

# Cudele: An API and Framework for Programmable Consistency and Durability in a Global Namespace

Michael A. Sevilla, Ivo Jimenez, Noah Watkins, Jeff LeFevre  
Peter Alvaro, Shel Finkelstein, Patrick Donnelley\*, Carlos Maltzahn

University of California, Santa Cruz; \*Red Hat

{msevilla, ivo, jayhawk, jlefevre}@soe.ucsc.edu, {palvaro, shel, carlosm}@ucsc.edu, pdonnell@redhat.com

**Abstract**—HPC and data center scale developers are abandoning POSIX IO because the file system metadata synchronization and serialization overheads of providing strong consistency and durability are too costly – and often unnecessary – for their applications. Unfortunately, designing file systems with weaker consistency or durability excludes applications that rely on stronger guarantees, forcing developers to re-write their applications or deploy them on a different system. Mounting multiple systems in the global namespace forces users to (1) provision separate storage clusters and (2) manually move data across systems. We present a framework and API that lets clients specify their consistency/durability requirements and dynamically assign them to subtrees in the same namespace, allowing users to optimize subtrees over time and space for different workloads. We confirm the performance benefits of techniques presented in related work but also explore new consistency/durability metadata designs, all integrated the same storage system. By custom fitting a subtree to a create-heavy application, we show 8× speedup and 2× scale when compared to our baseline system.

## I. INTRODUCTION

File system metadata services in HPC and data center applications have scalability problems because administrative tasks, like checkpointing [1] or scanning the file system [2], contend for the same directories and inodes. Applications perform better with dedicated metadata servers [3], [4] but provisioning a metadata server for every client is unreasonable. This problem is exacerbated by current trends; for example, HPC architectures are transitioning from complex storage stacks with burst buffer, file system, object store, and tape tiers to more simplified stacks with just a burst buffer and object store [5]. These types of trends put more pressure on data access because more requests end up hitting the same layer and latencies cannot be hidden while data migrates across tiers.

To address this, developers are relaxing the consistency and durability semantics in the file system because weaker guarantees are sufficient for their applications. For example, many batch style jobs do not need the strong consistency that the file system provides, so BatchFS [2] and DeltaFS [6] do more client-side processing and merge updates when the job is done. Developers in these domains are turning to these non-POSIX IO solutions because their applications are well-understood (e.g., well-defined read/write phases, synchronization only needed during certain phases, workflows describing computation, etc.) and because these applications wreak havoc on file systems designed for general-purpose workloads (e.g., checkpoint-restart’s N:N and N:1 create patterns [1]).

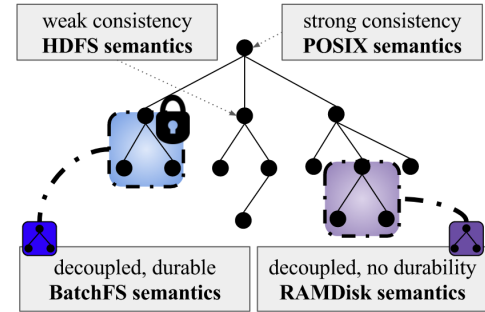


Fig. 1: Illustration of subtrees with different semantics co-existing in a global namespace. For performance, clients can relax consistency on their subtree (HDFS) or decouple the subtree and move it locally (BatchFS, RAMDisk). Decoupled subtrees can relax durability for even better performance.

One popular approach for relaxing consistency and durability is to “decouple the namespace”, where clients lock the subtree they want exclusive access to as a way to tell the file system that the subtree is important or may cause resource contention in the near-future [2], [4], [6]–[8]. Then the file system can change its internal structure to optimize performance. For example, the file system could enter a mode that prevents other clients from interfering with the decoupled directory. This delayed merge (*i.e.* a form of eventual consistency) and relaxed durability improves performance and scalability by avoiding the costs of remote procedure calls (RPCs), synchronization, false sharing, and serialization. While the performance benefits of decoupling the namespace are obvious, applications that rely on the file system’s guarantees must be deployed on an entirely different system or re-written to coordinate strong consistency/durability themselves.

To address this problem, we present an API and framework that lets developers dynamically control the consistency and durability guarantees for subtrees in the file system namespace. Figure 1 shows a potential setup in our proposed system where a single global namespace has subtrees for applications optimized with techniques from different state-of-the-art architectures. The HDFS<sup>1</sup> subtree has weaker than strong consistency because for performance it lets clients read files opened for writing [9], which means that not all updates are

<sup>1</sup>HDFS itself is not directly evaluated in this paper, although the semantics and their performance is explored in §V-B.

immediately seen by all servers; the BatchFS and RAMDisk subtrees are decoupled from the global namespace and have similar consistency/durability behavior to those systems; and the POSIX IO subtree retains the rigidity of POSIX IO’s strong consistency. Subtrees without policies inherit the consistency/durability semantics of the parent and future work will examine embeddable or inheritable policies.

Our prototype system, Cudele, achieves this by exposing “mechanisms” that developers use to specify their preferred semantics. Cudele supports 3 forms of consistency (invisible, weak, and strong) and 3 degrees of durability (none, local, and global) giving the user a wide range of policies and optimizations that can be custom fit to an application. We make the following contributions:

- 1) A framework and API for assigning consistency/durability policies to subtrees in the file system namespace; this lets users navigate the trade-offs of different metadata protocols on the same storage system.
- 2) This framework lets subtrees with different semantics co-exist in a global namespace. We show how this scales further and performs better than systems that use one strategy for the entire namespace .
- 3) A prototype that lets users custom fit subtrees to applications dynamically.

The last contribution lays the groundwork for future work on our prototype. It is better than the current practice of mounting different storage systems in a global namespace because there is no need to provision dedicated storage clusters to applications or move data between these systems. For example, the results of a Hadoop job do not need to be migrated into CephFS for other processing; instead the user can change the semantics of the HDFS subtree into a CephFS subtree. This may cause metadata/data movement to make things strongly consistent again but this is superior to moving all data across file system boundaries. We enable studies that adjust these semantics over *time and space*, where subtrees can change their semantics and migrate around the cluster without ever moving the data they reference.

The results in this paper lay the groundwork for such a system and confirm the assertions of “clean-state” research that decouple namespaces; specifically that the technique drastically improves performance ( $104\times$  speed up). We go a step further by quantifying the costs of merging updates ( $7\times$  slow down) and maintaining durability ( $10\times$  slow down). In our prototype, we get an  $8\times$  speedup and can scale to twice as many clients when we assign a more relaxed form of consistency and durability to a subtree with a create-heavy workload. We use Ceph as a prototyping platform because it is used in cloud-based and data center systems and has a presence in HPC [10].

## II. POSIX IO OVERHEADS

In our examination of the overheads of POSIX IO we benchmark and analyze CephFS, the file system that uses Ceph’s object store (*i.e.* RADOS) to store its data and metadata.

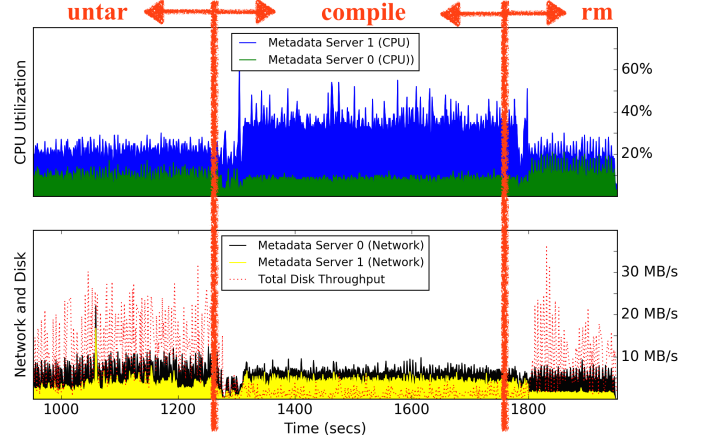


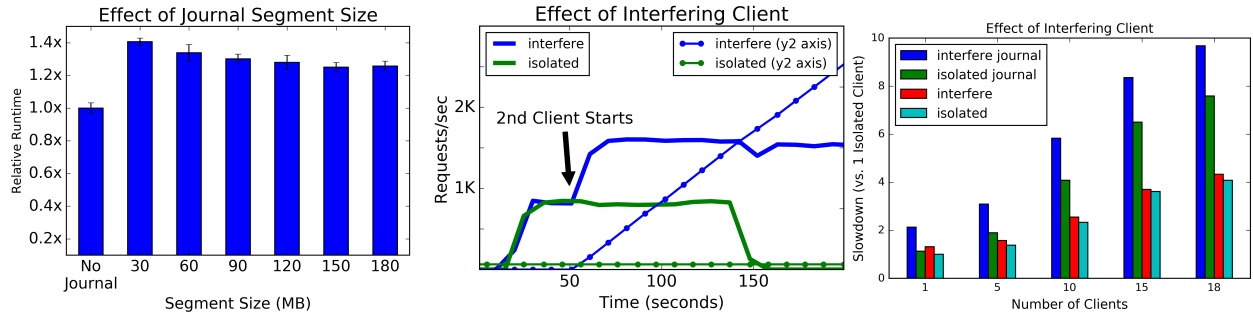
Fig. 2: Create-heavy workloads (such as `untar`) incur the highest disk, network, and CPU utilization because of the consistency and durability demands of CephFS.

To show how the file system behaves under high metadata load we use a create-heavy workload. During this process we discovered, based on the analysis and breakdown of costs, that durability and consistency have high overhead but we urge the reader to keep in mind that this file system is an implementation of one set of design decisions and our goal here is to highlight the effect that those decisions have on performance.

Figure 2 shows the resource utilization of compiling the Linux kernel. The `untar` phase, which is characterized by many creates, has the highest resource usage (combined CPU, network, and disk) because of the number of RPCs needed for consistency and durability. The RPCs are also serialized because the metadata server is single threaded; although naïve, this design is common because of the complexity of multi-threaded metadata servers [11]. Traditional file system techniques for improving performance, such as caching inodes, do not help for create-heavy workloads. In this section, we quantify the costs of consistency and durability in CephFS. At the end of each subsection we compare the approach to “decoupled namespaces”, the technique in related work that detaches subtrees from the global namespace to relax consistency/durability guarantees.

### A. Durability

While durability is not specified by POSIX IO, users expect that files they create or modify survive failures. We define three types of durability: global, local, and none. Global durability means that the client or server can fail at any time and metadata will not be lost because it is “safe” (*i.e.* striped or replicated across a cluster). Local durability means that metadata can be lost if the client or server stays down after a failure. None means that metadata is volatile and that the system provides no guarantees when clients or servers fail. None is different than local durability because regardless of the type failure, metadata will be lost when components die in a None configuration.



(a) Runtime normalized to no journal (b) Interference forces `lookup()`s (c) Interference hurts variability  
 Fig. 3: The overhead of durability and strong consistency in CephFS. (a) shows the effect of different journal segment sizes, which are streamed into the object store for fault tolerance. (b) and (c) show that when a second client “interferes”, capabilities are revoked and metadata servers do more work.

**CephFS Design:** A journal of metadata updates that streams into the resilient object store. Similar to LFS [12] and WAFL [13] the metadata journal is designed to be large (on the order of MBs) which ensures (1) sequential writes into the object store and (2) the ability for daemons to trim redundant or irrelevant journal entries. The journal is striped over objects where multiple journal updates can reside on the same object. There are two tunables for controlling the journal: the segment size and the number of segments that can be written in parallel. Unless the journal saturates memory or CPU resources, larger values for these tunables results in better performance.

In addition to the metadata journal, CephFS also represents metadata in RADOS as a metadata store, where directories and their file inodes are stored as objects. The metadata server applies the updates in the journal to the metadata store when the journal reaches a certain size. The metadata store is optimized for recovery (*i.e.* reading) while the metadata journal is write-optimized.

Figure 3a shows the effect of journaling of different journal segment sizes; the larger the segment size the bigger that the writes into the object store are. The trade-off for better performance is memory consumption because larger segments take up more space for buffering. When journaling is on, the metadata server periodically stops serving requests to flush (*i.e.* apply journal updates) to the metadata store. The journal overhead is sufficient enough to slow down metadata throughput but not so much as to overwhelm the bandwidth of the object store. We measured our peak bandwidth to be 100MB/s, which is the speed of our network link.

**Comparison to decoupled namespaces:** In BatchFS and DeltaFS, to the best of our knowledge, when a client or server fails there is no recovery scheme. For BatchFS, if a client fails when it is writing to the local log-structured merged tree (implemented as an SSTable [14]) then those batched metadata operations are lost. For DeltaFS, if the client fails then, on restart, the computation does the work again – since the snapshots of the namespace are never globally consistent and there is no ground truth. On the server side, BatchFS and DeltaFS use IndexFS. IndexFS writes metadata to SSTables, which initially reside in memory but are later flushed to the

underlying distributed file system.

### B. Strong Consistency

Access to metadata in a POSIX IO-compliant file system is strongly consistent, so reads and writes to the same inode or directory are globally ordered. The synchronization and serialization machinery needed to ensure that all clients see the same state has high overhead.

**CephFS Design:** Capabilities keep metadata strongly consistent. To reduce the number of RPCs needed for consistency, clients can obtain capabilities for reading, creating and updating inodes, as well as caching reads, buffering writes, changing the file size, and performing lazy IO.

To keep track of the read caching and write buffering capabilities, the clients and metadata servers agree on the state of each inode using an inode cache. If a client has the directory inode cached it can do metadata writes (*e.g.*, create) with a single RPC. If the client is not caching the directory inode then it must do an extra RPC to determine if the file exists. Unless the client immediately reads all the inodes in the cache (*i.e.* `ls -alR`), the inode cache is less useful for create-heavy workloads.

The benefits of caching the directory inode when creating files is shown in Figure 3b. The colors show the behavior of the client for two different runs. If only one client is creating files in a directory (“isolated” curve on  $y_1$  axis) then that client can lookup the existence of new files locally before issuing a create request to the metadata server. If another client starts creating files in the same directory (“interfere” curve on  $y_1$  axis) then the directory inode transitions out of read caching and the first client must send `lookup()`s to the metadata server (“interfere” curve on  $y_2$  axis). These extra requests increase the throughput of the “interfere” curve because the metadata server can handle the extra load but performance suffers. This degradation is shown in Figure 3c, where we scale the number of clients and show the slowdown of the slowest client. The results are normalized to a single isolated client without a metadata server journal. For the “interfere” bars, each client creates files in private directories and at 30 seconds we launch another process that creates files in those directories. 18 clients is the maximum number the metadata

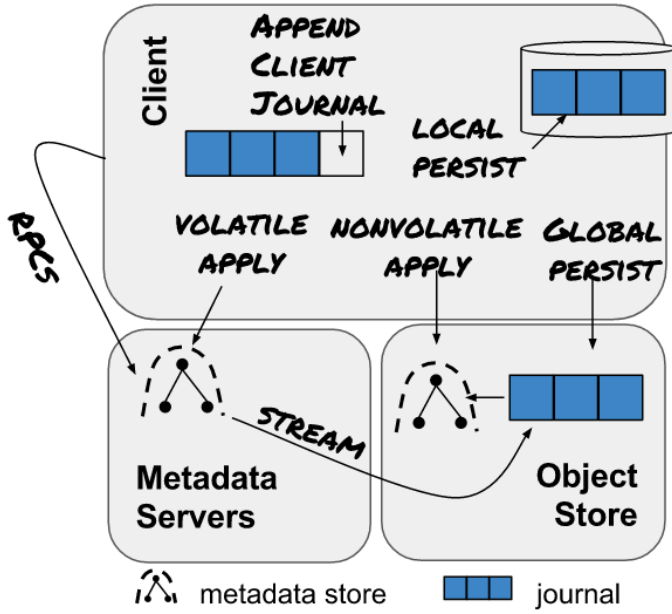


Fig. 4: Illustration of the **mechanisms** used by applications to build consistency/durability semantics. Descriptions are provided in Table I and the underlined words in Section §III-A.

| Mechanism             | Description                          |
|-----------------------|--------------------------------------|
| RPCs                  | round trip remote procedure calls    |
| Stream                | stream journal into object store     |
| Append Client Journal | events appended to in-memory journal |
| Volatile Apply        | apply to metadata store in memory    |
| Nonvolatile Apply     | apply to metadata store in obj store |
| Local Persist         | journal saved to client's disk       |
| Global Persist        | journal saved in object store        |

TABLE I: Descriptions of the mechanisms. Example compositions are shown in Table II.

server could handle for this metadata-intensive workload; at higher client load, the metadata server started complaining about laggy and unresponsive requests.

**Comparison to decoupled namespaces:** Decoupled namespaces merge batches of metadata operations into the global namespaces when the job completes. In BatchFS the merge is delayed by the application using an API to switch between asynchronous to synchronous mode. The merge itself is explicitly managed by the application but future work looks at more automated methodologies. In DeltaFS snapshots of the metadata subtrees stays on the client machines; there is no ground truth and consistent namespaces are constructed and resolved at application read time or when a 3rd party system (*e.g.*, middleware, scheduler, etc.) needs a view of the metadata. As a result, all the overheads of maintaining consistency that we showed above are delayed until the merge phase.

### III. METHODOLOGY: GLOBAL NAMESPACE, SUBTREE CONSISTENCY/DURABILITY

In this section we describe Cudele, our API and framework that lets users assign consistency and durability semantics to subtrees in the global namespace. A **mechanism** is an ab-

straction and basic building block for constructing consistency and durability guarantees. Cudele exposes these mechanisms and the user composes them together to construct **policies**. These policies are assigned to subtrees and they dictate how the file system should handle operations within that subtree. Below, we describe the mechanisms (which are underlined), the policies, and the API for assigning policies to subtrees.

#### A. Mechanisms: Building Guarantees

Figure 4 shows the mechanisms (labeled arrows) in Cudele and which daemon(s) they are performed by. Table I has a description of what each mechanism does. Decoupled clients use the Append Client Journal mechanism to append metadata updates to a local, in-memory journal. This is fast because clients do not need to check for consistency when writing events and the metadata server blindly applies the updates because it assumes the events were already checked for consistency. The trade-off here is performance; it is a dangerous approach but could be implemented safely if the clients or metadata server are configured to check the validity of events before writing them. Once the job is complete, the system calls Cudele mechanisms to achieve the desired consistency/durability semantics. Cudele provides a library for clients to link into and all operations are performed by the client.

1) *Mechanisms Used for Consistency:* RPCs send remote procedure calls for every metadata operation from the client to the metadata server, assuming the request cannot be satisfied by the inode cache. This mechanism is part of the default CephFS implementation and is the strongest form of consistency because clients see metadata updates right away. Nonvolatile Apply replays the client's journal onto the metadata cluster's metadata store. The client's in-memory journal is written into the object store and the metadata servers are restarted. When the metadata servers re-initialize, they notice new journal updates in the object store and replay the events onto their in-memory metadata stores. Volatile Apply takes the client's in-memory journal on the client and applies the updates directly to the in-memory metadata store maintained by the metadata servers. We say volatile because – in exchange for peak performance – Cudele makes no consistency or durability guarantees while Volatile Apply is executing. If a concurrent update from a client occurs there is no rule for resolving conflicts and if the client or metadata server crashes there may be no way to recover.

The biggest difference between Volatile Apply and Nonvolatile Apply is the medium they use to communicate. Volatile Apply applies updates directly to the metadata servers' metadata store while Nonvolatile Apply uses the object store to communicate the journal of updates from the client to the metadata servers. Nonvolatile Apply is safer but has a large performance overhead because objects in the metadata store need to be read from and written back to the object store.

2) *Mechanisms Used for Durability:* Stream is one of the mechanisms used by default in CephFS. Using existing configuration settings in Ceph we can turn Stream on and off.



| C →<br>D ↓ | invisible                                | weak  | strong                    |
|------------|--|---|---------------------------|
| none       | append client journal                    | append client journal<br>+volatile apply                    | RPCs                      |
| local      | append client journal<br>+local persist  | append client journal<br>+local persist<br>+volatile apply  | RPCs<br>+local<br>persist |
| global     | append client journal<br>+global persist | append client journal<br>+global persist<br>+volatile apply | RPCs<br>+stream           |

TABLE II: Users can explore the consistency (C) and durability (D) spectrum by composing Cudele mechanisms.

If it is off, then the metadata servers will not save journals in the object store. For Local Persist, clients write serialized log events to a file on local disk and for Global Persist, clients push the journal into the objects store. The overheads for both Local Persist and Global Persist is the write bandwidth of the local disk and object store, respectively. These persist mechanisms are part of the library that links into the client.

### B. Defining Policies in Cudele

The spectrum of consistency and durability guarantees that users can construct is shown in Table II. The columns are the different consistency semantics and the rows cover the spectrum of durability guarantees. For consistency: “invisible” means the system does not handle merging updates into a global namespace and it is assumed that middleware or the application manages consistency lazily; “weak” merges updates at some time in the future (e.g., when the system has time, when the number of updates reaches a certain threshold, when the client is done writing, etc.); and updates in “strong” consistency are seen immediately by all clients. For durability, “none” means that updates are volatile and will be lost on a failure. Stronger guarantees are made with “local”, which means updates will be retained if the client node recovers and reads the updates from local storage, and “global”, where all updates are always recoverable.

Existing, state-of-the-art systems in HPC can be represented by the cells in Table II (we construct the semantics for these systems later in Section §V-A). POSIX IO-compliant systems like CephFS and IndexFS have global consistency and durability<sup>2</sup>; DeltaFS uses “invisible” consistency and “local” durability and BatchFS uses “weak” consistency and “local” durability. These systems have other features that could push them into different semantics but we assign labels here based on the points emphasized in the papers. To compose the mechanisms users inject which mechanisms to run and which to use in parallel using a domain specific language. Although we can achieve all permutations of the different guarantees in Table II, not all of them make sense. For example, it makes little sense to do `append client journal+RPCs` since both mechanisms do the same thing or `stream+local persist` since “global” durability is stronger and has more

overhead than “local” durability. The cost of each mechanism and the semantics described above are quantified in Sections §V-A and §V-B.

The consistency and durability properties in Table II are not guaranteed until all mechanisms in the cell are complete. The compositions should be considered atomic and there are no guarantees while transitioning between policies. For example, updates are not deemed to have “global” durability until they are safely saved in the object store. If a failure occurs during Global Persist or if we inject a new policy that changes a subtree from Local Persist to Global Persist, Cudele makes no guarantee until the mechanisms are complete. Despite this, production systems that use Cudele should state up-front what the transition guarantees are for subtrees.

### C. Cudele Namespace API

Users control consistency and durability for subtrees by contacting a daemon in the system called a monitor, which manages cluster state changes. Users present a directory path and a policies configuration that gets distributed and versioned by the monitor to all daemons in the system. For example, (msevilla/mydir, policies.yml) would decouple the path “msevilla/mydir” and would apply the policies in “policies.yml”. The policies file supports the following parameters (default values are in parenthesis):

- `consistency`: which consistency model to use (RPCs)
- `durability`: which durability model to use (stream)
- `allocated_inodes`: the number of inodes to provision to the decoupled namespace (100)
- `interfere_policy`: how to handle a request from another client targeted at this subtree (allow)

The “Consistency” and “Durability” parameters are set to the Cudele mechanisms; they can be composed (+) or run in parallel (|||). “Allocated Inodes” is a way for the application to specify how many files it intends to create. It is a contract so that the file system can provision enough resources for the incumbent merge and so it can give valid inodes to other clients. The inodes can be used anywhere within the decoupled namespace (i.e. at any depth in the subtree).

“Interfere Policy” has two settings: `block` and `allow`. For `block`, any requests to this part of the namespace returns with “Device is busy”, which will spare the metadata server from wasting resources for updates that may get overwritten. If the application does not mind losing updates, for example it wants approximations for results that take too long to compute, it can select `allow`. In this case, metadata from the interfering client will be written and the computation from the decoupled namespace will take priority at merge time because the results are more accurate.

Given these default values decoupling the namespace with an empty policies file would give the application 100 inodes but the subtree would behave like the existing CephFS implementation. BatchFS merges updates into the global namespace after the job is complete. This can be achieved with by composing two mechanisms:

<sup>2</sup>This is the normal case. IndexFS also has bulk merge which would transition the system into “weak consistency”

## IV. IMPLEMENTATION

We use a programmable storage approach [15] to design Cudele; namely, we try to leverage components inside Ceph to inherit the robustness and correctness of the internal subsystems. Using this ‘dirty-slate’ approach, we only had to implement 4 of the 6 mechanisms from Figure 4 and just 1 required changes to the underlying storage system itself. In this section, we first describe a Ceph internal subsystem or component and then we show how we use it in Cudele.

**Metadata Store:** In CephFS, the metadata store is a data structure that represents the file system namespace. This data structure is stored in two places: in memory (*i.e.* in the collective memory of the metadata server cluster) and as objects in the object store. In the object store, directories and their inodes are stored together in objects to improve the performance of scans. The metadata store data structure is structured as a tree of directory fragments making it easier to read and traverse.

**Cudele:** the “RPCs” mechanism uses the metadata store to service requests. Using code designed for recovery, “Non-volatile Apply” and “Volatile Apply” replay updates onto the metadata store in memory and in the object store, respectively. When the clients are ready to merge their updates back into the global namespace, they pass a binary file of metadata updates to the metadata server.

**Journal:** The journal is the second way that CephFS represents the file system namespace; it is a log of metadata updates that can materialize the namespace when the updates are replayed onto the metadata store. The journal is a “pile system”; writes are fast but reads are slow because state must be reconstructed. Specifically, reads are slow because there is more state to read, it is unorganized, and many of the updates may be redundant.

**Cudele:** the journal format is used by “Stream”, “Append Client Journal”, “Local Persist”, and “Global Persist”. “Stream” is the default implementation for achieving global durability in CephFS but the rest of the mechanisms are implemented by writing with the journal format. By writing with the same format, the metadata servers can read and use the recovery code to materialize the updates from a client’s decoupled namespace (*i.e.* merge). These implementations required no changes to CephFS because the metadata servers know how to read the events the library is writing. By re-using the journal subsystem to implement the namespace decoupling, Cudele leverages the write/read optimized data structures, the formats for persisting events (similar to TableFS’s SSTables [14]), and the functions for replaying events onto the internal namespace data structures.

**Journal Tool:** The journal tool is used for disaster recovery and lets administrators view and modify the journal. It can read the journal, export the journal as a file, erase events, and apply updates to the metadata store. To apply journal updates to the metadata store, the journal tool reads the journal from object store objects and replays the updates on the metadata store in the object store.

**Cudele:** the external library the clients link into is based on the journal tool. It already had functions for importing, exporting, and modifying the updates in the journal so we re-purposed that code to implement “Append Client Journal”, “Volatile Apply”, and “Nonvolatile Apply”.

**Inode Cache:** In CephFS, the inode cache reduces the number of RPCs between clients and metadata servers. Without contention clients can resolve metadata reads locally thus reducing the number of operations (*e.g.*, `lookup()`s). For example, if a client or metadata server is not caching the directory inode, all creates within that directory will result in a lookup and a create request. If the directory inode is cached then only the create needs to be sent. The size of the inode cache is configurable so as not to saturate the memory on the metadata server – inodes in CephFS are about 1400 bytes<sup>3</sup>. The inode cache has code for manipulating inode numbers, such as pre-allocating inodes to clients.

**Cudele:** “Nonvolatile Apply” uses the internal inode cache code to allocate inodes to clients that decouple parts of the namespace and to skip inodes used by the client at merge time.

**Large Inodes:** In CephFS, inodes already store policies, like how the file is striped across the object store or for managing subtrees for load balancing. These policies adhere to logical partitionings of metadata or data, like Ceph pools and file system namespace subtrees. To implement this, the namespace data structure has the ability to recursively apply policies to subtrees and to isolate subtrees from each other.

**Cudele:** uses the large inodes to store consistency and durability policies in the directory inode. This approach uses the File Type interface from the Malacology programmable store system [15] and it tells clients how to access the underlying metadata. The underlying implementation stores executable code in the inode that calls the different Cudele mechanisms. Of course, there are many security and access control aspects of this approach but that is beyond the scope of this paper.

## V. EVALUATION

Cudele lets users construct consistency/durability guarantees using well-established research techniques from other systems; so instead of evaluating the scalability and performance of the techniques themselves against other file systems, we show that (1) the mechanisms we propose are useful for constructing semantics used by real systems, (2) the techniques can work side-by-side in the same namespace, and (3) the combination of these techniques can help the system scale further when subtrees are coupled to the correct type of application. We scope the evaluation to one metadata server and scale the number of parallel clients. This type of analysis shows the capacity and performance of a metadata server with superior metadata protocols, which should be used to inform load balancing across a metadata cluster. Load balancing across a cluster of metadata servers with partitioning and replication can be explored with something like Mantle [3].

<sup>3</sup>[http://docs.ceph.com/docs/jewel/dev/mds\\_internals/data-structures/](http://docs.ceph.com/docs/jewel/dev/mds_internals/data-structures/)

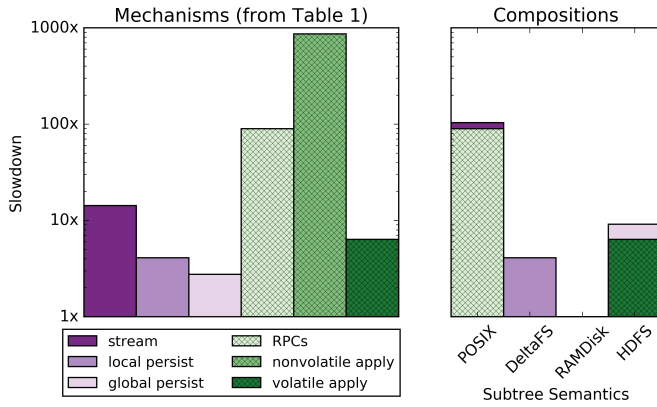


Fig. 5: Performance of each mechanism (left) and building the consistency/durability semantics of real-world systems (right) for 100K files creates from a single client. Results are normalized to the runtime of writing events to the client’s in-memory journal.

#### A. Microbenchmarks

Figure 5 shows the runtime of the Cudele mechanisms for a single client creating files in the same directory, normalized to the time it takes to write 100K file create updates to the client’s in-memory journal (*i.e.* the “Append Client Journal” mechanism). “Stream” is an approximation of the overhead and is calculated by subtracting the runtime of the job with the journal turned off from the runtime with the journal turned on. 100K is the maximum recommended size of a directory in CephFS; preliminary experiments with larger directory sizes show memory problems.

**Poorly Scaling Data Structures:** Despite doing the same amount of work, mechanisms that rely on poorly scaling data structures have large slowdowns for the larger number of creates. For example, “RPCs” which relies on the internal CephFS directory structures has a  $90\times$  slowdown for 100K files. It is a well-known problem that directory data structures do not scale when creating files in the same directory [4] and any mechanism that uses these data structures will experience similar slowdowns. Other mechanisms that write events to a journal experience a much less drastic slowdown because the journal data structure does not need to be scanned for every operation. Events are written to the end of the journal without even checking the validity (*e.g.*, if the file already exists for a create), which is another form of relaxed consistency because the file system assumes the application has resolved conflicting updates in a different way.

**Overhead of RPCs:** “RPCs” is  $66\times$  slower than “Volatile Apply” because sending individual metadata updates over the network is costly. While “RPCs” sends a request for every file create, “Volatile Apply” writes all the updates to the in-memory journal and applies them to the in-memory data structures in the metadata server. While communicating the decoupled namespace directly to the metadata server is faster, communicating through the object store (“Nonvolatile Apply”) is  $10\times$  slower.

**Overhead of “Nonvolatile Apply”:** The cost of “Nonvolatile apply” is much larger than all the other mechanisms. That mechanism was not implemented as part of Cudele – it was a debugging and recovery tool packaged with CephFS. It works by iterating over the updates in the journal and pulling all objects that *may* be affected by the update. This means that two objects are repeatedly pulled, updated, and pushed: the object that houses the experiment directory and the object that contains the root directory (*i.e.* /). The cost of communicating through the object store is shown by comparing the runtime of “Volatile apply” + “Global persist” to “Nonvolatile Apply”. These two operations end up with the same final metadata state but using “Nonvolatile Apply” is clearly inferior.

**Parallelism of the Object Store:** Comparing “Local” and “Global persist” demonstrates the bandwidth advantages of storing the journal in a distributed object store. The “Global Persist” performance is  $1.5\times$  faster because the object store is leveraging the collective bandwidth of the disks in the cluster. This benefit comes from the object store itself but should be acknowledged when making decisions for the application; the size of the object store can help mitigate the overheads of globally persisting metadata updates.

**Journal Size:** We measure the amount of storage per journal update to be about 2.5KB. The storage footprint scales linearly with the number of metadata creates and suggests that updates for a million files would be  $2.5\text{KB} * 1 \text{ million files} = 2.38\text{GB}$ . Transfer times for payloads of this size in most HPC/data center networks are reasonable.

**Takeaway:** measuring the mechanisms individually shows that their overheads and costs can differ *by orders of magnitude*. We also show that some mechanisms, like Nonvolatile Apply, are not worthwhile as currently implemented.

#### B. Use Case 1: Creates in the Same Directory

First we show how we can build application-specific subtrees by composing mechanisms and that this approach results in more a scalable global namespace. We start with clients creating files in private directories because this workload is heavily studied in HPC [2]–[4], [16], [17], mostly due to checkpoint-restart [1]. But the workload also appears in cloud workloads [18], where systems like Hadoop use the file abstraction to exchange work units to workers or to indicate when jobs complete [19]. A more familiar example is uncompressing an archive (*e.g.*, `tar xzf`), where the file system services a flash crowd of creates across all directories as shown in Figure 2. We use this as a microbenchmark because it allows us to control the size of the log, since each create creates a single journal event. The workload is clients creating 100K files in private directories in the same global namespace.

The graph on the right of Figure 5 shows how applications can compose mechanisms together to get the consistency/durability guarantees they need in a global namespace. We label the *x*-axis with systems that employ these semantics, as described in Figure 1. Again, the runtime is normalized to creating files in the client’s in-memory journal. We make no

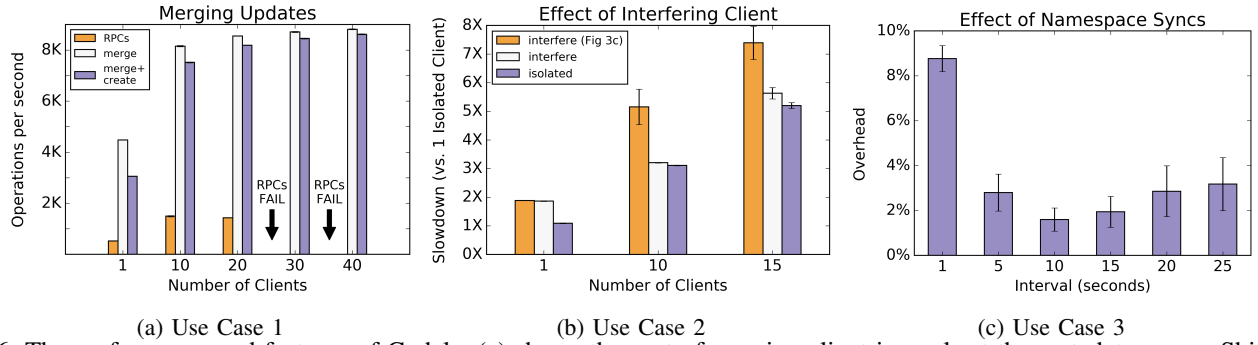


Fig. 6: The performance and features of Cudele. (a) shows the cost of merging client journals at the metadata server. Shipping and merging journals of updates scales better than RPCs because there are less messages and consistency/durability code paths are bypassed. In (b), the allow/block API isolates directories from interfering clients. (c) is the slowdown of a single client syncing updates to the global namespace. The inflection point is the trade-off of frequent updates vs. larger journal files.

guarantees during execution of the mechanisms or when transitioning semantics – the semantics are guaranteed *once the mechanism completes*. So if servers fail during a mechanism, metadata or data may be lost.

**Takeaway:** we confirm the performance benefits of other well-established research systems but, more importantly, we show that the Cudele mechanisms we propose are useful for building and evaluating different consistency/durability semantics.

To show the contention at the metadata server, in Figure 6a we scale the number of clients for the RPC per operation strategy used by many POSIX IO systems and the decoupled namespace strategy used by DeltaFS. For RPCs, we scale the number of clients making requests to the metadata server and for the decoupled bars, we scale the number of concurrent journal merge requests. So for the 30 clients data point, we are measuring the operations per second for 30 client journals that land on the metadata server at the same time. The “create & merge” bar adds the time it takes for clients to generate the journal of events in parallel.

This workload scales further with decoupled namespaces because the metadata server is not overwhelmed with RPC requests. Compared to Figure 3c we can scale to 40 clients which is more than double the capacity of a single metadata server using the RPC strategy to maintain strong consistency. Related work has pointed to on-disk metadata format as the primary factor for improved scalability; but our results show that the weakened consistency and durability enables the scalability improvements. The on-disk metadata formats in this experiment are the same for all schemes.

The performance of decoupled namespaces is better than RPCs because we place no restrictions on the validity of metadata inserted into the journal and avoid touching poorly scaling data structures. We can scale linearly if we weaken our consistency by never checking for the existence of files before creating files (*i.e.* only append `open()` requests to the journal). When the updates are merged by the metadata server into the in-memory metadata store they never scan growing data structures. Had we implemented the “Create” mechanism

to include `lookup()` commands before `open()` requests, we would have seen the poor scaling that we see with the “RPCs” mechanism.

**Takeaway:** the ability to couple well-established techniques to specific applications allows the global namespace to scale further and perform better. In our case, RPCs overwhelm the system but decoupling/shipping a journal of updates improves scalability by  $2\times$ .

### C. Use Case 2: Creates with Interfering Client

Next we show how Cudele can be programmed to block interfering clients using the same problematic workload from Figure 3b. In this workload, clients create files in private directories and a separate client creates files in each directory. This introduces false sharing and the metadata server revokes capabilities on directories touched by the interfering client. While HPC tries to avoid these situations with workflows [2], [6], it still happens in distributed file systems when users unintentionally access directories in a shared file system. We also have another client write to a decoupled namespace and merge its updates at 90 seconds. We only scale to 15 clients because the performance variability of a nearly overloaded metadata server, in our case 18 clients, is too high.

Figure 6b plots the overhead of the slowest client and the error bars are the standard deviation of 3 runs. “Interfere (Fig 3c)” is the result from Figure 3b; “Interfere” uses the allow/block API to return `-EBUSY` to interfering clients; and “Isolated” is the baseline performance without an interfering. Each subtree has “RPCs” and “Stream” enabled. Results results are normalized to the runtime of a single client creating files in a directory. Note that IndexFS has a similar behavior to “interfere” with leases, except clients block.

We draw three conclusions: (1) clients that use the API to block interfering clients get the same performance as isolated clients, (2) there is a negligible effect on performance for the extra work the metadata server does to return `-EBUSY`, and (3) merging updates from the decoupled client has a negligible effect on performance.



**Takeaway:** the API lets users isolate directories when applications need better and more reliable performance. Blocking updates is an effective way of controlling consistency.

#### D. Use Case 3: Read while Writing

The final use case shows how the API gives users fine-grained control of the consistency semantics for niche or *ad-hoc* workflows. Scientists use the file system to check the progress of jobs using `ls` even though this operation is notoriously heavy-weight [20], [21]. The number of files or size of the files is indicative of the progress. This practice is not too different from systems that use the file system to manage the progress of jobs; Hadoop writes to temporary files, renames them when complete, and creates a “DONE” file to indicate to the runtime that the task did not fail and should not be re-scheduled on another node. In this scenario, Cudele users will not see the progress of decoupled namespaces since their updates are not globally visible. To help scientists judge the progress of their jobs, Cudele clients have a “namespace sync” that sends batches of updates back to the global namespace at regular intervals.

Figure 6c shows the performance degradation of a single client writing updates to a decoupled namespace and pausing to send updates to the metadata server. We scale the namespace sync interval to show the trade-off of frequently pausing or writing large logs of updates. We use an idle core to log the updates and to do the network transfer. The client only pauses to fork off a background process, which is expensive as the address space needs to be copied. The alternative is to pause the client completely and write the update to disk but since this implementation is limited by the speed of the disk, we choose the memory-to-memory copy of the fork approach.

**Takeaway:** syncing namespace updates has up to a 9% overhead and can be tuned depending on the users preference but, more importantly, the API gives users fine-grain control of their consistency/durability to support current practices or experimental workflows.

#### E. Discussion and Future Work

Cudele and the experiments we present here prompt many avenues for future work. First is quantifying the benefits of dynamically changing the semantics of a subtree from stronger to weaker guarantees (or vice versa), as described in the introduction. Second is embeddable policies, where child subtrees have specialized features but still maintain the guarantees of their parent subtrees. We already show how applications can compose guarantees and control their performance, knobs that we did not have before, but embeddable policies are a way to compose the policies themselves. For example, a RAMDisk subtree is POSIX IO-compliant but relaxes durability constraints, so it can reside under a POSIX IO subtree that is strongly consistent. These embeddable policies should be controlled and enforced by Cudele. Finally, performance prediction benefits from our cost quantification of each mechanism. Applications can use Cudele to microbenchmark their components and software, similar to what we did in

Section §V-A. Using those results, they can predict how much slower their system will be if they adopt stronger consistency or durability. This is a form of co-design that takes a “dirty-slate” approach but building just the guarantees the application needs from existing implementations. This can also be a verification tool where performance that varies wildly from the predicted performance can be a red flag that something is wrong or that the bottleneck is not in the consistency or durability plane.

## VI. RELATED WORK

The bottlenecks associated with accessing POSIX IO file system metadata are not limited to HPC workloads and the same challenges that plagued these systems for years are finding their way into the cloud. Workloads that deal with many small files (*e.g.*, log processing and database queries [22]) and large numbers of simultaneous clients (*e.g.*, MapReduce jobs [23]), are subject to the scalability of the metadata service. The biggest challenge is that whenever a file is touched the client must access the file’s metadata and maintaining a file system namespace imposes small, frequent accesses on the underlying storage system [24]. Unfortunately, scaling file system metadata is a well-known problem and solutions for scaling data IO do not work for metadata IO [24]–[27].

POSIX IO workloads require strong consistency and many file systems improve performance by reducing the number of remote calls per operation (*i.e.* RPC amplification). As discussed in the previous section, caching with leases and replication are popular approaches to reducing the overheads of path traversals but their performance is subject to cache locality and the amount of available resources, respectively; for random workloads larger than the cache extra RPCs hurt performance [4], [16] and for write heavy workloads with more resources the RPCs for invalidations are harmful. Another approach to reducing RPCs is to use leases or capabilities.

High performance computing has unique requirements for file systems (*e.g.*, fast creates) and well-defined workloads (*e.g.*, workflows) that make relaxing POSIX IO sensible. BatchFS assumes the application coordinates accesses to the namespace, so the clients can batch local operations and merge with a global namespace image lazily. Similarly, DeltaFS eliminates RPC traffic using subtree snapshots for non-conflicting workloads and middleware for conflicting workloads. MarFS gives users the ability to lock “project directories” and allocate GPFS clusters for demanding metadata workloads. TwoTiers eliminates high-latencies by storing metadata in a flash tier; applications lock the namespace so that metadata can be accessed more quickly. Unfortunately, decoupling the namespace has costs: (1) merging metadata state back into the global namespace is slow; (2) failures are local to the failing node; and (3) the systems are not backwards compatible.

For (1), state-of-the-art systems manage consistency in non-traditional ways: IndexFS maintains the global namespace but blocks operations from other clients until the first client drops the lease, BatchFS does operations on a snapshot of the namespace and merges batches of operations into the global

namespace, and DeltaFS never merges back into the global namespace. The merging for BatchFS is done by an auxiliary metadata server running on the client and conflicts are resolved by the application. Although DeltaFS never explicitly merges, applications needing some degree of ground truth can either manage consistency themselves on a read or add a bolt-on service to manage the consistency. For (2), if the client fails and stays down, all metadata operations on the decoupled namespace are lost. If the client recovers, the on-disk structures (for BatchFS and DeltaFS this is the SSTables used in TableFS) can be recovered. In other words, the clients have state that cannot be recovered if the node stays failed and any progress will be lost. This scenario is a disaster for checkpoint-restart where missed cycles may cause the checkpoint to bleed over into computation time. For (3), decoupled namespace approaches sacrifice POSIX IO going as far as requiring the application to link against the systems they want to talk to. In today's world of software defined caching, this can be a problem for large data centers with many types and tiers of storage. Despite well-known performance problems POSIX IO and REST are the dominant APIs for data transfer.

## VII. CONCLUSION

Relaxing consistency/durability semantics in the file system is a double-edged sword. While it performs and scales better, it alienates applications that rely on strong consistency and durability. Mounting other systems to the global namespace is convenient but wastes resources and incurs data movement. Cudele lets users assign consistency/durability guarantees to subtrees in the global namespace, which can be custom fit to the application. Using Cudele, we show how applications can co-exist and perform well in a global namespace and lay the groundwork for new systems that dynamically change consistency/durability guarantees to avoid provisioning dedicated storage clusters and moving large amounts of data.

## REFERENCES

- [1] J. Bent, G. Gibson, G. Grider, B. McClelland, P. Nowoczynski, J. Nunez, M. Polte, and M. Wingate, "PLFS: a checkpoint filesystem for parallel applications," in *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, ser. SC '09.
- [2] Q. Zheng, K. Ren, and G. Gibson, "BatchFS: Scaling the File System Control Plane with Client-funded Metadata Servers," in *Proceedings of the Workshop on Parallel Data Storage*, ser. PDSW '14.
- [3] M. A. Sevilla, N. Watkins, C. Maltzahn, I. Nassi, S. A. Brandt, S. A. Weil, G. Farnum, and S. Fineberg, "Mantle: A Programmable Metadata Load Balancer for the Ceph File System," in *Proceedings of the Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC '15.
- [4] K. Ren, Q. Zheng, S. Patil, and G. Gibson, "IndexFS: Scaling File System Metadata Performance with Stateless Caching and Bulk Insertion," in *Proceedings of the Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC '14, 2014.
- [5] J. Bent, B. Settlemeyer, and G. Grider, "Serving Data to the Lunatic Fringe."
- [6] Q. Zheng, K. Ren, G. Gibson, B. W. Settlemeyer, and G. Grider, "DeltaFS: Exascale File Systems Scale Better Without Dedicated Servers," in *Proceedings of the Workshop on Parallel Data Storage*, ser. PDSW '15, 2015.
- [7] G. Grider, D. Montoya, H.-b. Chen, B. Kettering, J. Inman, C. De-Jager, A. Torrez, K. Lamb, C. Hoffman, D. Bonnie, R. Croonenberg, M. Broomfield, S. Leffler, P. Fields, J. Kuehn, and J. Bent, "MarFS - A Scalable Near-Posix Metadata File System with Cloud Based Object Backend," in *Work-in-Progress at Proceedings of the Workshop on Parallel Data Storage*, ser. PDSW '15, November 2015.
- [8] S. Faibish, J. Bent, U. Gupta, D. Ting, and P. Tzelnic, "Slides: 2 tier storage architecture." [Online]. Available: <http://www.pdl.cmu.edu/SDI/2015/slides/Faibish-CMU-PDL-Spring-2015-final.pdf>
- [9] J. D. Kamal Hakimzadeh, Hooman Peiro Sajjad, "Scaling HDFS with a Strongly Consistent Relational Model for Metadata," in *Proceedings of the International Conference on Distributed Applications and Interoperable Systems*, ser. DAIS '14.
- [10] F. Wang, M. Nelson, S. Oral, S. Atchley, S. Weil, B. W. Settlemeyer, B. Caldwell, and J. Hill, "Performance and Scalability Evaluation of the Ceph Parallel File System," in *Proceedings of the Workshop on Parallel Data Storage Workshop*, ser. PDSW '13.
- [11] K. Chasapis, M. F. Dolz, M. Kuhn, and T. Ludwig, "Evaluating Lustre's Metadata Server on a Multi-socket Platform," in *Proceedings of the 9th Parallel Data Storage Workshop*, ser. PDSW '14, 2014.
- [12] M. Rosenblum and J. Ousterhout, "The Design and Implementation of a Log-Structured File System," in *ACM Transactions on Computer Systems*, '92.
- [13] D. Hitz, J. Lau, and M. Malcolm, "File system design for an nfs file server appliance," in *Proceedings of the USENIX Technical Conference*, ser. WTEC '94.
- [14] K. Ren and G. Gibson, "TABLEFS: Enhancing Metadata Efficiency in the Local File System," in *Proceedings of the USENIX Annual Technical Conference*, ser. ATC '13, 2013.
- [15] M. A. Sevilla, N. Watkins, I. Jimenez, P. Alvaro, S. Finkelstein, J. LeFevre, and C. Maltzahn, "Malacology: A Programmable Storage System," in *Proceedings of the European Conference on Computer Systems*, ser. Eurosys '17, Belgrade, Serbia.
- [16] S. A. Weil, K. T. Pollack, S. A. Brandt, and E. L. Miller, "Dynamic Metadata Management for Petabyte-Scale File Systems," in *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, ser. SC '04.
- [17] S. V. Patil and G. A. Gibson, "Scale and Concurrency of GIGA+: File System Directories with Millions of Files," in *Proceedings of the Conference on File and Storage Technologies*, ser. FAST '11.
- [18] L. Xiao, K. Ren, Q. Zheng, and G. A. Gibson, "ShardFS vs. IndexFS: Replication vs. Caching Strategies for Distributed Metadata Management in Cloud Storage Systems," in *Proceedings of the Symposium on Cloud Computing*, ser. SoCC '15.
- [19] K. o. V. Shvachko, "HDFS Scalability: The Limits to Growth."
- [20] P. Carns, S. Lang, R. Ross, M. Vilayannur, J. Kunkel, and T. Ludwig, "Small-file Access in Parallel File Systems," in *Proceedings of the Symposium on Parallel and Distributed Processing*, ser. IPDPS '09.
- [21] M. Eshel, R. Haskin, D. Hildebrand, M. Naik, F. Schmuck, and R. Tewari, "Panache: A Parallel File System Cache for Global File Access," in *Proceedings of the Conference on File and Storage Technologies*, ser. FAST '10.
- [22] A. Thusoo, Z. Shao, S. Anthony, D. Borthakur, N. Jain, J. Sen Sarma, R. Murthy, and H. Liu, "Data Warehousing and Analytics Infrastructure at Facebook," in *Proceedings of the SIGMOD International Conference on Management of Data*, ser. SIGMOD '10.
- [23] K. McKusick and S. Quinlan, "GFS: Evolution on Fast-forward," *Communications ACM*, vol. 53, no. 3, pp. 42–49, Mar. 2010.
- [24] D. Roselli, J. R. Lorch, and T. E. Anderson, "A Comparison of File System Workloads," in *Proceedings of the USENIX Annual Technical Conference*, ser. ATC '00.
- [25] C. L. Abad, H. Luu, N. Roberts, K. Lee, Y. Lu, and R. H. Campbell, "Metadata Traces and Workload Models for Evaluating Big Storage Systems," in *Proceedings of the International Conference on Utility and Cloud Computing*, ser. UCC '12.
- [26] S. R. Alam, H. N. El-Harake, K. Howard, N. Stringfellow, and F. Verzeloni, "Parallel I/O and the Metadata Wall," in *Proceedings of the Workshop on Parallel Data Storage*, ser. PDSW '11.
- [27] S. A. Weil, S. A. Brandt, E. L. Miller, D. D. E. Long, and C. Maltzahn, "Ceph: A Scalable, High-Performance Distributed File System," in *Proceedings of the Symposium on Operating Systems Design and Implementation*, ser. OSDI '06.