

Storage Developer Conference September 22-23, 2020

SkyhookDM: Storage and Management of Tabular Data in Ceph

Jeff LeFevre, Carlos Maltzahn Center for Research in Open Source Software University of California, Santa Cruz



Skyhook Data Management Z

SD@

- Open source software, LGPL 2.1 License
- Built on Ceph distributed object storage
- Computational storage for tabular data
- Extensible, scalable, generic

 CROSS
 CENTER FOR RESEARCH IN
OPEN SOURCE SOFTWARE

 ARACHE
ARROW

 COOGLE FlatBuffers: Memory Efficient Serialization Library

Approach

- Programmable storage (<u>See programmability.us</u>)
 - Combine, expose, or extend existing storage services toward new functionality
- In-storage execution of data management tasks
 - Embed external libraries in storage
- Dynamically offload computation to storage servers
- Dynamically reorganize physical data configuration
- Reduce CPU and Network resources for client apps



User-defined extensions to Ceph

- Utilize Ceph's existing object class mechanism ('cls')
 - Extensible framework for objects
 - ceph/src/cls
- Methods executed directly by objects
 - Shared libraries available on all OSDs
- Utilized by Ceph internals
 - CephFS, rgw, rbd, others...

CLS growth in Ceph



Growth of object classes and methods in Ceph mainline

2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

SD@

Snapshot of 'cls' Classes in Ceph

master - ceph / src / cls /

Go to file Add file

ivancich Merge pull request #31058 from	n cbodley/wip-rgw-skip-bilog	86df8b9 4 days ago 🕚 History
2pc_queue	rgw/notifications: persistency - cleanup stale reservations	19 days ago
as cas	cls/cas: replace bool get() with void get()	3 months ago
cephfs	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
cmpomap	cls/cmpomap: add cls module for CMPXATTR-like functionality in omap	5 months ago
hello	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
🛅 journal	cls/journal: use EC pool stripe width for padding appends	5 months ago
lock	cls,rados,rbd,mds,common: Avoid name collision with Windows headers	2 months ago
log	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
🖿 lua	librados: add symbol versioning to the C++ API	2 years ago
numops	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
totp	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
aueue	cls/rgw_gc: Clearing off urgent data in bufferlist, before	3 months ago
🖿 rbd	librbd: track in-progress migration aborting operation	22 days ago
refcount	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
rgw	Merge pull request #31058 from cbodley/wip-rgw-skip-bilog	4 days ago
rgw_gc	cls/rgw_gc: Fixing carriage returns in log statement.	3 months ago
sdk	cls/sdk: Update cls_sdk.cc to work without using namespace	2 years ago
timeindex	cls: Build ceph-osd without using namespace declarations in headers	5 months ago
🔲 user	rgw: introduce safe user-reset-stats	2 months ago
>version	de: Build canh-ord without using namernace declarations in headers	5 months ago

2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

SD@

Read/Write Interface

- Ceph objects can access their local data via two interfaces within 'cls'
 - 1. Chunkstore raw device access
 - KVstore Local instance of RocksDB on OSD (omap interface)
- For us
 - 1. Map tabular data to a device and offset
 - 2. Consider storing tabular data and/or metadata

CLS Interface Examples

20

Some functions available within a cls method

- // read/write
 - cls_cxx_read(ctx, off, len, buf)
 - cls_cxx_write(ctx off, len, buf)
 - cls_cxx_replace(ctx, off, len, buf)
- // metadata
 - cls_cxx_setaxxtr(ctx, name, buf)
 - cls_cxx_stat(ctx, size_t, NULL)
- // utilize local RocksDB instance on the OSD
 - cls_cxx_map_setvals(ctx, map<str, buf>)
 - cls_cxx_map_getval(ctx, key, buf)
 - cls_cxx_map_getvals(ctx, keystart, nkeys, map<str,buf>, more)

CLS Interface Examples

Some functions available within a cls method

- // read/write
 - cls_cxx_read(ctx, off, len, buf)
 - cls_cxx_write(ctx off, len, buf)
 - cls_cxx_replace(ctx, off, len, buf)
- // metadata
 - cls_cxx_setaxxtr(ctx, name, buf)
 - cls_cxx_stat(ctx, size_t, NULL)
- // utilize local RocksDB instance on the OSD
 - cls_cxx_map_setvals(ctx, map<str, buf>)
 - cls_cxx_map_getval(ctx, key, buf)
 - cls_cxx_map_getvals(ctx, keystart, nkeys, map<str,buf>, more)

2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

Notice the use of offset/length allows partial read/write of objects SD@

 Enables great flexibility for physical data layout within each object



SD@

 Enables great flexibility for physical data layout within each object



SD₂₀

- Enables great flexibility for physical data layout within each object
- RocksDB enables query-able metadata



SD @

- Enables great flexibility for physical data layout within each object
- RocksDB enables query-able metadata



Example Use Case

Example Use-case: Custom cls Write

SD@

• Write original image, create thumbnails during write



Example Use-case: Custom cls Metadata

 Create metadata – generate image labels, store as metadata in local RocksDB



Example Use-case: Custom cls Read

SD@

Filter data by label="XYZ"



SD@

Create and Register Class/Method

C++ snippet

|--|

CLS_NAME(tabular)

cls_handle_t h_class; cls_method_handle_t h_exec_query_op; cls_method_handle_t h_build_index;

void __cls_init()

CLS_LOG(20, "Loaded tabular class!"); cls_register("tabular", &h_class);

cls_register_cxx_method(h_class, "exec_query_op", CLS_METHOD_RD, exec_query_op, &h_exec_query_op);

 {

// contains the serialized user request.
query_op op;

// decode the query op to get the query params bufferlist::const_iterator it = in->begin(); ceph::decode(op, it);

...

ceph::bufferlist buf;

// lookup metadata
cls_cxx_map_getval(hctx, key, buf)

...

// read local data
int ret = cls_cxx_read(hctx, off, len, &buf);

// process data

...

out->append(result, sizeof(result))

return 0;

'cls' for SkyhookDM

- Note that CLS mechanism *already exists* in Ceph
 - Used heavily by Ceph internals as shown
- We create custom read/write methods
 - Our methods are *not* Ceph specific
 - C++ code, Arrow library
- We simply utilize Chunk store and KV store interfaces
 - Approach is applicable to any system that offers such interfaces for objects

SkyhookDM Architecture



SkyhookDM Architecture

SD@



Data Management Application

20

- Client-side interface to SkyhookDM's LIBRADOS object classes
- Can consider several approaches
 - Distributed processing application frameworks
 - Spark, Dask, others
 - Database External Table interface (widely avail)
 - e.g., PostgreSQL foreign data wrapper
 - FileAPIs that map onto themselves/pass thru
 - e.g., HDF5 Virtual Object Layer (VOL)





Data Format, Partitioning, Access Methods

Data Formats in SkyhookDM

- Utilize fast in-memory serialization formats
 - Google Flatbuffers/<u>Flexbuffers</u>
 - Fields can be accessed without parsing / copying / object allocation.

20

- Apache Arrow Provide "Arrow Native" storage + processing!
 - Very popular in big data world and for data exchange
 - Recent stable release of version 1.0 (July 24, 2020)
 - Compute API for Arrow tables
 - Recent *Dataset API*, provides table abstraction over a collection of remote files/fragments

Partitioning in SkyhookDM



Vertical (col) Partitioning

2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

SD@

Partitioning in SkyhookDM

SD@



Horizontal (row) Partitioning

*uses JumpConsistentHash

Partitioning in SkyhookDM



Horizontal (row) Partitioning

*uses JumpConsistentHash

2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

Key properties of partitions

SD₂₀

- Format retains data's semantics (data schema)
- Object names are generated
- Objects are distributed by Ceph based on name
- Object location not stored by Skyhook

Data Processing

SD₂₀

• SELECT, PROJECT, Aggregate, Groupby, Sort



Moving to Arrow Compute API, will support more

What Metadata to Store in RocksDB?

SD @

- Physical offsets
- Logical content location
 - Create index on various columns
 - Can consider text indexing as well
- Column statistics, access patterns
 - Value distribution important for query optimization
 - Estimate selectivity, then choose scan or index
- Object-local metadata, so each object can optimize itself
 2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

When to Pushdown (offload) Processing

- Currently this is binary to pushdown into storage or not
- Can be a runtime decision by query optimizer based on cluster knowledge
- What if a storage server is overloaded?
 - Object may reject processing, "pushback"
 - We are working on this mechanism
 - Key consideration when offloading in our framework

Physical Design Optimizations

Physical Design

- Long-studied problem in databases
- Good physical design is crucial for workload performance
- Includes partitioning, data format, data layout, indexing, materialized views
- In our case
 - Map table data structure to physical layout on disk

Data Layouts within Object

SD@

Consider Arrow table format





Data Layouts within Object

Consider Arrow table format





SD@

Data Layouts within Object

Consider Arrow table format





SD@



Experimental Results

Experimental Setup

- Data: TPC-H Lineitem table 750M rows
 - 10,000 objects, 75,000 rows/object
- Formats: Flatbuffer/Flexbuffer (row), Arrow (col)
- Queries: select and project
 - SELECT * FROM lineitem WHERE extended_price > 91,500.00
 - SELECT extended_price FROM lineitem
- Hardware: NSF Cloudlab 40 core, 10GbE, 1TB HDD
- Ceph with SkyhookDM extensions, 8 OSDs, 1 client machine
- Simple client side driver, process in client or pushdown to storage using SkyhookDM extensions



BASELINE Read cost only, no query



SD[®]



SD@

SD@



2020 Storage Developer Conference. © CROSS@UCSC

Client-side Server-side (stacked) Client machine - standard standard read - no processing osd7 **Receive all data** osd6 osd5 osd4 osd3 BASELINE osd2 NT CPU PERCENT 10 10 osd1 osd0 Read cost only, no query Timostor Timostor Receive + process Client machine - without without pushdown processing osd7 all data osd6 osd5 osd4 Ehzo osd2 osd1 osd0 CPU Timestep Timester Client machine - with pushdown processing Storage machines - with pushdown processing osd7 osd6 **Receive processed** osd5 osd4 osd3 data ENT osd2 osd1 PER(osd0 U PER Ы 10 50 40

SD@

SELECT 1% No offload to storage

SELECT 1% Offload to storage

2020 Storage Developer Conference. © CROSS@UCSC



Ongoing Work

- Adapting to Apache Arrow Dataset API
 - Interacting with Arrow community for feedback
 - Creating a LIBRADOS Fragment, hopefully push upstream (just starting implementation now)
- RADOS read from remote object (collect)
- Deployment now via Kubernetes and Rook

PING CRESS CENTE OPEN

CENTER FOR RESEARCH IN OPEN SOURCE SOFTWARE

- Bridges gap between student research & open source projects
- Funded by endowment from Sage Weil (Ceph founder) & corporate memberships.
 - Fujitsu Laboratories, Kioxia, Seagate
- Supports graduate research & Incubates work beyond graduation to reach critical mass
 - Skyhookdm project (skyhookdm.com) store & manage tabular data in Ceph
 - Popper project (getpopper.io) container native workflow execution engine
- Directed by Carlos Maltzahn <u>carlosm@ucsc.edu</u>
- <u>cross.ucsc.edu</u> for more information



Summing it up - SkyhookDM

20

- Data partitioning and layout
 - Physical mapping of data onto objects
- Offload processing
 - Custom `cls' methods to execute query ops
- Offload physical design
 - Format transformations, indexing, and query-able metadata
- Not hacking Ceph, just using existing mechanism

• CLS methods updated by copying .so file to OSDs 2020 Storage Developer Conference. © CROSS@UCSC. All Rights Reserved.

Acknowledgements

- Center for Research in Open Source Software at UCSC (CROSS)
- NSF Grant OAC-1836650, CNS-1764102, CNS-1705021
- IRIS-HEP Software Institute
- Current and previous CROSS Corporate Member companies
- Everyone who has contributed to SkyhookDM project!
 - Esp. Noah Watkins, Michael Sevilla, Ivo Jimenez, Ken lizawa
 - Many internal and external students, Google Summer of Code fellows, IRIS-HEP fellows, Master's projects and theses
- Thank you!





CENTER FOR RESEARCH IN OPEN SOURCE SOFTWARE

20

Please take a moment to rate this session.

Your feedback matters to us.







2020 Storage Deve

ELAPSED TIME (seconds)

ELAPSED TIME (seconds)





SD@

Average execution time for point query (unique record), client-side vs. serverside processing (pushdown)

Cloudlab c220g2 machines, SSDs, 1 Billion rows, 10K objects, 1x replication, cold cache



Number of OSDs