRES. In *Proceedings of the 1986 ACM SIGMOD International Conference on Management of Data (SIGMOD '86)*. ACM, New York, NY, USA, 340–355. <https://doi.org/10.1145/16894.16888>
- [92] ZAR team. 2019. "Write hole" phenomenon. <http://www.raid-recovery-guide.com/raid5-write-hole.aspx>.
- [93] ThinkParQ. 2018. An introduction to BeeGFS. https://www.beegfs.io/docs/whitepapers/Introduction_to_BeeGFS_by_ThinkParQ.pdf.
- [94] Stephen C Tweedie. 1998. Journaling the Linux ext2fs Filesystem. In *The Fourth Annual Linux Expo*. Durham, NC, USA.
- [95] Sage Weil. 2009. Re: [RFC] big fat transaction ioctl. <https://lwn.net/Articles/361472/>.
- [96] Sage Weil. 2009. [RFC] big fat transaction ioctl. <https://lwn.net/Articles/361439/>.
- [97] Sage Weil. 2011. [PATCH v3] introduce sys_syncfs to sync a single file system. <https://lwn.net/Articles/433384/>.
- [98] 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 (OSDI '06)*. USENIX Association, Berkeley, CA, USA, 307–320. <http://dl.acm.org/citation.cfm?id=1298455.1298485>
- [99] Sage A. Weil, Scott A. Brandt, Ethan L. Miller, and Carlos Maltzahn. 2006. CRUSH: Controlled, Scalable, Decentralized Placement of Replicated Data. In *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing (SC '06)*. ACM, New York, NY, USA, Article 122. <https://doi.org/10.1145/1188455.1188582>
- [100] Sage A. Weil, Andrew W. Leung, Scott A. Brandt, and Carlos Maltzahn. 2007. RADOS: A Scalable, Reliable Storage Service for Petabyte-scale Storage Clusters. In *Proceedings of the 2Nd International Workshop on Petascale Data Storage: Held in Conjunction with Supercomputing '07 (PDSW '07)*. ACM, New York, NY, USA, 35–44. <https://doi.org/10.1145/1374596.1374606>
- [101] Brent Welch, Marc Unangst, Zainul Abbasi, Garth Gibson, Brian Mueller, Jason Small, Jim Zelenka, and Bin Zhou. 2008. Scalable Performance of the Panasas Parallel File System. In *Proceedings of the 6th USENIX Conference on File and Storage Technologies (FAST'08)*. USENIX Association, Berkeley, CA, USA, Article 2, 17 pages. <http://dl.acm.org/citation.cfm?id=1364813.1364815>
- [102] Lustre Wiki. 2017. Introduction to Lustre Architecture. <http://wiki.lustre.org/images/6/64/LustreArchitecture-v4.pdf>.
- [103] Wikipedia. 2018. Btrfs History. <https://en.wikipedia.org/wiki/Btrfs#History>.
- [104] Wikipedia. 2018. XFS History. <https://en.wikipedia.org/wiki/XFS#History>.
- [105] Wikipedia. 2019. Cache flushing. https://en.wikipedia.org/wiki/Disk_buffer#Cache_flushing.
- [106] Charles P. Wright, Richard Spillane, Gopalan Sivathanu, and Erez Zadok. 2007. Extending ACID Semantics to the File System. *Trans. Storage* 3, 2, Article 4 (June 2007). <https://doi.org/10.1145/1242520.1242521>
- [107] Fengguang Wu. 2012. I/O-less Dirty Throttling. https://events.linuxfoundation.org/images/stories/pdf/lcjp2012_wu.pdf.
- [108] Shiqin Yan, Huaicheng Li, Mingzhe Hao, Michael Hao Tong, Swaminathan Sundararaman, Andrew A. Chien, and Haryadi S. Gunawi. 2017. Tiny-Tail Flash: Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs. In *15th USENIX Conference on File and Storage Technologies (FAST 17)*. USENIX Association, Santa Clara, CA, 15–28. <https://www.usenix.org/conference/fast17/technical-sessions/presentation/yan>
- [109] Lawrence Ying and Theodore Ts'o. 2017. Dynamic Hybrid-SMR: an OCP proposal to improve data center disk drives. <https://www.blog.google/products/google-cloud/dynamic-hybrid-smr-ocp-proposal-improve-data-center-disk-drives/>.
- [110] Zhihui Zhang and Kanad Ghose. 2007. hFS: A Hybrid File System Prototype for Improving Small File and Metadata Performance. In *Proceedings of the 2Nd ACM SIGOPS/EuroSys European Conference*

- on *Computer Systems 2007 (EuroSys '07)*. ACM, New York, NY, USA, 175–187. <https://doi.org/10.1145/1272996.1273016>
- [111] Qing Zheng, Charles D. Cranor, Danhao Guo, Gregory R. Ganger, George Amvrosiadis, Garth A. Gibson, Bradley W. Settlemyer, Gary Grider, and Fan Guo. 2018. Scaling Embedded In-situ Indexing with deltaFS. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis (SC '18)*. IEEE Press, Piscataway, NJ, USA, Article 3, 15 pages. <http://dl.acm.org/citation.cfm?id=3291656.3291660>
- [112] Alexey Zhuravlev. 2016. ZFS: Metadata Performance. https://www.eofs.eu/_media/events/lad16/02_zfs_md_performance_improvements_zhuravlev.pdf.