

Tintenfisch: Namespace Schemas for Scalable File System Metadata Generation

Michael A. Sevilla, Other People, Carlos Maltzahn
University of California, Santa Cruz
 {msevilla, people, carlosm}@ucsc.edu

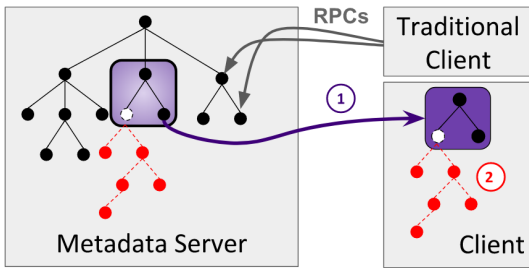


Figure 1: In (1), clients decouple file system subtrees and interact with their copies locally for high performance. In (2), clients and metadata servers generate subtrees using “namespace schemas”, thus reducing RPC load.

1 Introduction

We propose Tintenfisch, a file system that allows users to succinctly express the structure and patterns of the metadata they intend to create. They can also merge new metadata (that they did not explicitly state up front) into the global namespace. Using this semantic knowledge, Tintenfisch can optimize performance by reducing the number of RPCs needed for (1) metadata writes because clients/servers can create metadata independently and (2) metadata reads because clients can construct metadata and pull data directly from the object store. Figure 1 provides an architectural overview; clients first decouple the file system subtree they want to operate on¹ then clients and metadata servers lazily generate subtrees as needed using a namespace schema (described in §3). The namespace schema is stored in the root inode of the decoupled subtree.

The fundamental insight is that the client and server both understand the final structure of the file system metadata so there is no need to communicate. The idea uses concepts from decoupled namespaces [16, 17] and patterned IO [7] to build a scalable global namespace.

¹This is not a contribution as it was presented in [12].

Less work is done on the metadata servers and clients pick up some of the metadata load. This approach is similar to predicate push downs in databases, where structure is described to lower storage layers using XML or JSON [5]. It is our hope that Tintenfisch will also be able to change the representation or structure of the file system metadata according to the file type or workload. We have the following contributions:

- namespace descriptions and overheads for examples from different domains (high performance computing, high energy physics, and large scale simulations)
- namespace schemas: a technique to compact metadata, thus reducing RPC amplification and facilitating lazy metadata generation when needed
- a programmable storage approach that pushes user-defined functionality into the storage system, facilitating application-specific storage stacks using a ‘dirty-slate’ approach

2 Motivating Examples

We look at the namespaces of 3 applications. Each is from different domains and this list is not meant to be exhaustive, as similar organizations exist for many domains, even something as distant as the mail application on a Mac. To highlight the scalability challenges for file system metadata management, we focus on large scale systems in high performance computing, and high energy physics, and large scale simulations.

We benchmark over Ceph (Jewel version) with n object storage daemons (OSDs), 1 metadata server (MDS), 1 monitor server (MON), and 1 client. We use 3 OSDs because it sustains 16 concurrent writes of 4MB at 600MB/s for 2 minutes. 250MB/s is the max speed of the SSDs, so the setup achieves 80% of the cluster SSD bandwidth. We use CephFS, the POSIX-compliant file system that uses Ceph’s RADOS object store [14], as the

underlying file system. This analysis focuses on the file system metadata RPCs between the client and metadata server and does not include the RPCs needed to write and read actual data. CephFS uses a cluster of metadata servers to service file system metadata requests [15] and to characterize the workload, we instrumented the metadata server to track the number of each request type².

2.1 High Performance Computing: PLFS

Checkpointing performs small writes to a single shared file but because file systems are optimized for large writes, performance is poor. To be specific, it is easier for applications to write checkpoints to a single file with small, unaligned, writes of varying length varying write (N-1) but general-purpose distributed file systems are designed for writes to different files (N-N).

The general problem is that the application understands the workload but it cannot communicate a solution to the storage system. The common solution is to add middleware (i.e. software that sits between the application and the storage system) to translate the data into a format the storage system performs well at.

The problem is that the underlying file system cannot keep up with the metadata load imposed by PLFS. PLFS creates an index entry for every write, which results in large per-processes tables [6]. This makes reading or scanning a logical file slow because PLFS must construct a global index by reading each process's local index. This process incurs a `readdir` and, if the file is open by another process, an additional `stat()` because metadata cannot be cached in the container [2].

2.1.1 System Architecture

PLFS [2] solved the checkpoint-restart problem by mapping logical files to physical files on the underlying file system. The solution targets N-1 strided checkpoints, where many processes write small IOs to offsets in the same logical file. The key insight of PLFS is that general purpose file systems perform well for applications that use N-N checkpoints and that the N-1 strided checkpoint style can be transformed with a thin interposition layer. To map offsets in the logical file to physical files each process maintains an index of {logical offset, physical offset, length, physical block id}.

PLFS maps an application's preferred data layout into one that the file system performs well on. Each process appends writes to a different data file in the hierarchical file system and records an offset and length are recorded in an index file. Reads aggregate per-process index files into a global index file, which it uses as lookup table for logical file.

²This code was merged into the Ceph project.

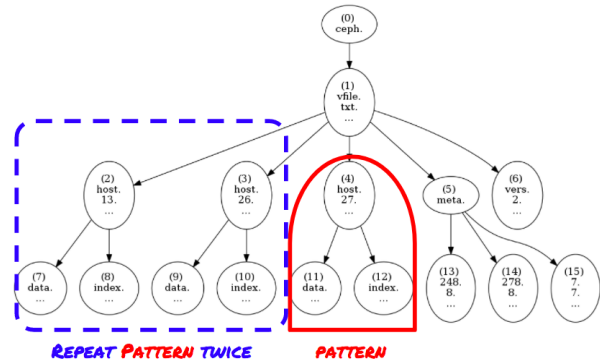


Figure 2: The PLFS file system namespace is structured and predictable; the pattern (solid line) is repeated for each hosts. In this case, there are three hosts so the pattern is repeated two more times.

2.1.2 Namespace Description

When PLFS maps a logical file to many physical files, it deterministically creates the file system namespace in the backend file system. For metadata writes, the number of directories is dependent on the number of PLFS servers. When a file is created in a PLFS mount, a directory called a container is created in the underlying file system. Inside the container, directories are created per server resulting in m requests.

The namespace structure for 3 processes writing to variable offsets in a single PLFS file is shown in Figure 2. The `host*` directory and `data*/index` files (denoted by the solid red line) are created for every node in the system. The pattern is repeated twice (denoted by the dashed blue line) in the Figure, representing 2 additional hosts.

2.1.3 Overhead: Processing RPCs

For writing, the number of files is dependent on the number of PLFS processes. Processes write data and index files to the directory assigned to their server. The number of requests, which can range from $2n$ to $4n$, is a function of files per directory. If each server only has one process, then the number of file create requests will be $2n$ because each process can cache the directory inode and create files in one request. With multiple processes per server, clients write to the same directory and the number of requests doubles because processes cannot cache the directory inode.

For metadata reads, transforming write IO in this way has space and read overheads. In PLFS, this is a problem because index files need to be coalesced on reads. Patterned PLFS [7] reduces the space overheads by storing formulas, instead of index files, to represent

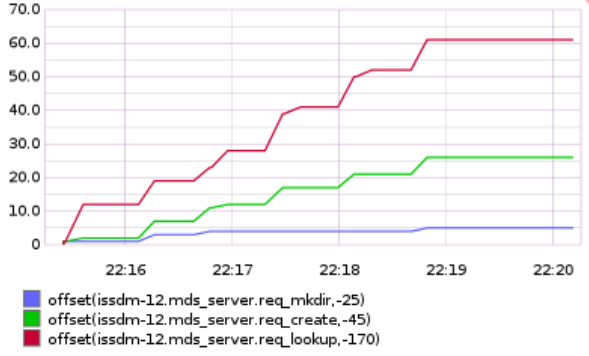


Figure 3: The requests scale linearly with the number of clients. `lookup()` requests dominate because clients share the root and must do path traversal.

write behavior. Diddlings [6] transfers index files after each write to absorb the transfer overheads up front. While these approaches help alleviate read overheads, they do not reduce the file system metadata load, which is the real problem. Reading the index file still requires a file system metadata operation.

Takeaway: PLFS has namespaces with many requests, so the scalability is limited by the throughput of the metadata server in the underlying distributed file system.

2.2 High Energy Physics: ROOT

The High Energy Physics (HEP) community uses a framework called ROOT to manipulate, manage, and visualize data about proton-proton collisions collected at the large hadron collider (LHC). The data is used to re-simulate phenomena of interest for analysis and there are different types of reconstructions each with various granularities. The data is organized as nested, object oriented event data and the length of the runs (and thus the number of events) are of arbitrary length and type (e.g., particle objects, records of low-level detector hit properties, etc.). Reconstruction takes detector conditions (e.g., alignment, position of the beam, etc.) as input. Data is streamed from the LHC into large immutable datasets, stored publicly in data centers around the world. Physicists analyze the dataset by downloading interesting events stored in ROOT files.

2.2.1 System Architecture

A ROOT file is a list of objects and data is accessed by consulting metadata in the header and seeking to a location in the bytestream, as shown in Figure 4a. Scattered in the ROOT file is both data and ROOT file specific metadata called Logical Record Headers (LRH). For

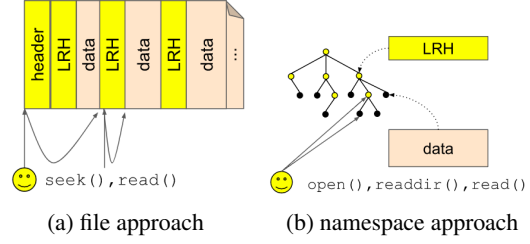


Figure 4: (a) shows how the file approach for handling HEP data stores everything in a single ROOT file, where the client reads the header and seeks to metadata (LRH) and data. For ROOT files stored in distributed file systems reads will have amplification because the system’s striping strategies are not aligned to Baskets. (b) shows how the namespace approach stores Baskets as files in the file system namespace so clients read only the data they need.

this discussion, the following objects are of interest: a “Tree” is a table of a collection of events, listed sequentially and stored in a flat namespace; a “Branch” is a data container representing columns of a Tree; and “Baskets” are byte ranges partitioned by events and indexed by LRHs. Other objects, such as “Keys”, “Directories”, and “Leaves”, contain HEP-specific metadata. Clients request Branches and data is transferred as Baskets; so Branches are the logical view of the data for users and Baskets are the compression, parallelization, and transfer unit. In summary, ROOT files are self-describing files containing data located with metadata and serialized/deserialized with the ROOT framework. Much of the development was done at CERN in parallel with other HPC ventures. As a result, the strategies are similar to techniques used in HDF5, Parquet, and Avro.

The advantages of the ROOT framework is the ability to (1) read only parts of the data and (2) easily ingest remote data over the network. Unfortunately, the HEP community is running into scalability problems. The current effort is to integrate the ROOT framework with Ceph. But naive approaches such as storing ROOT files as objects in an object store or files in a file system have read amplification (*i.e.* read more than is necessary). Users would pull the entire GB-sized blob to their laptop instead of reading metadata and Branches a-la-cart. An alternative strategy that attempts to reduce read amplification is the “namespace approach” [11].

2.2.2 Namespace Description

The namespace approach views a ROOT file as a namespace of data. At the top of the namespace are Keys, each containing pointers to groups of Branches. For example, “MetaData” has data about the run and “Events” has all

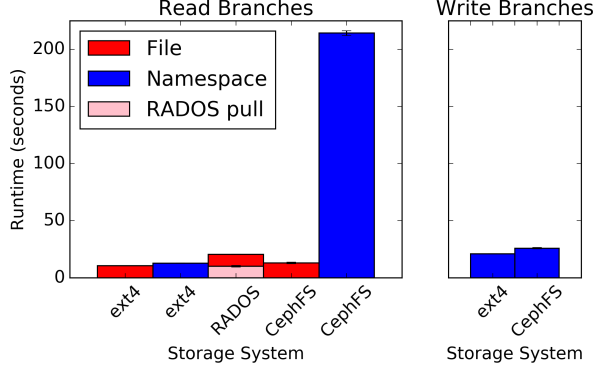


Figure 5: The “File” bars show the runtime of storing HEP data as a single ROOT file and the “Namespace” bars show the runtime of storing a file per Basket. The baseline is the runtime of “ext4” on a single node, “RADOS” is reading HEP data as a Ceph object, and “CephFS” is reading HEP data as files in a CephFS mount. “CephFS” is far slower because of the file system metadata load.

the proton-proton activity. Physicists ask for Branches, where each Branch can be made up of multiple sub-Branches (*i.e.* Events/Branch0/Branch1), similar to pathname components in a file system file name. To accommodate this model, the namespace approach partitions the ROOT file into a file system namespace. As shown in Figure 4b, file system directories hold Branch metadata and files contain Baskets. Clients only pull baskets they care about, which prevents read amplification. Unfortunately, storing this metadata in a file system would overwhelm most file systems in two ways: (1) too many inodes and (2) per-file overhead. To quantify (1), consider the Analysis Object Dataset which has a petabyte of data sets made up of a million ROOT files. If these ROOT files are split by Branches, there would be a billion files. To quantify (2), we benchmark a simple ROOT workload over CephFS.

2.2.3 Overheads

We benchmark the write and read overhead of storing HEP data with the file approach stored as one object in an object store, with the file approach stored as one file in a file system, and with the namespace approach stored as many files in a file system. The file approaches are deployed without any changes to the ROOT framework. For the namespace approach, HEP-specific metadata is mapped onto the file system namespace. In CephFS, Baskets are stored in Ceph objects and the Branch hierarchy is managed by the metadata server. Clients contact the metadata server with a Branch request, receive back the Branch hierarchy necessary to name the Ceph object

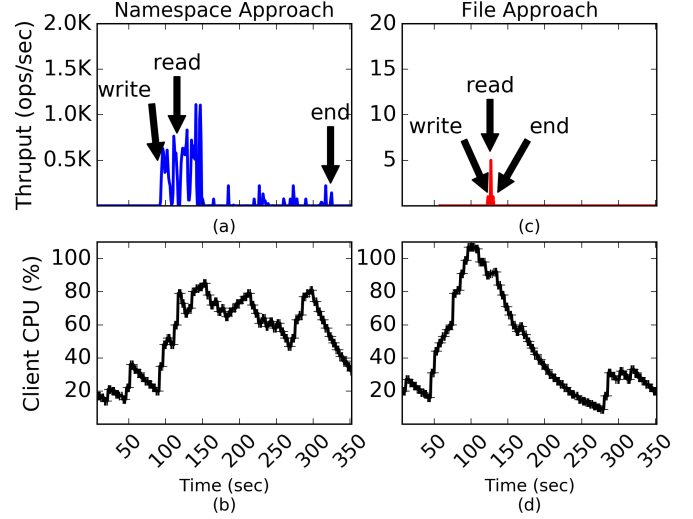


Figure 6: Reading and writing high-energy physics (HEP) data as many files allows physicists to read just the data they care about. But using this namespace approach sends many RPCs to the metadata server (a), resulting in worse performance and lower CPU utilization at the client (b). Alternatively, using the traditional file approach has IO amplification because all data moves over the network but less RPCs (c), better performance, and higher client CPU utilization (d).

containing the Basket as well as the deserialization metadata necessary to read the object. The workload is a list of Branch accesses from a trace of the NPTupleMaker high energy physics application. Each Branch access is:

Branch0/Branch1,3,1740718587,5847,97,136

where the tuple is the full Branch name, Basket number, offset into the ROOT file, size of the Basket, start entry of the Basket, and end entry of the Basket. For the file approach, we use the offset into the ROOT file and the size of the Basket. In setup 1, the ROOT file is pulled locally and the Branches are read from the file. In setup 2, the offset and size of the read are sent to the CephFS metadata server. For setup 3, the full Branch name and Basket number are used to traverse the file system namespace.

The read and write performance for the different approaches are shown in Figure 5, where the x -axis is different storage backends, the y -axis is runtime, and the error bars are the standard deviations for six runs. Compared to the single node baseline (ext4), the namespace approach on CephFS is $16.8\times$ slower. The reason is shown in Figure 6. The file system metadata accesses, characterized by many open() requests, incur many RPCs. This causes worse performance and lower client CPU utilization compared to reading a single ROOT file. So the cost of read amplification in the

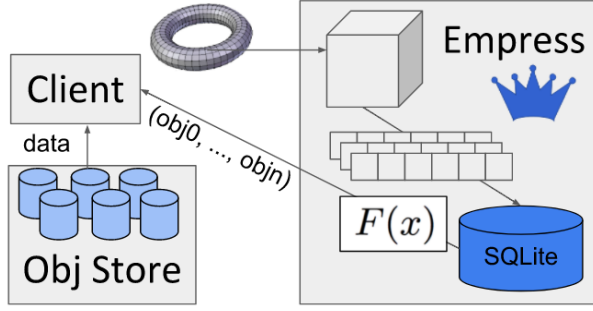


Figure 7: The SIRIUS project uses Empress to store metadata for bounding boxes in a 3D torus. Bounding box coordinates and a list of object names are stored in SQLite. The object names are calculated using an objecter function, labeled $F(x)$ above. For reads, clients query Empress for the object list before reading from the object store.

file approach is offset by the cost of doing namespace operations. For this experiment, the ROOT file is 1.7GB and 65% of the file is accessed so the namespace approach might be more scalable for different workloads.

Takeaway: The ROOT framework has large namespaces and namespaces with many requests, so the scalability is limited by the capacity and request throughput of the metadata server in the underlying distributed file system.

2.3 Large Scale Simulations: SIRIUS

SIRIUS [8] is the Exascale storage system being designed for the Storage System and I/O (SSIO) initiative. The core tenant of the SIRIUS project is application hints that allow the storage to reconfigure itself for higher performance. Techniques include tiering, management policies, data layout, quality of service, and load balancing.

2.3.1 System Architecture

SIRIUS includes a metadata service called Empress [9], which is a query-able SQLite instance that stores metadata for bounding boxes (*i.e.* a 3-dimensional coordinate space). Figure 7 shows how metadata is stored in and retrieved from Empress; the entire simulation space, represented as a 3D torus, is partitioned into bounding boxes, whose coordinates are stored in SQLite for each variable (*e.g.*, temperature, pressure, etc.). The objecter function, $F(x)$, translates the bounding box coordinates into a list of object names, which are used to write/read data to/from the scalable object store. SIRIUS aims to satisfy queries structured like the “Six Patterns in HPC” [10].

2.3.2 Namespace Description

The seven columns in the database are six coordinates, which are integers representing the contents at a location in the global space, and a string representing the object name. The global space is partitioned into non-overlapping, regular shaped cells. Each row in the database is duplicated for each variable in the simulation because variables may have a different partitioning of the global space; for example temperature is computed for every cell while pressure is computed for ever n cells. Figure 7 only has three variables, represented by the rows of boxes but most simulations will have a minimum of 10 variables.

2.3.3 Overheads

Empress lists all objects names for a bounding box and clients filter for data of interest. This is a problem because of the size of the object name list. Listing large numbers of items in file systems is notoriously slow and studies on 1s have shown the operation to be especially heavy-weight [3, 4]. As currently designed, Empress would store too many object names while the target query types only request subsets of the resulting list.

of object names: a back-of-the-envelope calculation for the number of object names in the system is:

$$= \frac{(\# \text{ processes}) \times (\text{data/process}) \times (\text{timesteps}) \times (\text{variables})}{(\text{object size})}$$

where reasonable values would be one million processes, each writing 8GB of data for 100 timesteps and 10 variables. This represents a simulation space of $10K \times 10K \times 10K$ cells containing 8 byte floats. Using an 8MB stripe size (the optimal object size in RADOS), the object name list size is:

$$= \frac{(1 * 10^6) \times (8 * 10^9) \times (100) \times (10)}{(8 * 10^6)} = 1 * 10^{12} \text{ objects}$$

Queries use a subset of object names. When running queries characterized by the “Six Patterns in HPC”, Empress must exhaustively list objects and filter the coordinates of interest. In addition to the time it would take to search 1 trillion objects in SQLite, another problem is the actual storage footprint of storing this many items.

Takeaway: the SIRIUS storage system has large namespaces so the scalability is limited by the capacity of the metadata server in the underlying distributed file system; namely by the performance of scanning lists.

3 Namespace Schemas

For n processes on m servers:

```
# of dirs =  $m \times \text{mkdir}()$ 
# of file =  $2 \times n \times m$ 
# of file per dir =  $n/m$ 
```

(a) Function schema for PLFS

```
local box require 'box2d'
for i=_x,_x+x do -- iterate over the
  for j=_y,_y+y do -- given bounding
    for k=_z,_z+z do -- box coordinates
      if temperature>30 then
        -- partition object list one way, e.g.,
        b0, b1 = box.nsplrit(2)
      else
        -- partition a different way, e.g.,
        b0, b1, b2, b3 = box.nsplrit(4)
      end
    end
  end
end end end
return obj_list
```

(b) Code schema for SIRIUS

pointer_schema: (o0, o2, o9)

code_schema:

```
void recurseBranch(TObjArray *o) {
  TIter i(o);
  for(TBranch *b=i.Next(); i.Next()!=0; b=i.Next()) {
    processBranch(b);
    recurseBranch(b->GetListOfBranches());
  }
}
```

(c) Pointer and Code schemas for HEP

Figure 8: Namespace schemas that generate subtrees for 3 motivating examples.

For three domain-specific applications and use-cases, we have identified different scalability challenges:

1. namespaces with many requests
2. namespaces managed by remote servers
3. namespaces that are too large

Tintenfisch addresses all three challenges by having clients/servers store namespace schemas, which generate file system metadata. This approach reduces RPC load (addresses challenges 1 and 2) and facilitates lazy file system metadata generation when the metadata is needed (addresses challenge 3). Tintenfisch relies on the user to design effective namespace schemas that leverage domain-specific knowledge to get the highest performance. This programmable storage approach [?] helps application developers tailor the storage system to the use case without having to design a new storage system from scratch.

Tintenfisch is built on Cudele [12] so a centralized, globally consistent metadata service (either a single metadata server or a cluster of active-active metadata servers) provides clients with the root inode of the subtree of interest and clients can do metadata IO locally with the consistency/durability semantics they require. In Tintenfisch, the namespace schema is stored in the directory inode of the root of the subtree that the client cares about. We use the “file type” interface from the Malacology [13] project to facilitate this domain-specific functionality. This is similar to push-down predicates in databases, where the application is providing domain-specific knowledge that the storage system knows how to leverage. We have defined three types of namespace schemas: formula, code, and pointers.

3.1 Formula Schema

For this schema, a formula that generates the file system namespace is stored in the inode. This formula takes domain-specific information as input and produces a list of files and directories. For example, because PLFS deterministically creates files and directories based on the number of clients, Tintenfisch can use a formula schema like the one in Figure 8a. The function takes as input the number of processes and hosts in the cluster and outputs the number of directories, the number of files, and the number of files per directory. For the example, the namespace drawn in Figure 2 can be generated with the formula in Figure 8a using an input of 3 hosts each with 1 process. The output is 3 directories and 6 files, with 2 files per directory.

3.2 Code Schema

Sometimes the namespace schema logic is too complex to store as a single function or requires external libraries to interpret metadata. For example, the SIRIUS use case constructs the namespace using domain-specific partitioning logic written in Lua (for sandboxing purposes). Tintenfisch provides a code schema that gives users the flexibility to write programs that generate the namespace.

A code schema for the SIRIUS project is shown in Figure 8b. A namespace is constructed by iterating through the bounding box coordinates and checking if a threshold temperature is eclipsed. If it is, then extra names are generated using the box2d Lua package. Although the partitioning function itself is not realistic, it shows how code schemas can accommodate domain-specific data layout policies that are complex and/or require external libraries.

3.3 Pointer Schema

Sometimes there is no formal specification for the namespace. For example, the ROOT framework uses self-describing files so headers and metadata need to be read for each ROOT file. In these scenarios, a code schema is insufficient for generating the file system namespace because all necessary metadata is in objects scattered in the object store. Pointer schemas reference data in scalable storage and avoids storing large amounts of metadata in inodes, which is a frowned upon in distributed file systems like CephFS [1].

For example, a code schema containing library code for the ROOT framework *and* a pointer schema for referencing the input to the code can be used to describe a ROOT file system namespace. The pointer schema would point to objects in the object store that contain the necessary metadata for constructing the file system namespace. This is shown in Figure 8c; clients requesting Branches would follow the pointer schema to the objects containing metadata and would read them to locate Baskets using the code schema. An added benefit to this solution is that Tintenfisch can lazily construct parts of the namespace as needed, which avoids the problem of running out of inodes discussed in Section 2.2.2.

A Fresh, Unorthodox, Unexpected, Controversial, and Counterintuitive Idea

Global file systems can be scalable if programmed correctly. For example, the following notions are out-dated:

- robust so they are fast... but we show that today's apps are so large that we need to specialized storage systems
- general because they have been around for so long... but we show that most apps don't need fs metadata
- subject to data IO performance... but we show that metadata is slow

References

- [1] Ceph Documentation. http://docs.ceph.com/docs/jewel/dev/mds_internals/data-structures/, December 2017.
- [2] BENT, J., GIBSON, G., GRIDER, G., MCCLELLAND, B., NOWOCZYNSKI, P., NUNEZ, J., POLTE, M., AND WINGATE, M. PLFS: A Checkpoint Filesystem for Parallel Applications. In *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, SC '09.
- [3] CARNS, P., LANG, S., ROSS, R., VILAYANNUR, M., KUNKEL, J., AND LUDWIG, T. Small-file Access in Parallel File Systems. In *Proceedings of the Symposium on Parallel and Distributed Processing*, IPDPS '09.
- [4] ESHEL, M., HASKIN, R., HILDEBRAND, D., NAIK, M., SCHMUCK, F., AND TEWARI, R. Panache: A Parallel File System Cache for Global File Access. In *Proceedings of the Conference on File and Storage Technologies*, FAST '10.
- [5] FINKELSTEIN, S. Personal communication.
- [6] GRIDER, G. If the plfs index is still too large too efficiently read/distribute pull out the hammerplfs collectives, reduces plfs index to nearly nothing disabling strategy for optimizing reads and active analysis.
- [7] HE, J. Io acceleration with pattern detection.
- [8] KLASKY, S. A., ABBASI, H., AINSWORTH, M., CHOI, J., CURRY, M., KURC, T., LIU, Q., LOFSTEAD, J., MALTZAHN, C., PARASHAR, M., PODHORSZKI, N., SUCHYTA, E., WANG, F., WOLF, M., CHANG, C. S., CHURCHILL, M., AND ETHIER, S. Exascale Storage Systems the SIRIUS Way. *Journal of Physics: Conference Series*.
- [9] LAWSON, M., ULMER, C., MUKHERJEE, S., TEMPLET, G., LOFSTEAD, J. F., LEVY, S., WIDENER, P. M., AND KORDENBROCK, T. Empress: Extensible Metadata Provider for Extreme-Scale Scientific Simulations. In *Proceedings of the Joint International Workshop on Parallel Data Storage & Data Intensive Scalable Computing Systems*, PDSW '17.
- [10] LOFSTEAD, J., POLTE, M., GIBSON, G., KLASKY, S., SCHWAN, K., OLDFIELD, R., WOLF, M., AND LIU, Q. Six Degrees of Scientific Data: Reading Patterns for Extreme Scale Science IO. In *Proceedings of High Performance and Distributed Computing*, HPDC '11.
- [11] PIVARSKI, J. How to make a petabyte ROOT file: proposal for managing data with columnar granularity.
- [12] SEVILLA, M. A., JIMENEZ, I., WATKINS, N., FINKELSTEIN, S., LEFEVRE, J., ALVARO, P., AND MALTZAHN, C. Cud-ele: An API and Framework for Programmable Consistency and Durability in a Global Namespace. In *Proceedings of the 32nd IEEE International Parallel and Distributed Processing Symposium* (2018), IPDPS '18.
- [13] SEVILLA, M. A., WATKINS, N., JIMENEZ, I., ALVARO, P., FINKELSTEIN, S., LEFEVRE, J., AND MALTZAHN, C. Malacology: A Programmable Storage System. In *Proceedings of the European Conference on Computer Systems* (2017), EuroSys '17.
- [14] WEIL, S. A., BRANDT, S. A., MILLER, E. L., LONG, D. D. E., AND MALTZAHN, C. Ceph: A Scalable, High-Performance Distributed File System. In *Proceedings of the Symposium on Operating Systems Design and Implementation*, OSDI '06.
- [15] WEIL, S. A., POLLACK, K. T., BRANDT, S. A., AND MILLER, E. L. Dynamic Metadata Management for Petabyte-Scale File Systems. In *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, SC '04.
- [16] ZHENG, Q., REN, K., AND GIBSON, G. BatchFS: Scaling the File System Control Plane with Client-funded Metadata Servers. In *Proceedings of the Workshop on Parallel Data Storage*, PDSW '14.
- [17] ZHENG, Q., REN, K., GIBSON, G., SETTLEMYER, B. W., AND GRIDER, G. DeltaFS: Exascale File Systems Scale Better Without Dedicated Servers. In *Proceedings of the Workshop on Parallel Data Storage* (2015), PDSW '15.