

Large-Scale Data and Computation

Yifu Huang

School of Computer Science, Fudan University
huangyifu@fudan.edu.cn

COMP620003 Advanced Computer Networks Report, 2013

Outline

- 1 Motivation
- 2 GFS
- 3 MapReduce
- 4 Bigtable
- 5 Discussion
- 6 BBS Search Engine

Motivation

- Why we choose these papers?
 - Big data
 - Google
- Why we need GFS, MapReduce, Bigtable?
 - A scalable distributed file system for large distributed data-intensive applications
 - A programming model for processing and generating large datasets that is amenable to a broad variety of real-world tasks
 - A distributed storage system for managing structured data that is designed to scale to a very large size
- What can we learn from these papers?
 - The design and implementation ideas behind GFS, MapReduce, Bigtable
 - Get ready to enjoy open source equivalents HDFS, YARN, HBase

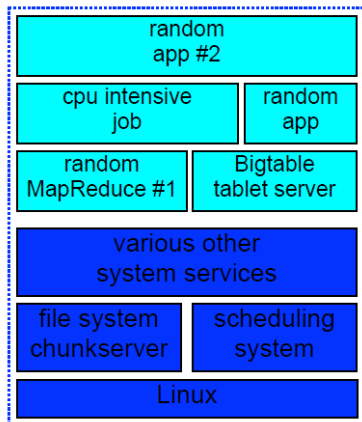
Motivation

- Why we choose these papers?
 - Big data
 - Google
- Why we need GFS, MapReduce, Bigtable?
 - A scalable distributed file system for large distributed data-intensive applications
 - A programming model for processing and generating large datasets that is amenable to a broad variety of real-world tasks
 - A distributed storage system for managing structured data that is designed to scale to a very large size
- What can we learn from these papers?
 - The design and implementation ideas behind GFS, MapReduce, Bigtable
 - Get ready to enjoy open source equivalents HDFS, YARN, HBase

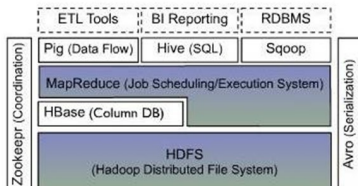
Motivation

- Why we choose these papers?
 - Big data
 - Google
- Why we need GFS, MapReduce, Bigtable?
 - A scalable distributed file system for large distributed data-intensive applications
 - A programming model for processing and generating large datasets that is amenable to a broad variety of real-world tasks
 - A distributed storage system for managing structured data that is designed to scale to a very large size
- What can we learn from these papers?
 - The design and implementation ideas behind GFS, MapReduce, Bigtable
 - Get ready to enjoy open source equivalents HDFS, YARN, HBase

Ecosystem



The Hadoop Ecosystem



The Google File System

The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung

Google, Inc.

Design Overview

- The system is built from many inexpensive commodity components that often fail
- The system stores a modest number of large files
- The workloads primarily consist of two kinds of reads: large streaming reads and small random reads
- The workloads also have many large, sequential writes that append data to files
- The system must efficiently implement well-defined semantics or multiple clients that concurrently append to the same file
- High sustained bandwidth is more important than low latency

Design Overview

- The system is built from many inexpensive commodity components that often fail
- The system stores a modest number of large files
- The workloads primarily consist of two kinds of reads: large streaming reads and small random reads
- The workloads also have many large, sequential writes that append data to files
- The system must efficiently implement well-defined semantics or multiple clients that concurrently append to the same file
- High sustained bandwidth is more important than low latency

Design Overview

- The system is built from many inexpensive commodity components that often fail
- The system stores a modest number of large files
- The workloads primarily consist of two kinds of reads: large streaming reads and small random reads
- The workloads also have many large, sequential writes that append data to files
- The system must efficiently implement well-defined semantics or multiple clients that concurrently append to the same file
- High sustained bandwidth is more important than low latency

Design Overview

- The system is built from many inexpensive commodity components that often fail
- The system stores a modest number of large files
- The workloads primarily consist of two kinds of reads: large streaming reads and small random reads
- The workloads also have many large, sequential writes that append data to files
- The system must efficiently implement well-defined semantics or multiple clients that concurrently append to the same file
- High sustained bandwidth is more important than low latency

Design Overview

- The system is built from many inexpensive commodity components that often fail
- The system stores a modest number of large files
- The workloads primarily consist of two kinds of reads: large streaming reads and small random reads
- The workloads also have many large, sequential writes that append data to files
- The system must efficiently implement well-defined semantics or multiple clients that concurrently append to the same file
- High sustained bandwidth is more important than low latency

Design Overview

- The system is built from many inexpensive commodity components that often fail
- The system stores a modest number of large files
- The workloads primarily consist of two kinds of reads: large streaming reads and small random reads
- The workloads also have many large, sequential writes that append data to files
- The system must efficiently implement well-defined semantics or multiple clients that concurrently append to the same file
- High sustained bandwidth is more important than low latency

Design Overview (cont.)

Architecture

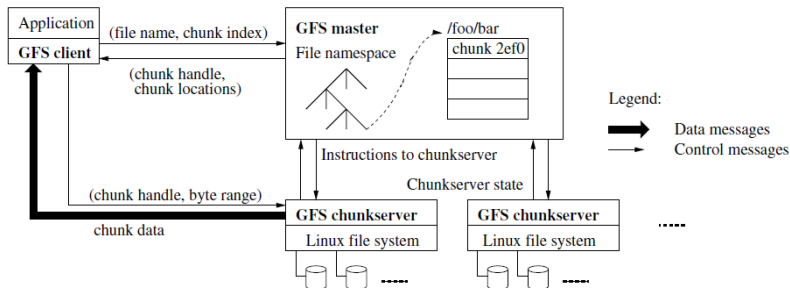


Figure 1: GFS Architecture

Design Overview (cont.)

- **Chunk Size: 64 MB**
- Metadata
 - The file and chunk namespaces
 - The mapping from files to chunks
 - The locations of each chunk's replicas
- Chunk Locations
- Operation Log
- Consistency Model
 - Guarantees by GFS
 - Implications for Applications

Design Overview (cont.)

- Chunk Size: 64 MB
- Metadata
 - The file and chunk namespaces
 - The mapping from files to chunks
 - The locations of each chunk's replicas
- Chunk Locations
- Operation Log
- Consistency Model
 - Guarantees by GFS
 - Implications for Applications

Design Overview (cont.)

- Chunk Size: 64 MB
- Metadata
 - The file and chunk namespaces
 - The mapping from files to chunks
 - The locations of each chunk's replicas
- Chunk Locations
- Operation Log
- Consistency Model
 - Guarantees by GFS
 - Implications for Applications

Design Overview (cont.)

- Chunk Size: 64 MB
- Metadata
 - The file and chunk namespaces
 - The mapping from files to chunks
 - The locations of each chunk's replicas
- Chunk Locations
- Operation Log
- Consistency Model
 - Guarantees by GFS
 - Implications for Applications

Design Overview (cont.)

- Chunk Size: 64 MB
- Metadata
 - The file and chunk namespaces
 - The mapping from files to chunks
 - The locations of each chunk's replicas
- Chunk Locations
- Operation Log
- Consistency Model
 - Guarantees by GFS
 - Implications for Applications

System Interactions

- Leases and Mutation Order

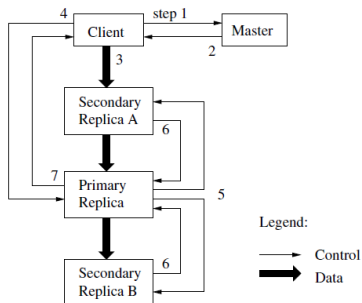


Figure 2: Write Control and Data Flow

System Interactions (cont.)

- Data Flow

- We decouple the flow of data from the flow of control to use the network efficiently
- Our goals are to fully utilize each machine's network bandwidth, avoid network bottlenecks and high-latency links, and minimize the latency to push through all the data

- Atomic Record Appends

- GFS provides an atomic append operation called record append

- Snapshot

- The snapshot operation makes a copy of a file or a directory tree almost instantaneously, while minimizing any interruptions of ongoing mutations

System Interactions (cont.)

- Data Flow

- We decouple the flow of data from the flow of control to use the network efficiently
- Our goals are to fully utilize each machine's network bandwidth, avoid network bottlenecks and high-latency links, and minimize the latency to push through all the data

- Atomic Record Appends

- GFS provides an atomic append operation called record append

- Snapshot

- The snapshot operation makes a copy of a file or a directory tree almost instantaneously, while minimizing any interruptions of ongoing mutations

System Interactions (cont.)

- Data Flow

- We decouple the flow of data from the flow of control to use the network efficiently
- Our goals are to fully utilize each machine's network bandwidth, avoid network bottlenecks and high-latency links, and minimize the latency to push through all the data

- Atomic Record Appends

- GFS provides an atomic append operation called record append

- Snapshot

- The snapshot operation makes a copy of a file or a directory tree almost instantaneously, while minimizing any interruptions of ongoing mutations

Fault Tolerance And Diagnosis

- High Availability
 - Fast Recovery
 - Chunk Replication
 - Master Replication
- Data Integrity
- Diagnostic Tools

Fault Tolerance And Diagnosis

- High Availability
 - Fast Recovery
 - Chunk Replication
 - Master Replication
- Data Integrity
- Diagnostic Tools

Fault Tolerance And Diagnosis

- High Availability
 - Fast Recovery
 - Chunk Replication
 - Master Replication
- Data Integrity
- Diagnostic Tools

MapReduce: Simplified Data Processing on Large Clusters

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

Google, Inc.

Model

- Map

- Takes an input pair and produces a set of intermediate key/value pairs
- The MapReduce library groups together all intermediate values associated with the same intermediate key and passes them to the reduce function

- Reduce

- Accepts an intermediate key and a set of values for that key and merges these values together to form a possibly smaller set of values

Model

- Map

- Takes an input pair and produces a set of intermediate key/value pairs
- The MapReduce library groups together all intermediate values associated with the same intermediate key and passes them to the reduce function

- Reduce

- Accepts an intermediate key and a set of values for that key and merges these values together to form a possibly smaller set of values

Implementation

User Program

Master

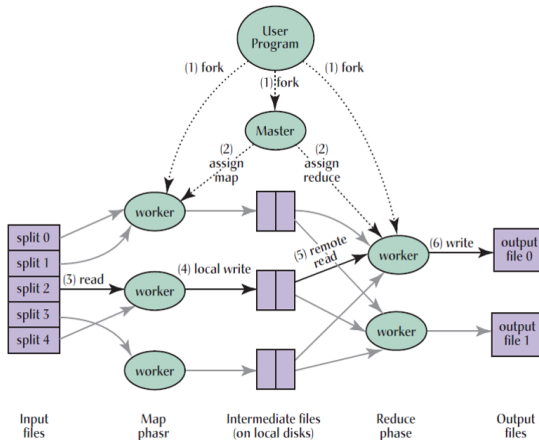
Worker (Map/Reduce)

Fault tolerance

Locality

Task granularity

Backup tools



Performance

Grep

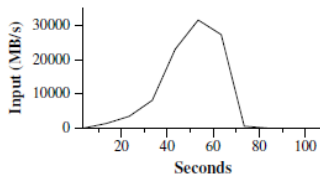
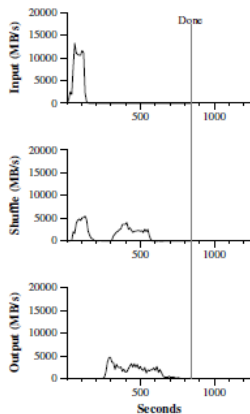


Figure 2: Data transfer rate over time

Performance (cont.)

Sort



(a) Normal execution

Bigtable: A Distributed Storage System for Structured Data

Bigtable: A Distributed Storage System for Structured Data

Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach Mike Burrows, Tushar Chandra, Andrew Fikes, Robert E. Gruber

Google, Inc.

Model

- Rows
 - Atomic
 - Tablets
- Column Families
 - Family : qualifier
 - Access control
- Timestamps
 - Decreasing order
 - Garbage collect
- Cell
 - (row : string, column : string, time : int64) -> string

Model

- Rows
 - Atomic
 - Tablets
- Column Families
 - Family : qualifier
 - Access control
- Timestamps
 - Decreasing order
 - Garbage collect
- Cell
 - (row : string, column : string, time : int64) -> string

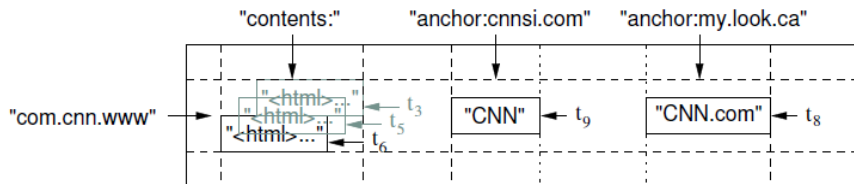
Model

- Rows
 - Atomic
 - Tablets
- Column Families
 - Family : qualifier
 - Access control
- Timestamps
 - Decreasing order
 - Garbage collect
- Cell
 - (row : string, column : string, time : int64) -> string

Model

- Rows
 - Atomic
 - Tablets
- Column Families
 - Family : qualifier
 - Access control
- Timestamps
 - Decreasing order
 - Garbage collect
- Cell
 - (row : string, column : string, time : int64) -> string

Model (cont.)



Storage

Region

–Store

—memStore

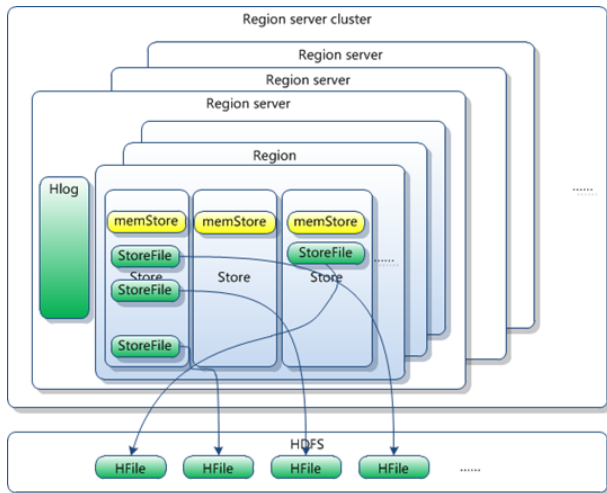
—StoreFile

File

–Blocks

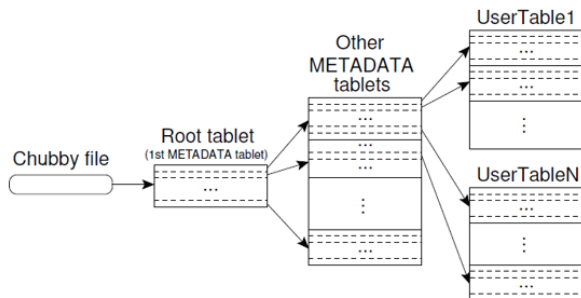
Log

–Write ahead log



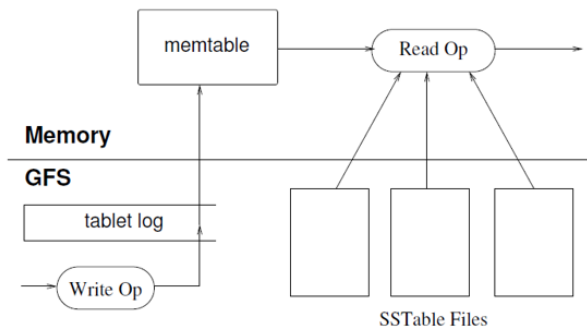
Implementation

- Tablet Location



Implementation (cont.)

- Tablet Serving



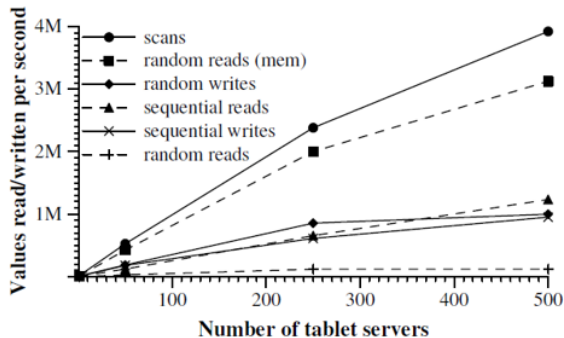
Evaluation

- Benchmarks

Experiment	# of Tablet Servers			
	1	50	250	500
random reads	1212	593	479	241
random reads (mem)	10811	8511	8000	6250
random writes	8850	3745	3425	2000
sequential reads	4425	2463	2625	2469
sequential writes	8547	3623	2451	1905
scans	15385	10526	9524	7843

Evaluation (cont.)

• Scaling



Evaluation (cont.)

- Applications

Project name	Table size (TB)	Compression ratio	# Cells (billions)	# Column Families	# Locality Groups	% in memory	Latency-sensitive?
<i>Crawl</i>	800	11%	1000	16	8	0%	No
<i>Crawl</i>	50	33%	200	2	2	0%	No
<i>Google Analytics</i>	20	29%	10	1	1	0%	Yes
<i>Google Analytics</i>	200	14%	80	1	1	0%	Yes
<i>Google Base</i>	2	31%	10	29	3	15%	Yes
<i>Google Earth</i>	0.5	64%	8	7	2	33%	Yes
<i>Google Earth</i>	70	–	9	8	3	0%	No
<i>Orkut</i>	9	–	0.9	8	5	1%	Yes
<i>Personalized Search</i>	4	47%	6	93	11	5%	Yes

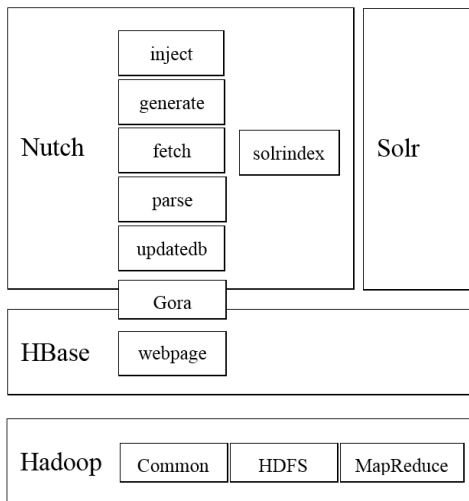
Discussion

- Contributions of GFS, MapReduce, Bigtable
 - Supporting large-scale data processing workloads on commodity hardware
 - Easy to use, hides details of parallelization, fault tolerance, locality optimization, and load balancing
 - Scale well both in terms of data size (from URLs to web pages to satellite imagery) and latency requirements (from backend bulk processing to real-time data serving)
- Drawbacks of GFS, MapReduce, Bigtable
 - Single master may be still a potential bottleneck
 - Disk I/O, time consuming
 - Do not support many database operations (multi-row transaction, secondary index ...)

Discussion

- Contributions of GFS, MapReduce, Bigtable
 - Supporting large-scale data processing workloads on commodity hardware
 - Easy to use, hides details of parallelization, fault tolerance, locality optimization, and load balancing
 - Scale well both in terms of data size (from URLs to web pages to satellite imagery) and latency requirements (from backend bulk processing to real-time data serving)
- Drawbacks of GFS, MapReduce, Bigtable
 - Single master may be still a potential bottleneck
 - Disk I/O, time consuming
 - Do not support many database operations (multi-row transaction, secondary index ...)

Architecture



Architecture (cont.)

- Cluster
- Intel(R) Core(TM)2 Duo CPU E7500 @ 2.93GHz
 - 4GB RAM
 - Hadoop 1.2.1
 - Namenode/Jobtracker: 1, datanode/tasktracker: 24
 - HBase 0.90.4
 - Master: 1, zookeeper: 3, regionserver: 24
- Data
 - [Http://bbs.fudan.edu.cn/bbs/all](http://bbs.fudan.edu.cn/bbs/all)
 - Board: 376
 - Post: 3111945

Architecture (cont.)

- Cluster
- Intel(R) Core(TM)2 Duo CPU E7500 @ 2.93GHz
 - 4GB RAM
 - Hadoop 1.2.1
 - Namenode/Jobtracker: 1, datanode/tasktracker: 24
 - HBase 0.90.4
 - Master: 1, zookeeper: 3, regionserver: 24
- Data
 - [Http://bbs.fudan.edu.cn/bbs/all](http://bbs.fudan.edu.cn/bbs/all)
 - Board: 376
 - Post: 3111945

Demo

10.171.5.222:3000

搜索

日月光华

搜索

© If

Demo (cont.)

搜索

日月光华

王雪平

搜索

phrack La Vie En Rose CS_Graduate 高级网络考试地点 2008年12月26日23:38:27 星期五 请问大家一下高级网络 王雪平的考试在哪个教室啊？多谢～ -- 群 ※ 来源:日月光华 bbs.fudan.edu.cn [FROM: 10.100.129.115]

[原帖 »](#)

brucelee brucelee CS_Graduate Re: 高级网络考试地点 2008年12月28日18:32:24 星期天 【在 phrack 的大作中提到】：请问大家一下高级网络 王雪平的考试在哪个教室啊？：多谢～ 老逸夫楼 502 31号下午 -- ※ 来源:日月光华 bbs.fudan.edu.cn-HTTP [FROM: 10.100.126.96]

[原帖 »](#)

maodingding blueice CS_Graduate Re: 高级网络考试地点 2008年12月29日10:35:43 星期一 【在 phrack 的大作中提到】：请问大家一下高级网络 王雪平的考试在哪个教室啊？：多谢～ 补充：下午1点到3点，并带好网络论文。 -- ※ 来源:日月光华 bbs.fudan.edu.cn-HTTP [FROM: 10.11.12.72]

[原帖 »](#)

phrack La Vie En Rose CS_Graduate Re: 高级网络考试地点 2008年12月29日12:50:06 星期一 多谢！

References I

- [1] The Google File System. SOSP. 2003.
- [2] MapReduce: Simplified Data Processing on Large Clusters. OSDI. 2004.
- [3] Bigtable: A Distributed Storage System for Structured Data. OSDI. 2006.
- [4] The Hadoop Distributed File System. MSST. 2010.
- [5] Hadoop: the definitive guide. 2012.
- [6] HBase: the definitive guide. 2011.
- [7] Data-intensive text processing with MapReduce. 2010.
- [8] The Chubby lock service for loosely-coupled distributed systems. OSDI. 2006.