

Malacology: A Programmable Storage System Built on Ceph

Carlos Maltzahn, Michael Sevilla, Noah Watkins, Ivo Jimenez

UC Santa Cruz

{carlosm,msevilla,jayhawk,ivo}@soe.ucsc.edu

Abstract—Ceph is a distributed storage system with many code-hardened components, yet many of these subsystems are never re-used or re-purposed for other tasks. Monitor processes (MONs) maintain systems in the cluster using consensus, versioning, and consistency protocols. Object Storage Daemons (OSDs) store data on disk using sophisticated peer-to-peer techniques like replication, load balancing, and consensus. Metadata daemons (MDSs) act as gateways for file-based storage using protocols for consistency, balancing load, and mediating shared data. In this work, we take a “programmable storage” approach to building systems that emphasizes the re-use and re-purposement of different Ceph subsystems. Using this framework, we build two services, a POSIX metadata load balancer and a distributed shared commit-log, that re-use load balancing, recovery, cache coherence/capabilities, and object class functionality already present in Ceph.

I. INTRODUCTIONS

Large distributed systems tackle difficult distributed systems problems with code-hardened subsystems. For example, Ceph~[?] addresses *durability* with its RADOS object store (e.g., replication, erasure coding, and data scrubbing), *consistent versioning* by having daemons exchange “maps” of the cluster configuration, and *consensus* by having monitor daemons (MONs) use PAXOS. We contend that re-using and re-purposing these code-hardened subsystems is paramount to (1) improving the longevity and community uptake of “research quality” code and (2) avoiding duplication of the same protocols and algorithms throughout the system. Unfortunately, many internal subsystems are not exposed to other parts of the systems.

In this paper, we examine the programmability of Ceph, the open-source storage system solution backed by Red Hat. Ceph is known as the swiss army knife of storage, offering file, object, and block APIs for applications. The main draw of Ceph is the flexibility to use all three layers on the same storage system. While Ceph is one of the most flexible distributed systems out there, we contend that the system could be *even more programmable*. With minimal changes to the architecture, and building on many of the subsystems already baked into Ceph, we can build large research-quality systems.

We present Malacology, a programmable storage system capable of incorporating new functionality and re-purposing existing subsystems. Administrators inject functionality into system as Lua scripts. We build the framework on Ceph~[?] by adding a command to the monitor daemons (MON) and a Lua interpreter to the object storage daemon (OSD) and metadata

server daemon (MDS). As shown in Figure~??, this framework is expressive enough to provide the functionality necessary for implementing other research-quality systems and services. Our contributions are:

- a programmable storage system implementation
- re-using code-hardened systems: sandboxed and vetted
- example systems that use this framework
 1. shared log service based on CORFU~[?]
 2. metadata load balancer based on Mantle~[?]

In the remainder of this paper we demonstrate the power of programmable storage. First we describe programmable storage in more depth (§IV). Next we introduce our implementation of a programmable storage framework that exposes and re-uses many structures in Ceph (§IV-A). We conclude with descriptions and evaluations of these ideas by synthesizing entirely new storage services on an existing system through configuration and small changes: a distributed shared log (§??) and a programmable metadata load balancer (§??).

II. HIGHLY TAILORED AND APPLICATION-SPECIFIC STORAGE SYSTEMS

When application goals are not met by a storage system the most common reaction is to design a workaround. Workarounds roughly fall into one of two categories: so called “bolt-on” services that introduce a 3rd party system (e.g. a metadata service), or expanded application responsibility in the form of data management (e.g. a new data layout).

Extra Services “Bolt-on” services are designed to improve overall application performance, but come at the expensive of additional sub-systems and dependencies that the application must manage, as well as trust. For example, it is well understood that MapReduce performs poorly for iterative and interactive computation due to its failure model that heavily relies on on-disk storage of intermediate data. Many have added services to Hadoop to keep more data in the runtime (e.g., HaLoop[?], Twister[?], CGL-MapReduce[?], MixApart[?]). While performance improves, it comes at a cost: “bolt-on” services frequently result in overly complex systems that re-implement functionality and re-execute redundant code, unnecessarily increasing the likelihood of bugs.

Application Changes The second approach to adapting to a storage system deficiency is to change the application itself by adding more data management intelligence, often into the application itself, or as domain-specific middleware. For

instance, an application may change to exploit data locality or I/O parallelism in a distributed storage system. This is not a bad proposition, but creates a coupling that is highly tied to the underlying physical properties of the system, making it difficult to adapt to future changes at the storage system level.

Storage Changes When these two approaches fail to meet the needs of the application, developers turn their attention to the storage system itself. For example, HDFS has been the focus of scalability concerns, especially for metadata-intensive workloads~[?]. This has led to modifications to its architecture or API~[?] to improve performance. Yet another approach is to “modify” a storage system using auto-tuning techniques that attempt to find a good solution among a huge space of available system configurations. However, in practice auto-tuning is limited to only the configuration “knobs” that the storage system exposes (e.g. block size). For instance, auto-tuning may be capable of identifying instances in which new data layouts would benefit a workload, but unless the system can provide such a transformation, the option is left off the table.

We advocate a new approach that we refer to as *storage programmability* which is a method by which an application communicates its requirements to the storage system in a way that allows the application to realize a new behavior without sacrificing the correctness of the underlying system.

Active storage is a hybrid approach to changing the application and storage system. Pushing computation closer to the data is not a new idea. Active storage techniques are used in production Ceph environments to have object storage devices (OSDs) process data before sending over the network or storing it. The code that does the processing in the OSD is called an `object storage interface` and it can be loaded at runtime and customized with user-defined functionality. The basic idea is shown in Figure~??, where the `libcls_md5.so` shared library performs the MD5 hash on an object at the OSD instead of transferring data over the network. This ability to carry out arbitrary operations on objects stored on OSDs helps applications improve performance by optimizing things like network round trips, data movement, and remote resources which may be idle.

Developers and users actively develop new object storage interfaces. Table~?? shows the wide-variety of object interfaces that have been co-designed with applications that run on top of Ceph and Figure~?? shows a dramatic growth in the use of co-designed interfaces since 2010. Figure~?? examines this growth in interfaces further by showing the lines of code that are changed ($\{y\}$ axis) over time ($\{x\}$ axis). The prevalence of blue dots indicates that users frequently update their interfaces. This interface churn reflects both the popularity of object interfaces and the complexity of the processing being done by the OSDs. While we consider active storage to be an excellent example of *programmability*, what separates our proposal from previous work is the observation that so much more of the storage system can be reused to construct advanced, domain-specific interfaces.

The popularity of object interfaces hints at there trends

in the Ceph community: (1) increasingly, the default algorithms/tunables of the storage system are insufficient for the application’s performance goals, (2) programmers are becoming more aware of their application’s behavior, and (3) programmers know how to manage resources to improve performance. Programmers gravitate towards object interfaces because it gives them ability to tell the storage system about their application: if it is CPU or IO bound, if it has locality, if its size has the potential to overload a single proxy node, etc. The programmers know what the problem is and how to solve it, but until object interfaces, had no way to tell the storage system how to handle their data.

III. CEPH IS A PRODUCTION-QUALITY DISTRIBUTE SYSTEM

A. Active Storage

Object storage interfaces are compiled into shared libraries and loaded into a running OSD daemon using `dlopen()`. Because the shared library has to link against symbols found in the running executable, the interfaces are stored on the local file systems of each OSD so they can be versioned and distributed alongside the same Ceph binaries. This policy is for safety, since the OSDs could be on different distros that provide different versioning capabilities. Also, this approach allows bug fixes to co-evolve with the rest of the Ceph installation.

Although customizable OSD interfaces are powerful, the current implementation has drawbacks. OSD interfaces are written in C/C++ and compiled into a shared library, so the developer must account for different target architectures. Second, C/C++ provides more functionality than is necessary since the snippets of interface code mainly set policies and perform simple operations. Finally, developing these interfaces in C/C++ has high overhead. The developer must learn how the OSD dynamically loads the shared libraries ({e.g., OSDs rely on a strict naming convention to find shared libraries}*), how to get their interfaces compiled using Ceph’s `make` system, and how to debug issues that are not related to their specific interface (*{e.g.} the OSD cannot find the shared library*) - this learning curve is unacceptable for non-Ceph developers, especially since most interfaces are one-off solutions specific to their applications.*

Ceph has a whole infrastructure for getting dynamic code into the OSD. The `ClassHandler` class safely opens object interfaces, executes commands defined by the shared library, and fails gracefully with helpful errors if anything goes wrong. Instead of injecting strings directly into Ceph (like we did with the original `Mantle` implementation), the `ClassHandler` class safely opens shared code, even code with dependencies that aren’t in Ceph itself (e.g., `Lua`). With the `ClassHandler`, we get the safety and robustness of loading dynamic code, the ease of transferring state between object interfaces and the OSD internals, integration with testing and correctness suites, and `structs` for interface data and handlers.

B. Durability

Ceph provides storage by striping and replicating data across RADOS, Ceph's reliable distributed object store. RADOS uses many techniques to ensure that data is not corrupted or lost, such as erasure coding, replication, and data scrubbing. Furthermore, many of these techniques try to be autonomous so that work is distributed across the cluster. For example, when placement groups change the OSDs rebalance and re-shard data in the background in a process called placement group splitting.

C. Consistency and Versioning

Ceph needs to keep track of cluster state and it does this with separate "services": client authentication, logging, metadata server (MDS) maps, monitor (MON) maps, object storage device (OSD) maps and placement group (PG) maps. These services are all managed in Ceph monitor processes (MONs) and they all talk to a PAXOS instance; this PAXOS instance goes and talks to other PAXOS instances on other MONs so that they can all agree on the correct version of the service; essentially they reach an agreement about what the state of the cluster is.

IV. PROGRAMMABLE STORAGE

A. Implementation

We present a framework for managing Lua object classes comprised of 3 parts:

- 1) general and customizable shared libraries for the OSDs and MDSs using Lua
- 2) storing Lua object classes as RADOS objects
- 3) monitor command for specifying the Lua class

B. Distributing Work with Customizable Interfaces

Our framework has a mechanism for defining and running object and metadata balancer classes using Lua. Our Lua bindings expose functions and symbols both ways; the host program can call functions defined in Lua and the Lua scripts can call functions defined in native C++. These bindings are merged upstream. We choose Lua for 4 reasons: performance, portability, size, and security.

Lua is a fast scripting language. It was designed to be an embedded language and the LuaJIT virtual machines boasts near-native performance~[?]. Lua is frequently used in game engines to set policies but we use it here because most of the user-defined classes in Ceph are policies as well! We do not want to provide specific implementations, like pulling data from objects or transferring them over the network, but instead strive to say {what to do} with the data once we have it. Separating policy from mechanism is a driving factor in using Lua.

Lua is also portable. The object interfaces are Lua is secure (sandboxing).

For object interfaces, clients send Lua classes to the OSD and then they can invoke operations in that class remotely on the OSD. In reality, this feature is a statically loaded object class written in C/C++ that runs dynamically defined object interfaces, where clients access the classes using an

`exec()` wrapper. Although sending the class with each client request is stateless (which has its own set of nice features), it is costly for network bandwidth and difficult to organize at the application level. Also, while the `execute` wrapper simplifies the implementation, it burdens the applications, since they need to be recompiled to use the `exec()` function.

For balancer interfaces, clients put balancers into RADOS and the CephFS metadata balancer invokes the operations in the user-defined class remotely on the MDS.

1) Generalizing the Lua VM: We generalize the Lua VM since it will be used in both the OSD and MDS. We put the core sandbox wrapper for the Lua object interface in a common directory in Ceph and link against it. This core Lua wrapper contains just the dependencies, symbols, and functions, need to run the Lua VM. Some of these functions include `clslua_log()` for transferring logs to the daemon logs and `clslua_pcall()` for calling Lua functions.

Dependencies, symbols, and functions specific to the object or balancer interface are put in the `cls` or `bal` directories, respectively. For example, the Lua object interface uses functions that have placement group filters, cryptography functions, and object metadata (i.e. `cxx_omap_getvals()`), object data, object extended attributes, and object versions – all these are dependencies, symbols, and functions are part of the OSD process but not the MDS process. As a result, we put interface code that uses these in `cls` directory. With this scheme, The OSD will `dlopen()` the shared library created by Lua object interface shared library while the MDS will `dlopen()` the Lua balancer interface shared library; both shared libraries will the Lua core.

Both the object and balancer interfaces define functions and attach them to the core Lua sandbox. For example, the Lua object storage interface class attaches data, extended attribute, object map, and object version functions, to the the Lua core functions; the exact functions are shown in Listing-??. The advantages of this is that we avoid duplicating code, we provide a framework for putting Lua code in other parts of the system, and we remove components and APIs that are too integrated with the OSD.

2) Generalizing the Class Handler:

- Re-Used Components: Class Handler, Lua Class
- Durability with RADOS
- Send functionality with request
- Loading from the file system

Our framework can also load Lua interfaces from the local file system – the same technique as the C/C++ object interfaces. First, the OSD looks for C/C++ shared library. If the OSD cannot find the file, it looks for a Lua script of the same name. The script is read into a string of the C++ object class; this string is later forwarded to the native Lua class. When a requests comes in, the method name is passed with the `exec()` function and the corresponding function in the Lua class is called. Now clients need to only send their Lua object class once and the OSDs will store them locally. Also, clients can execute any Lua handler they want.

3) Storing Lua object classes as RADOS objects: [7:31] Nothing would stop C++ from being stashed in objects and recompiled on whatever platform they are being loaded on dynamically (like a DB compiles each query). There just hasn't been a need for such a feature. The next blog post, which I didn't at all allude to in the one you are reading, shows how the Lua scripts can be put into the OSD map and then monitors manage and version the script, distributing them to each OSD. Stored in objects vs stored in the OSD map is a debate, I guess. I am not sure which makes more sense.

C. Monitor command for specifying Lua class

We added a command, `ceph osd pool set-class <pool> <class> <script>`, that the user uses to inject the Lua object interface into the cluster. By leveraging the monitor daemons, we get the consistency from the monitors' version management, the distribution from the monitors' data structures (which are already distributed), and the durability from the robustness of the monitors and persistence with Paxos. The implementation uses the placement group pool data structure which maintains the policies (e.g., erasure coding) and organization details (e.g., snapshots) for each pool. While stuffing the information into the monitor map, the structure that describes the the monitor topology, was an option, the placement group pool allows finer grained control over different types of data. In regards to the consistency, each time the user injects a new Lua object class it enters the monitor quorum as a proposal; updates to the placement group pool need to wait until accepted, at which point the state will be propagated and versioned. We hooked up the monitor command to the `cls_lua` class by having the OSD pull the script from the placement group instead of trying to dynamically load the shared library with `dlopen()` or the Lua script.

Advantages:

- lets us store/manage interfaces
- does most of the work (versioning, consistency, durability) for us
- gives a new abstraction for dealing/mutating interfaces

The `cls_lua` branch lets users write object classes in Lua, the `lua-rados` client library lets users write applications that talk to RADOS with Lua, and the `cls-client` uses the `lua-rados` library to send Lua object classes to the OSDs equipped with the LuaJIT VM.

D. Specifying the Lua Class

Maintain versions and consistency

- `cls_lua` branch: write object classes in Lua
- `lua-rados`: talk to RADOS with Lua
- `cls-client`: use `lua-rados` + `cls_lua` branch to send Lua object classes to OSDs equipped with LuaJIT VM

V. RESEARCH QUALITY SYSTEMS BUILT ON MALACOLOGY

A. Mantle: A Programmable Metadata Load Balancer

Many distributed file systems decouple metadata and data I/O so that these services can scale independently~[?], [?], [?],

[?], [?], [?]. Despite this optimization, scaling the metadata services is still difficult because metadata accesses impose small and frequent requests on the underlying storage system~[?]. Many techniques for designing the metadata services have been proposed to accommodate this workload: Lustre[?], GFS[?], and HDFS~[?] keep all metadata on one server; GIGA+~[?], IndexFS~[?], Lazy Hybrid~[?], GPFS~[?], and pNFS~[?] hash the file system namespace across a dedicated cluster and cache metadata; Panasas~[?] and CephFS~[?] partition the file system namespace into subtrees and assign them to servers. These systems have novel mechanisms for scalable metadata but the techniques are often "locked-in" to the systems they are implemented on.

Mantle~[?] is a programmable metadata balancer that separates the metadata balancing policies from their mechanisms. Administrators inject code to change how the metadata cluster distributes metadata. In the paper, we showed how to implement a single node metadata service, a distributed metadata services with hashing, and a distributed metadata service with dynamic subtree partitioning.

The Ceph team wants to merge Mantle because the scriptability is useful for debugging, controlling the metadata balancer, and examining trade-offs for different balancers. Unfortunately, this research quality system is not as robust as Ceph and the Ceph team wants more safety, durability, and consistency for the new functionality.

VI. ZLOG: EXPLORING CONSISTENCY, TIERING, AND STRIPING STRATEGIES

Malacology uses the same Active and Typed Storage module presented in `DataMods`~[?]; Asynchronous Service and File Manifolds can be implemented with small changes to the Malacology framework, namely asynchronous object calls and Lua stubs in the inode, respectively.

VII. CONCLUSION AND FUTURE WORK

Programmable storage is a viable method for eliminating duplication of complex error prone software that are used as workarounds for storage system deficiencies. However, this duplication has real-world problems related to reliability. We propose that system expose their services in a safe way allowing application developers to customize system behavior to meet their needs while not sacrificing correctness.

We intend to pursue this work towards the goal of constructing a set of customization points that allow a wide variety of storage system services to be configured on-the-fly in existing systems. This work is one point along that path in which we have looked at a target special-purpose storage system. Ultimately we want to utilize declarative methods for expressing new services.

In conclusion, this paper should be accepted.

not sure why ./build fails without this

VIII. BIBLIOGRAPHY