



Unlock Bigdata Analytic Efficiency With Ceph Data Lake

Jian Zhang, Yong Fu,

March, 2018



Ceph Agenda

- Background & Motivations
- The Workloads, Reference Architecture Evolution and Performance Optimization
- Performance Comparison with Remote HDFS
- Summary & Next Step





BACKGROUND AND MOTIVATION





BOUNDED Storage and Compute resources on Hadoop Nodes brings challenges



Source: 451 Research, Voice of the Enterprise: Storage Q4 2015

*Other names and brands may be claimed as the property of others.



© Ceph Options To Address The Challenges

Large Cluster

- Lacks isolation noisy neighbors hinder SLAs
- Lacks elasticity rigid cluster size
- Can't scale compute/storage costs separately

More Clusters

 Cost of duplicating datasets across clusters

- Lacks on-demand provisioning
- Can't scale compute/storage costs separately

Compute and Storage Disaggregation

- Isolation of highpriority workloads
- Shared big datasets
- On-demand provisioning
- compute/storage costs scale separately

Compute and Storage disaggregation provides Simplicity, Elasticity, Isolation



Ceph+≣社区 Ceph Unified Hadoop* File System and API for cloud storage



<u>H</u>adoop <u>C</u>ompatible <u>F</u>ile <u>System</u> abstraction layer: Unified storage API interface Hadoop fs –ls s3a://job/





*Other names and brands may be claimed as the property of others.





THE WORKLOADS, REFERENCE ARCHITECTURE AND PERFORMANCE OVERVIEW



Workloads

ceph

Simple Read/Write

- DFSIO: TestDFSIO is the canonical example of a benchmark that attempts to measure the Storage's capacity for reading and writing bulk data.
- Terasort: a popular benchmark that measures the amount of time to sort one terabyte of randomly distributed data on a given computer system.

Data Transformation

• ETL: Taking data as it is originally generated and transforming it to a format (Parquet, ORC) that more tuned for analytical workloads.

Batch Analytics

- To consistently executing analytical process to process large set of data.
- Leveraging 54 derived from TPC-DS * queries with intensive reads across objects in different buckets



ceph Bigdata on Object Storage Performance overview --Batch analytics



- Significant performance improvement from Hadoop 2.7.3/Spark 2.1.1 to Hadoop 2.8.1/Spark 2.2.0 (improvement in s3a)
- Batch analytics performance of 10-node Intel AFA is almost on-par with 60-node HDD cluster





*Other names and brands may b claimed as the property of others



(?) ceph Improve Query Success Ratio with trouble-shootings







- 100% selected TPC-DS query passed with tunings
- Improper Default configuration
 - small capacity size,
 - wrong middleware configuration
 - improper Hadoop/Spark configuration for different size and format data issues



Optimizing HTTP Requests ceph -- The bottlenecks ESTAB

| | The Dottlenecks | ESTAB ESTAB | 0 0 | 0 0 | ::ffff:10.0.2.36:44448 ::ffff:10.0.2.36:44338 | ::ffff:10.0.2.254:80 ::ffff:10.0.2.254:80 |
|-----|--|-------------------|--------|--------|--|--|
| | | ESTAB | 0 | 0 | | |
| | Task Deserialization Time Shuffle Write Time | ESTAB | 0 | 0 | | |
| | Shuffle Read Time Result Serialization Time | ESTAB | 0 | 480 | ::1111:10.0.2.36:44450 | ::111:10.0.2.254:80 |
| | | ESTAD | 0 | 0 | | ··##+10.0.2.2E4:90 |
| | 1 / d> ¬t02 | ESTAD | 0 | 0 | ··ffff-10.0.2.30.44442 | |
| | | ESTAD | 0 | 0 | ··ffff-10.0.2.30.44390 | |
| | Compute tim | e ESTAB | 0 | 0 | | ··ffff:10.0.2.254.80 |
| | take the big | ESTAB | 0 | 0 | ··ffff·10.0.2.30.44432 | ··ffff:10.0.2.254:80 |
| | no tate of the second s | ESTAB | 0 | 0 | ··ffff·10.0.2.36·44444 | ··ffff:10.0.2.254:80 |
| R | eturn 500 | ESTAB 0 | 0 | fi | ff·10 0 2 36·44456 | ··ffff·10.0.2.254.80 |
| | (compute time | E LOTAD U | | econds | interval | |
| | read data +sc | ESTAB | 0 | 0 | "ffff:10.0.2.36:44508 | ···ffff·10 0 2 254·80 |
| Rad | | ESTAB | ő | Ő | "ffff:10.0.2.36:44476 | ··ffff·10.0.2.254:80 |
| | | ESTAB | õ | õ | ··ffff·10.0.2.36·44524 | ··ffff·10.0.2.254·80 |
| | | ESTAB | õ | õ | ··ffff·10.0.2.36·44374 | ··ffff·10.0.2.254·80 |
| | | FSTAR | Ő | õ | ··ffff·10.0.2.36·44500 | ··ffff·10.0.2.254·80 |
| | 20 / alient05 * | | - | 0 | ::ffff:10.0.2.36:44504 | ::ffff:10.0.2.254:80 |
| | New connection | ons out every tim | e, | 0 | ::ffff:10.0.2.36:44512 | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.259976 7fdau **c700 1 ===== starting new request req=0x7fddd6 | on not reused | | 0 | ::ffff:10.0.2.36:44506 | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.271829 7fddd5ffb, 1 ===== starting new request req=0x7fddd5ff5710 ===== | ESTAB | 0 | 0 | ::ffff:10.0.2.36:44464 | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.273940 7fddd7fff700 0 ERROR: flush_read_list(): d->client_c->handle_data() returned - | ESTAB | 0 | 0 | ::ffff:10.0.2.36:44518 | ::ffff:10.0.2.254:80 |
| | 5 | ESTAB | 0 | 0 | ::ffff:10.0.2.36:44510 | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.274223 7fddd7fff700 0 WARNING: set_req_state_err err_no=5 resorting to 500 | ESTAB | 0 | 0 | ::ffff:10.0.2.36:44442 | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.274253 7fddd7fff700 0 ERROR: s->cio->send_content_length() returned err=-5 | ESTAB | 0 | 0 | ::ffff:10.0.2.3 <mark>6:44526</mark> | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.274257 7fddd7fff700 0 ERROR: s->cio->print() returned err=-5 | ESTAB | 0 | 0 | ::ffff:10.0.2.3 <mark>6:44472</mark> | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.274258 7fddd7fff700 0 ERROR: STREAM_IO(s)->print() returned err=-5 | ESTAB | 0 | 0 | ::ffff:10.0.2.36:44466 | ::ffff:10.0.2.254:80 |
| | 2017-07-18 14:53:52.274267 7fddd7fff700 0 ERROR: STREAM_IO(s)->complete_header() returned err=-5 | | | | | |
| | | | | | | |

Unresued connection cause high read time and block performance



::ffff:10.0.2.254:80

::ffff:10.0.2.254:80

::ffff:10.0.2.254:80

::ffff:10.0.2.254:80

::ffff:10.0.2.36:44446

::ffff:10.0.2.36:44454

::ffff:10.0.2.36:44374

::ffff:10.0.2 36:44436

ESTAB

ESTAB

ESTAB

0 0

0 0

0 0

1597240



Background

- The S3A filesystem client supports the notion of input policies, similar to that of the POSIX fadvise() API call. This tunes the behavior of the S3A client to optimize HTTP GET requests for various use cases. To optimize HTTP GET requests, you can take advantage of the S3A experimental input policy fs.s3a.experimental.input.fadvise.
- Ticket: <u>https://issues.apache.org/jira/browse/HADOOP-13203</u>



Solution

Enable random read policy in core-site.xml:

<property>

<name>fs.s3a.experimental.input.fadvise</name>

<value>random</value>

</property>

<property>

<name>fs.s3a.readahead.range</name>

<value>64K</value>

</property>

 By reducing the cost of closing existing HTTP requests, this is highly efficient for file IO accessing a binary file through a series of `PositionedReadable.read()` and `PositionedReadable.readFully()` calls.







Optimizing HTTP Requests -- Performance



- Readahead feature support from Hadoop 2.8.1, but not enabled by default. Apply random read policy, 500 issue gone and performance improved 3x than before
- All Flash storage architecture also show great performance benefit and low TCO which compared with HDD storage



Ceph Optimizing HTTP Requests --Resource Utilization Comparison





ceph Resolving RGW BW bottleneck --The bottlenecks



- LB BW has became the bottleneck
- Observed many messages blocked at load balancer server(send to s3a driver), but not much blocked at receiving on s3a driver side





Hardware Configuration

ceph -- No LB with more RGWs and round-robin DNS



*Other names and brands may be claimed as the property of others.





- 18% performance improvement with more RGWs and round-robin DNS
- Query42(has less shuffle) is 1.64x faster in the new architecture





Key Resource Utilization Comparisor



 Compute side(Hadoop s3a driver) can read more data from OSD, which represent DNS deployment really can gain network throughput performance than single gateway with bonding/teaming technology









ceph Performance evaluation with RGW & OSD collocated



- No need extra dedicate RGW servers, RGW instance and OSD go through different network interface by enable ECMP
- No performance degradation, but more less TCO







- Scale out RGWs can improve performance before OSD(storage) saturating
- So How many RGWs can win the best performance should be decided by the bandwidth of each RGW server and throughput of OSDs









- Mount two RBDs on compute node remotely instead of physical shuffle device
- Ideally, the latency on remote RBDs larger than local physical device, but the bandwidth of remote RBD is not smaller than local physical device too
- So final performance of a TPC-DS query set on RBDs maybe close or even better than on physical device while there are heavy shuffles





PERFORMANCE COMPARISON WITH REMOTE HDFS





Ceph Hardware Configuration --Remote HDFS





ceph Bigdata on Cloud vs. Remote HDFS --Batch Analytics



- On-par performance compared with remote HDFS
 - With optimizations, bigdata analytics on object storage is onpar with remote, especially on parquet format data
 - performance of s3a driver close to native dfsclient, and demonstrate compute and storage separate solution has a considerable performance compare with combination solution



Ceph Bigdata on Cloud vs. Remote HDFS --DFSIO





Ceph Bigdata on Cloud vs. Remote HDFS --Terasort



| Job Name: | TeraSort | Time cost at |
|----------------------|------------------------------|-----------------|
| User Name: | root | Reduce stade is |
| Queue: | root.root | |
| State: | SUCCEEDED | big part |
| Uberized: | false | |
| Submitted: | Thu Nov 02 09:38:18 CST 2017 | |
| Started: | Thu Nov 02 09:38:50 CST 2017 | |
| Finished: | Thu Nov 02 10:24:56 CST 2017 | |
| Elapsed: | 46mins, 5sec | |
| Diagnostics: | | |
| Average Map Time | 1mins, 37sec | |
| Average Shuffle Time | 8mins, 53sec | |
| Average Merge Time | 30sec | Pood and write |
| Average Reduce Time | 23mins, 35sec | Reau and write |
| | | concurrently |





ceph Bigdata on Cloud vs. Remote HDFS --Ongoing rename optimizations

| DirectOutputCom mitter | An implementation in Spark 1.6, that return the destination address as working directory then no need to rename/move task output, no good robustness for failures, removed in Spark 2.0 |
|-------------------------------|---|
| IBM's "Stocator" committer | Targets Openstack Swift, good robustness, but it is another file system for s3a |
| Staging committer | A choice of new s3a committer*, need large capacity of hard disk for staged data |
| Magic committer | A choice of new s3a committer*, if you know your object store is consistent or use s3gurad, this committer has higher performance |

Rename operation can be improved!

* New s3a committer has merged into trunk, and will release in Hadoop 3.1 or later





SUMMARY & NEXT STEP



Summary and Next Step

Summary

- Bigdata on Ceph data lake is functionality ready validated by industry standard decision making workloads TPC-DS
- Bigdata on the Cloud delivers on-par performance with remote HDFS for batch analytics, intensive write operations still need further optimizations
- All flash solutions demonstrated significant TCO benefit compared with HDD solutions

Next

- Expand analytic workloads scope
- Rename operations optimizations to improve the performance
- Accelerating the performance with speed up layer for shuffle





Q&A





BACKUP



Ceph Experiment environment

| Cluster | Hadoop head | Hadoop slave | Load balancer | OSD | RGW |
|-----------|--|--|------------------------------|---|---|
| Roles | Hadoop name node Secondary name node Resource manager Data node Node manager Hive metastore service Yarn history server Spark history server Presto server | Data node Node manager Presto server | Haproxy | Ceph osd | Ceph rados gateway |
| # of node | 1 | 5 | 1 | 5 | 5 |
| Processor | Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz 44 cores HT enabled | | | | Intel(R) Xeon(R) CPU E31280 @ 3.50GH 4 cores HT enabled |
| Memory | 128 | BGB | 64GB | 128GB | 32GB |
| Storage | 4x 1TE 2x Intel S3510 480GB | 3 HDD SSD(vs s3700 metrics) | 1x Intel S3510 480 GB SSD | 1x Intel® P3700 1.6TB as jounal 4x 1.6TB Intel® SSD DC S3510 2X 400GB s370 as data store | 1x Intel S3510 480 GB SSD |
| Network | 10 | GB | 40GB | 10GB+10GB | 10GB |



| SW Configuration | | | |
|------------------|-------------|--|--|
| Hadoop version | 2.7.3/2.8.1 | | |
| Spark version | 2.1.1/2.2.0 | | |
| Hive version | 2.2.1 | | |
| Presto version | 0.177 | | |
| Executor memory | 22GB | | |
| Executor cores | 5 | | |
| # of executor | 24 | | |
| JDK version | 1.8.0_131 | | |
| Memory.overhead | 5GB | | |



S3A Key Performance Configuration

| fs.s3a.connection.maximu m | 10 |
|---------------------------------------|-------|
| fs.s3a.threads.max | 30 |
| fs.s3a.socket.send.buffer | 8192 |
| fs.s3a.socket.recv.buffer | 8192 |
| fs.s3a.threads.keepalivetim e | 60 |
| fs.s3a.max.total.tasks | 1000 |
| fs.s3a.multipart.size | 100M |
| fs.s3a.block.size | 32M |
| fs.s3a.readahead.range | 64k |
| fs.s3a.fast.upload | true |
| fs.s3a.fast.upload.buffer | array |
| fs.s3a.fast.upload.active.bl ocks | 4 |
| fs.s3a.experimental.input.f advise | radom |



(?)

Legal Disclaimer & Optimization Notice

ceph

- Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit www.intel.com/benchmarks.
- INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.
- Copyright © 2018, Intel Corporation. All rights reserved. Intel, Pentium, Xeon, Xeon Phi, Core, VTune, Cilk, and the Intel logo are trademarks of Intel Corporation in the U.S. and other countries.

Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804