Microsoft Windows Azure Storage, Facebook's Data Storage and Google Search Architecture

Hung-Wei Tseng



Recap: GFS architecture

decoupled data and control paths only control path goes through master





Chunk Server		
Linux FS		
Disk	Disk	
Disk	Disk	
Disk Disk		

load balancing, replicas among chunkservers

Recap: How does GFS achieve its goals?

- Storage based on inexpensive disks that fail frequently master/chunkserver/client — MapReduce is fault tolerant
- Many large files in contrast to small files for personal data
- Primarily reading streams of data large chunk size MapReduce reads chunks of large files
- Sequential writes appending to the end of existing files large chunk size — Output file keep growing as workers keep writing
- Must support multiple concurrent operations flat structure
- Bandwidth is more critical than latency large chunk size

-MapReduce only wants to finish tasks within "reasonable" amount of time

Recap: Why Windows Azure Storage

- Must tolerate many different data abstractions: blobs, tables and queues
- Learning from feedbacks in existing cloud storage
 - Strong consistency
 - Global and scalable namespace/storage
 - Disaster recovery
 - Multi-tenancy and cost of storage



Current scoreboard





Outline

- Windows Azure Storage: A Highly Available Cloud Storage Service with Strong Consistency
- f4: Facebook's Warm BLOB Storage System
- Google Search



GFS v.s. stamp in WAS





What is a stream?

- Regarding a stream in WAS, please identify how many of the following statements is/are true
 - ① A stream is a list of extents, in which an extent consists of consecutive blocks Similar to an extent-base file system. Shares the same benefits with EXT-based systems Each block in the stream contains a checksum to ensure the data integrity (2) As a result, we need to read a whole block every time.... But not a big issue because ... An update to a stream can only be appended to the end of the stream Append only, copy-on-write ... (Doesn't this sound familiar?) Improved bandwidth, data locality (3) Two streams can share the same set of extents (4)
 - Minimize the time when creating a new file De-duplication to save disk space Α.
 - B. 1 C. 2 D. 3



LogFS

GFS v.s. stamp in WAS





Partition layer			
		layer	
	Stream N	<i>l</i> anager	
tent ode	Extent node	Extent node	Extent node

node

node

node

Poll close in 1:30

Why "append-only" and "sealing"?

- In WAS, the stream is append only. The stamp will "seal" extents and extents will become immutable once sealed. How many of the following can sealing contribute to?
 - ① Must tolerate many different data abstractions: blobs, tables and queues
 - ② Strong consistency
 - ③ Global and scalable namespace/storage
 - ④ Disaster recovery
 - Multi-tenancy and cost of storage
 - A. 1
 - B. 2
 - C. 3

D. 4



Poll close in 1:30

Why "append-only" and "sealing"?

- In WAS, the stream is append only. The stamp will "seal" extents and extents will become immutable once sealed. How many of the following can sealing contribute to?
 - ① Must tolerate many different data abstractions: blobs, tables and queues
 - ② Strong consistency
 - ③ Global and scalable namespace/storage
 - ④ Disaster recovery
 - Multi-tenancy and cost of storage
 - A. 1
 - B. 2
 - C. 3

D. 4



Why "append-only" and "sealing"?

- In WAS, the stream is append only. The stamp will "seal" extents and extents will become immutable once sealed. How many of the following can sealing contribute to?
 - ① Must tolerate many different data abstractions: blobs, tables and queues
 - Strong consistency
 - ③ Global and scalable namespace/storage

Multi-tenancy and cost of storage

Ø Disaster recovery

2. Once an extent is sealed, any reads from any sealed replica will always see the same contents of the extent.

Append-only System - Having an append-only system and sealing an extent upon failure have greatly simplified the replication protocol and handling of failure scenarios. In this

F

A. 1



Erasure coding sealed extents is an important optimization, given the amount of data we are storing. It reduces the cost of storing data from three full replicas within a stamp, which is three times the original data, to only 1.3x - 1.5x the original data, depending





Write failure

- Consider the case where 1 of 3 nodes handling a write fails and the current extent is sealed at latest commit boundary (end of extent) — that data will be on failed node
- new extent created
- SM chooses three new replicas to store extents
- client retries via new primary among the three new replicas
- failed node, upon restart, will coord w/ SM to synchronize its extent to the commit length decided upon

GFS v.s. stamp in WAS





Partition layer			
Stream layer			
	Stream N	Janager	
ktent ode	Extent node	Extent node	Extent node
ctent ode	Extent node	Extent node	Extent node

What's the goal of partition layer?

- In WAS, streams are handled and allocated by the partition layer. How many of the following goals in WAS are achieved by the partition layer?
 - ① Must tolerate many different data abstractions: blobs, tables and queues
 - ② Strong consistency
 - Global and scalable namespace/storage (3)
 - Disaster recovery (4)
 - S Multi-tenancy and cost of storage
 - A. 1
 - B. 2
 - C. 3

D. 4



What's the goal of partition layer?

- In WAS, streams are handled and allocated by the partition layer. How many of the following goals in WAS are achieved by the partition layer?
 - ① Must tolerate many different data abstractions: blobs, tables and queues
 - ② Strong consistency
 - Global and scalable namespace/storage (3)
 - Disaster recovery (4)
 - S Multi-tenancy and cost of storage
 - A. 1
 - B. 2
 - C. 3

D. 4



What's the goal of partition layer?

- In WAS, streams are handled and allocated by the partition layer. How many of the following goals in WAS are achieved by the partition layer?
 - Must tolerate many different data abstractions: blobs, tables and queues
 - Strong consistency
 - Global and scalable namespace/storage
 - **Oisaster recovery** inter-stamp replication
 - S Multi-tenancy and cost of storage
 - A. 1
 - B. 2
 - C. 3



partition manager

Partition layer

- Managing high-level data abstractions
- Providing scalable object namespaces
- Providing transaction ordering and strong consistency for objects
- Storing object data on top of the stream layer
- Cache object data to reduce disk I/O

GFS v.s. stamp in WAS





Front-ends

Partition layer	
Stream layer	
Stream Manager	

tent	Extent	E
ode	node	n

ctent ode

Extent node

ent	Extent	Extent	Extent
de	node	node	node

Front-end layer

- A set of stateless servers taking incoming requests
 - Think about the benefits of stateless in NFS
- Keep partition maps to forward the request to the right server
 - A stamp can contain 10—20 racks with 18 disk-heavy storage node per rack
- Stream large objects directly from the stream layer and cache frequently accessed data for efficiency

Are they doing well?



22

Number of VMs



GFS v.s. WAS

	GFS (OSDI 2003)	
File organizations	file chunk block	
System architecture	master chunkserver	
Data updates		
Consistency models	relaxed consistency	
Data formats	files	
Replications	intra-cluster replication	
Usage of nodes	chunk server can perform both	se

WAS (SOSP 2011)

stream extent record

stream manager extent nodes

append only updates

strong consistency

multiple types of objects

geo-replication

parate computation and storage

f4: Facebook's Warm BLOB Storage System

Subramanian Muralidhar, Wyatt Lloyd, Sabyasachi Roy, Cory Hill, Ernest Lin, Weiwen Liu, Satadru Pan, Shiva Shankar, Viswanath Sivakumar, Linpeng Tang, and Sanjeev Kumar.



The original NFS-based FB storage

- Within a data center with high-speed network, the round-trip latency of network accesses is not really a big deal
- However, the amount of metadata, especially directory metadata, is huge — cannot be cached
- As a result, each file access still requires ~ 10 inode/data requests from disks/network nodes — kill performance





http://<CDN>/<Cache>/<Machine ID>/<Logical volume,Photo>

Finding a needle in Haystack: Facebook's photo storage. Doug Beaver, Sanjeev Kumar, Harry C. Li, Jason Sobel, Peter Vajgel. OSDI 2010





- Each storage unit provides 10TB of usable space, using RAID-6 20%redundancy for parity bits
 - Each storage split into 100 physical volumes (100GB)
 - Physical volumes on different machines grouped into logical volumes
 - A photo saved to a logical volume is written to all corresponding physical volumes **3** replicas
- Each volume is actually just a large file
 - Needle represents a photo
 - Each needle is identified through the offset
 - Sealed (the same as WAS) one the file reaches 100GB



Header Magic Number	
Cookie	
Key	
Alternate Key	
Flags	
Size	
Data	
Footer Magic Number	
Data Checksum	
Padding	

GFS v.s. WAS

	GFS (OSDI 2003)	Facebook Haystack (OSDI 2010)
File organizations	file chunk block	volume needle
System architecture	master chunkserver	directory haystack store
Data updates		
Consistency models	relaxed consistency	
Data formats	files	photo/needle
Replications	intra-cluster replication	RIAD-6 & geo-replication
Usage of nodes	chunk server can perform both	

WAS (SOSP 2011)
stream extent record
stream manager extent nodes
append only updates
strong consistency
multiple types of objects
geo-replication
separate computation and storage

Poll close in 1:30

Why FB wants f4 in addition to Haystack

- Regarding the optimization goals of f4, please identify how many the following statements is/are correct.
 - ① f4 is optimized for the throughput of accessing data
 - ② f4 is optimized for the latency of accessing data
 - ③ f4 is designed to reduce the degree of data replications in the system
 - ④ f4 is designed to tolerate various kind of failure in the system
 - A. 0
 - B. 1
 - C. 2

D. 3



Poll close in 1:30

Why FB wants f4 in addition to Haystack

- Regarding the optimization goals of f4, please identify how many the following statements is/are correct.
 - ① f4 is optimized for the throughput of accessing data
 - ② f4 is optimized for the latency of accessing data
 - ③ f4 is designed to reduce the degree of data replications in the system
 - ④ f4 is designed to tolerate various kind of failure in the system
 - A. 0
 - B. 1
 - C. 2

D. 3



Why FB wants f4 in addition to Haystack

- Regarding the optimization goals of f4, please identify how many the following statements is/are correct.
 - ① f4 is optimized for the throughput of accessing data
 - ② f4 is optimized for the latency of accessing data
 - If 4 is designed to reduce the degree of data replications in the system
 - If 4 is designed to tolerate various kind of failure in the system

A. 0 B. 1 D. 3 F. 4

Storage Efficiency One of the key goals of our new system is to improve storage efficiency, i.e., reduce the effective-replication-factor while still maintaining a high degree of reliability and performance.

> Fault Tolerance Another important goal for our storage system is fault tolerance to a hierarchy of faults to ensure we do not lose data and that storage is always available for client requests. We explicitly consider four types of failures:



Poll close in 1:30

What kind of data is f4 optimized for?

- Regarding the type of data that f4 aims at, please identify how many the following statements is/are correct.
 - ① f4 is optimized for most frequently requested data in Facebook services
 - ② f4 is optimized for frequently created, deleted data
 - ③ f4 is optimized for reducing the access latency of long-term storage
 - ④ f4 is optimized for read-only data
 - f4 is optimized for data that are not accessed very frequently
 - A. 1
 - B. 2
 - C. 3

D. 4



Poll close in 1:30

What kind of data is f4 optimized for?

- Regarding the type of data that f4 aims at, please identify how many the following statements is/are correct.
 - ① f4 is optimized for most frequently requested data in Facebook services
 - ② f4 is optimized for frequently created, deleted data
 - ③ f4 is optimized for reducing the access latency of long-term storage
 - ④ f4 is optimized for read-only data
 - f4 is optimized for data that are not accessed very frequently
 - A. 1
 - B. 2
 - C. 3

D. 4



What kind of data is f4 optimized for?

- Regarding the type of data that f4 aims at, please identify how many the following statements is/are correct.
 - ① f4 is optimized for most frequently requested data in Facebook services
 - ② f4 is optimized for frequently created, deleted data
 - ③ f4 is optimized for reducing the access latency of long-term storage
 - If 4 is optimized for read-only data
 - If 4 is optimized for data that are not accessed very frequently.
 - A. 1
 - B. 2
 - C. 3
 - D. 4
 - E. 5



"Temperature" of data





"Temperature" of data

	Hot	Warm
Access Frequency	Most frequent	Less frequent
Pattern	Created often, delete often	Not so frequently read Not so frequently deleted Maybe read-only
Size		65PB in 2014 and growing rapidl





Facebook storage architecture





Storage efficiency

- Reed-Solomon erasure coding
 - Strips: 10GB data + 4GB parity 1.4x space efficiency
 - One volume contains 10 strips
- XOR Geo-replication
 - Use XOR to reduce overhead further (e.g., Azure makes full copies)
 - Block A in DC1 + block B in DC2 -> parity block P in DC3
 - Any two blocks can be used to generate the third
 - 1.5x space efficiency
- 1.4*1.5 = 2.1x space efficiency in total



data center 3

Announcement

- Project revision
 - Allows you to revise your project with 20% of penalty on the unsatisfactory parts/test cases after the first-round of grading (firm deadline 3/11)
 - Say you got only 60% in the first-round, and you fixed everything before 3/11 - you can still get 60% + 80% * 40% = 92%
- iEVAL count as an extra, full-credit reading quiz
- Final contains two parts (each account for 50%)
 - Part 1: unlimited time between 3/12-3/17, open-ended questions
 - Part 2: 80 minute multiple choices/answers questions + two problem sets of comprehensive exam questions

Computer Science & Engineering





