

One Ring Cannot Rule Them All

Gary Ogasawara Cloudian Inc.



Motivation



Why Object Storage?

To the application user, the logical "object" matters, Not how it's physically stored (e.g., pieces, versions, location).





How Object Storage?



To optimize overall system need to be able to have finegrained and dynamic control inputs.



HyperStore Policy Engine

Store and tier objects intelligently using rules and current information

Dynamic storage policies optimized for each object.

Cassandra vs. Replicas vs. Erasure Coding Consistency Level Quality of Service Tiering Policy



System Overview



System Overview

- <u>S3 Compatible</u> Basic and advanced Amazon S3 features implemented.
- 2. Optimized for any workload
- **3.** <u>Small Start/Elastic:</u> True P2P Architecture, 10TB to Exabytes
- **4.** <u>**Reliable**</u>: True Multi-DC, Dynamic Consistency, AWS Location Constraint
- 5. <u>Turnkey:</u> Rich, flexible management features (QoS, Reporting, Chargeback, Admin GUI & API)



loudian					admin@cloudian.co	om (SystemA	idimin) (Log out
localan			Admin	Data Explorer	Account.	Report	System
ect System Category	Node Status (up	odated each minute) - Aug	- 19-2014 05:29 PM				
ystem Status		Host:	snalk-pc2 •	9			
ode Status			ononi per				
lode Management		Disk Used %			CPU Utilization	n %	
ystem Info					-1-		
otification Rules		41			1		
		114.8 GB/279.3 GB					
	Disk Reats bytes	114.8 GB(279.3 GB	0 Disi	. Writes bytes/sec		25	314
	Disk Reads bytes Network Status	114.8 08/279.3 09	0 Dist	Writes bytes/sec		25	314
	Disk Reads lytes Network Status Network Status	114.8 GB(279.3 GB	0 Disi	Writes bytes/sec	Transmitted:	25	11338
	Disk Reads byter Network Status Network bytes sy Transactions liser Beijung Through	114.9 08/279.3 08	0 Dist	Writes bytes/sec 11053 0.0	Transmitted. Put	25	314 11338 0.0
	Disk Reads byter Network Status Network bytes in Transciencisses Request Through Average Request	114.8 GB/279.2 GB issec set: 2 guit: Latency (ms):	0 Disi Received. Oet Oet	11053 0.0 0.0	Transmilled Put Put Put	25	314 11338 0.0 0 0
	Disk Reads løter Network Status Notwork bytesis Transactionsken Request Through Average Request Services Status	114.8 OB/279.3 OB	0 Disi Received Oet Oet	11063 0. 0.	Transmitted Put Put Put	25	314 11338 0.0 0 0
	Disk Reads løter Network Status Notwork bytesis Transactionsken Reguest Through Average Reguest Services Status Service	114.8 OB/279.3 OB	0 Disi Received Oct Oct Oct	11053 0.0 0 0	Transmitted Put Put Put Status	25	314 11338 0.0 0 0



This talk

High-level system view



Object Storage Cluster



Elastic, distributed and reliable



9

Cloudian

Logical architecture



© 2014 Cloudian, Inc. All rights reserved.

Cloudian

Multi data center



Single-Data Center Single Region Multi-Data Center Single Region Multi-Data Center Multi Region



HyperStore: Storage Type



Policy: Storage Type

- Policies tailored for different object types
- Large object support
- 30% faster reads
- 400% faster writes
- 4x more TPS
- 1.6x better disk utilization





Average latency vs. object size



Optimal storage type selection depends on object size.



Better performance with less variance





Replication





Erasure coding





Configurable

 Storage type per-bucket and per-object: {EC, Replicas, Cassandra}

Cloudian

- Number of data & coding fragments: k and m.
- Consistency level: {QUORUM(k), ALL}
- Fragment placement on unique nodes: {T, F}

Large object support

- Chunking
 - Break single objects into smaller chunks when storing
 - Distribute chunks across the cluster for better performance
- Multi-part
 - Parts uploaded independently and in any order
 - Single parts can be re-transmitted
 - After all parts are uploaded, then presented as a single object
 - All parts are stored in file system, regardless of threshold configuration



HyperStore: Consistency Level



Consistency Level Example

CL=ALL

After "some" objects are replicated, After "all" objects are replicated, Client receives acknowledgement Client receives acknowledgement Client Client QUORUM(R+W>N) example ALL example Replications=3 Replications=3 (N) (N) (W) Write=3 (W) Write=2 Read=2 Read=any (R) "ack" (R) "ack" Server Server Server Server Server 1 Server 1 Client Client Data is always consistent Data is consistent by QUORUM Cloudian © 2014 Cloudian, Inc. All rights reserved.

CL=QUORUM

Consistency levels

Level	Description
ONE	A write has been written to at least 1 replica's commit log and memory table before responding to the client.
QUORUM	Ensure that the write has been written to N / 2 + 1 replicas before responding to the client.
LOCAL_QUORUM	Similar to QUORUM but replicas are in same data center as coordinator node.
EACH_QUORUM	Ensure that the write has been written to a quorum of replicas in <i>each datacenter in the cluster</i> before responding to the client
ALL	All replicas must have received the write, otherwise the operation will fail.

Level	Description
ONE	Returns the response from the first replica causing a consistency check in a background thread.
QUORUM	Returns the record with the most recent timestamp once $(N/2 + 1)$ replicas has responded.
LOCAL_QUORUM	Returns the record with the most recent timestamp once $(N/2 + 1)$ replicas in the same data center as the coordinator node has reported.
ALL	Returns the record with the most recent timestamp once all replicas have responded. The read operation will fail if a replica does not respond.

READ

WRITE



Configurable consistency across DCs



Multi-DataCenter, Single Region

- 2 DCs, 1 Region, 5 replicas
- # of replicas per DC
- Consistency level: {QUORUM, LOCAL_QUORUM, EACH_QUORUM, ONE, ALL}

Multi-DataCenter, Multi Region

- 2 DCs, 2 Regions, 3 Replicas
- Location constraint at bucket level.
- Objects not shared across regions.
- Span buckets across regions (virtual buckets)
- QoS & Billing based on region
- In each region, standard multi-DC configuration.

Cloudian



Problem: CL=ALL for Multi-DC

- "All" or "QUORUM" cannot be completed in the case of failure
- "All" or "QUORUM" is slower for clients/applications



Problem: CL=LOCAL_QUORUM for Multi-DC

"Local _Quorum" cannot provide "data consistency" in the whole system



Cloudian

Automatic dynamic consistency level



HyperStore: QoS



Quality of Service (QoS)

Operator:

- Prevent a user or group from impacting other users/groups.
- Guarantee a specific level of performance to a user/group.

User/Group:

• Don't exceed a preset limit.



Quality of Service(QoS) Management

- Configurable maximum limits perregion at per-user, per-group, system level.
 - Requests/minute
 - Storage bytes
 - Storage objects
 - Data Bytes Inbound
 - Data Bytes Outbound
- While limit is reached, requests are rejected.



Storage Byte Limiter

Inbound/Outbound Requests Limiter





Limiter Off

Ioudian

Inbound/Outbound Data Byte Limiter

© 2014 Cloudian, Inc. All rights reserved.

28

HyperStore: Tiering



Tiering: Hybrid private/public

- Some data on-premise, some off-premise
- S3 bucket lifecycle policies (e.g., age) to migrate data to Amazon S3 and Glacier (or any S3 system)
- Read options:
 - Streaming
 - HTTP redirect
 - Restore
- Consolidated reports and bill © 2014 Cloudian, Inc. All rights reserved.



Cloudian products

HyperStore Software

- The software only version
- Runs on commodity hardware
- Runs on commodity software: Linux, POSIX filesystem

HyperStore Appliance

- Sold as an appliance by Cloudian or by a Cloudian Partner
- End user gets complete hardware/software solution
- No software installation needed
- 3 Models
 - HSA1024: 1U, 32GB RAM, 4xGigE NIC, 24 TB,
 - HSA1048: 1U, 32GB RAM, 4xGigE NIC, 48 TB
 - HSA2048: 2U, 64GB RAM, 2x1 GigE + 2x10GigE NIC, 48 TB, Flash Optimized









Closing

Optimize overall system by providing fine-grained controls, both manually and automatically changed.

More Info, free trial, demo, PoC:

- www.cloudian.com
- @CloudianStorage, @go10
- www.facebook.com/cloudian.cloudstorage





Backup



About Cloudian

- Object storage startup in Silicon Valley
- Production hardened product
- Target market: mid- to large-enterprises & regional service providers

CLOUDIAN PARTNERS





Objects as a higher layer of abstraction



Block & File vs. Object

Emergence of a two-tier enterprise storage architecture

Faster

- For 'hot' data
- Flash-optimized
- IOPS-centric
- VM/VDI optimized
- Variety of approaches





Bigger

Object

- For cool/cold data
- Object-based
- Scale-out (multi-PB)
- Software-centric
- Cloud-compatible





Object Storage Use Cases

CV7 commvault Enterprise backup SIMPAN CV7 commvault Long term archiving SIMPA Sync and share ctera ctera Remote office file storage Maginatics **On-premise** opens cloudstack Cloud storage (PaaS) **S3**

© 2014 Cloudian, Inc. All rights reserved.

Cloudian

S3 advanced features:

- Multi-part uploads: allows uploading large objects in multiple parts
- Versioning: multiple versions of same object
- Bucket Lifecycle: auto-expiration using rules
- Server side encryption: enhance confidentiality

eneral	Security Logging Requester Pays Versioning	
Ti	Versioning	
🖌 Ver	sioning	

- Location constraint: Assign data to specific region (eg for HIPAA Compliance)
- Bucket Website: Create buckets as websites to host web content
- Access control lists (ACLs) define access rights to bucket and object
- And more...

uthenticated URL		
one		
ccess Control List (AC	L) Permission	
davidykocher	FULL CONTROL	
log Delivery	READ_ACP	:
log Delivery	WRITE	:





