


The Google File System


Sanjay Ghemawat, Howard Gobioff, Shun-Tak Leung
Google
Vijay Kumar



Outline

- Introduction
- Goals
- Architecture
- Operations Supported
- Master Operations
- Performance
- Conclusion


10/21/2005 2



File System for Google

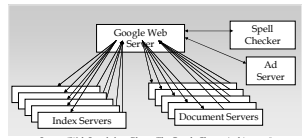
- A new distributed file system developed to meet the demands of Google's application workloads and technological environment
- Shares many of the same goals as the other distributed file systems

10/21/2005 3




Technological environment

- Clusters of more than 15,000 commodity class PCs
- Hardware prone to failure
- Fault tolerant software



Source: "Web Search for a Planet: The Google Cluster Architecture"


10/21/2005 4



Application Workloads

- Sequential reads
 - Indexer reading the contents of web pages
- Frequent Appends
 - Crawler appending new pages
- Files used as Producer - Consumer queues
 - Indexer waits for the crawler to retrieve contents
- Files used for multi-way merging


10/21/2005 5



GFS - Motivation

- Component failures - norm rather than exception
- Efficient management of large files
- Optimization of frequently performed operations
- Flexibility of co-designing application and file system

10/21/2005 6



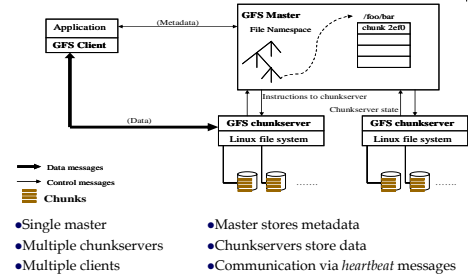
GFS - Goals

- Reliability, availability, scalability...
- Tolerance to hardware failures
- Managing numerous files of large size
- Optimizing commonly performed operations

10/21/2005

7

GFS Architecture



10/21/2005

8

GFS Architecture (contd)

- Files are broken into fixed sized chunks
- Chunks are identified by unique chunk handles
- Chunks are stored in chunkservers
- Fixed chunk size - 64 MB

10/21/2005

9

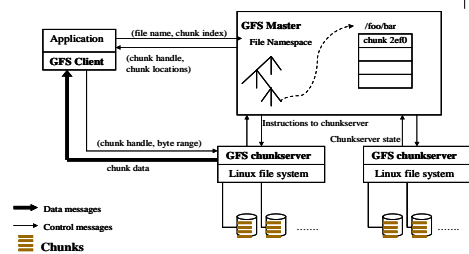
GFS Operations

- File open, close ...
- Data reads
- Data Mutations
 - Random writes
 - Record appends
- Snapshot

10/21/2005

10

Data Read



10/21/2005

11

Data Mutations

- Two kinds of data mutations are supported
 - Random writes
 - Record appends
- Record appends much more frequent than random writes
- Concurrent appends are common
- Leases used to maintain consistent mutation order

10/21/2005

12

Data Mutations (contd)

- Relaxed consistency model to support applications
- Relatively simple and efficient to implement
- GFS guarantees
 - Atomic file namespace mutations
 - State of a file region after a data mutation depends on type of mutation, success or failure, presence or absence of concurrent mutations

10/21/2005

13

Data Mutations (contd)

	Write	Record Append
Serial success	<i>defined</i>	<i>defined</i>
Concurrent successes	<i>consistent</i> but <i>undefined</i>	<i>interspersed with inconsistent</i>
Failure	<i>inconsistent</i>	

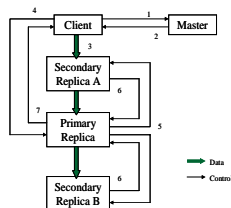
- Consistent file region - All client will see the same data, regardless of the replicas they read from
- Defined file region - It is consistent and clients see what mutations writes in its entirety
- Applications have to deal with the relaxed consistency

10/21/2005

14

Data Writes

1. Client requests for lease holder and secondary replicas
2. Master responds
3. Client pushes out data to the replicas
4. Client issues write instruction to primary replica

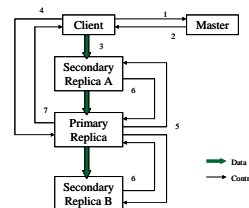


10/21/2005

15

Data Writes (contd)

5. Primary forwards write request to other replicas
6. Secondaries respond to primary reporting the status of the write
7. Primary reports the status to the client



10/21/2005

16

Record Appends

- Appends record to the file atomically (i.e., as a continuous sequence of bytes)
- Appended at an offset chosen by primary
- Pads chunks if append is expected to cross the boundary. Secondaries also do so. New chunk has to be allocated in this case
- Offset returned to the client on success
- Client retries on failure

10/21/2005

17

Snapshot

- Makes a copy of a file or directory tree
- Based on copy-on-write technique
- Minimal overhead involved

10/21/2005

18

Master Operation



- Master is responsible for
 - Metadata management
 - Namespace management
 - Replica management
 - Garbage Collection

10/21/2005

19

Metadata Management



- Master stores the following metadata
 - File and chunk namespaces - stored persistently
 - File to chunk mapping - stored persistently
 - Chunk replica locations - queries chunkserver
- Operation Log - historic record of metadata changes
- Replicated on multiple remote machines
- Size kept small by checkpointing

10/21/2005

20

Name Space Management



- Does not have a per directory data structure
- No support for aliases
- Namespace represented as lookup
- Files and directories have associated *read* and *write* locks

10/21/2005

21

Replica Management



- Chunk replicas created for
 - New chunk creation
 - Chunk re-replication
 - Chunk rebalancing
- Chunk replicas placement based on
 - Maximizing data reliability and availability
 - Maximizing network utilization
- Chunkserver's *disk utilization, count of recent chunks, rack position* affect decisions

10/21/2005

22

Garbage Collection



- Reclamation of physical storage
 - Files renamed to hidden names upon deletion
 - File metadata deleted during namespace scan
 - Physical storage freed during exchange of *Heart Beat* messages
- Stale Replica Detection
 - Stale chunks detected using *chunk version number*

10/21/2005

23

GFS Goals revisited



- Availability
 - Fast recovery
 - Chunk replication
 - Master replication
- Scalability
 - Master replication
 - Keeping master's involvement limited in data transactions

10/21/2005

24

GFS Goals revisited (contd)

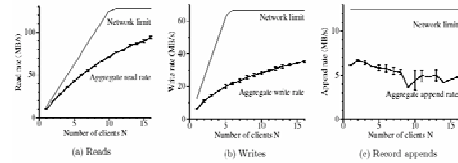
- Fault Tolerance
 - Replication, constant monitoring, fast recovery
- Data Integrity
 - Checksumming used to detect corruption
- Optimization for frequent operations
 - Relaxed consistency model

10/21/2005

25

Performance

- Micro benchmarks
 - Aggregate read rate - 75% of theoretical limit
 - Aggregate write rate - 50% of theoretical limit



10/21/2005

26

Performance (contd)

- Real world clusters

Cluster	A	B	Cluster	A	B
Chomberservers	512	227	Read rate (last minute)	368 MB/s	380 MB/s
Available disk space	72 TB	180 TB	Read rate (last hour)	562 MB/s	384 MB/s
Used disk space	35 TB	155 TB	Read rate (ignore restart)	589 MB/s	49 MB/ops
Number of files	735 k	757 k	Write rate (last minute)	1 MB/s	103 MB/s
Number of Dead files	22 k	232 k	Write rate (last hour)	2 MB/s	117 MB/s
Number of Chunks	002 k	1560 k	Write rate (since restart)	25 MB/s	13 MB/s
Operations at chomberservers	13 CRB	21 CRB	Master ops (last minute)	328 Ops/s	833 Ops/s
Operations at master	48 MB	60 MB	Master ops (last hour)	381 Ops/s	518 Ops/s
			Master ops (since restart)	292 Ops/s	347 Ops/s

10/21/2005

27

Performance (contd)

- Workload breakdown

Operation	Read		Write		Record Append	
	X	Y	X	Y	X	Y
CRB	0.1	2.6	0	0	0	0
1B_1K	0.1	4.1	6.6	4.9	0.2	9.2
1K_8K	65.2	38.5	0.4	1.0	18.9	15.2
8K_64K	29.9	45.1	17.8	43.0	78.0	2.8
64K_128K	0.1	0.7	2.3	1.9	< .1	4.3
128K_256K	0.2	0.3	31.6	0.4	< .1	10.6
256K_512K	0.1	0.1	4.2	7.7	< .1	31.2
512K_1M	3.9	6.9	35.3	28.7	2.2	25.5
1M_inf	0.1	1.8	1.5	12.3	0.7	2.2

Cluster	X	Y
Open	26.1	16.3
Delete	0.7	1.5
FindLocation	64.3	65.8
FindLeaseHolder	7.8	13.4
FindMatchingFiles	0.6	2.2
All other combined	0.5	0.8

10/21/2005

28

Conclusion

- Demonstrates qualities needed to support large scale data processing workloads on commodity hardware
- Delivers high throughput
- Successfully meets Google's storage needs

10/21/2005

29

Questions?

How many times have you used Google?

[button]

10/21/2005

30