

IOWA STATE UNIVERSITY
Department of Electrical and Computer Engineering

Lectures 37 - 39: Advanced Topics: Google File System & Beyond



1

Agenda

- Recap
- Advanced Topics
 - The Google File System (GFS)
 - The World After GFS

IOWA STATE UNIVERSITY

2

2

1

Recap

- What's cloud computing
 - Informal:
 - computing with large datacenters
 - users generally do not even know the exact location of “their” resources or even which country they are located in
 - NIST's five essential characteristics
- Types of cloud
 - Infrastructure as a Service (IaaS)
 - Platform as a Service (PaaS)
 - Software as a Service (SaaS)
 - Public vs. private clouds



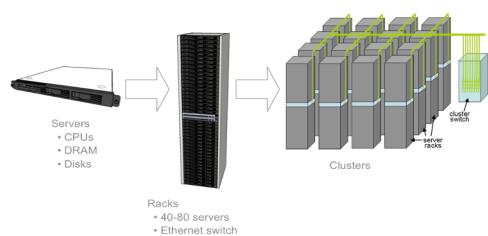
IOWA STATE UNIVERSITY

3

3

Recap

- Cloud Hardware
 - “Warehouse-Scale Computer” (WSC)
 - Datacenter as a computer
 - Rows of rack-mounted servers
 - Often organized as a few mostly independent clusters



IOWA STATE UNIVERSITY

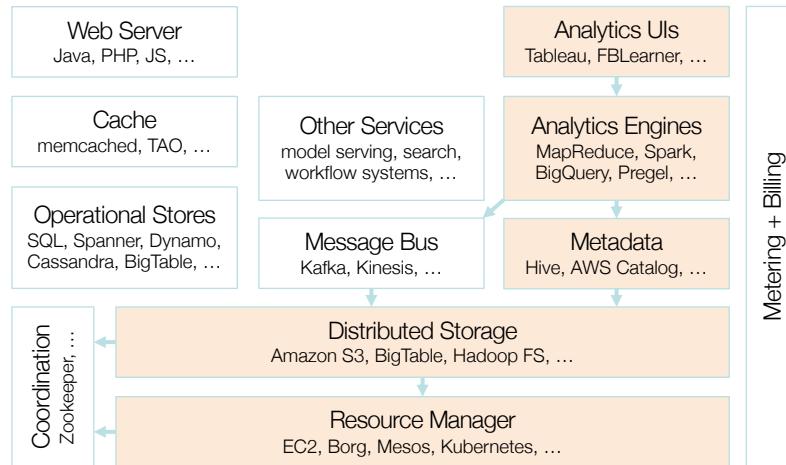
4

4

2

Recap

- Cloud Software



IOWA STATE UNIVERSITY

5

5

Agenda

▪ Recap

- Advanced Topics

- The Google File System (GFS)
- The World After GFS

IOWA STATE UNIVERSITY

6

6

3

Ordinary User's View of Google

- Search engine, Gmail, Google Maps, Ads, ...



IOWA STATE UNIVERSITY

7

7

CPRE 308/563X Students' View of Google

- A wide range of knowledge/techniques/problems/opportunities

Product design	...and much, much more!
User interfaces	
Machine learning, Statistics, Information retrieval, AI	
Data structures, Algorithms	
Compilers, Programming languages	
Networking, Distributed systems, Fault tolerance	
Hardware, Mechanical engineering	

IOWA STATE UNIVERSITY

8

8

Hardware Design Philosophy

- “Scale out” vs. “Scale up”
 - Google prefers low-end server/PC-class designs
 - Build lots of them!
 - Why?
 - Single machine performance is not interesting
 - Even smaller problems at Google are too large for any single system
 - Large problems have lots of available parallelism
 - In contrast to the supercomputer designs in national labs
 - Highly customized/specialized

IOWA STATE UNIVERSITY

9

9

“Google” Circa 1997

- Started by 2 graduate students at Stanford University



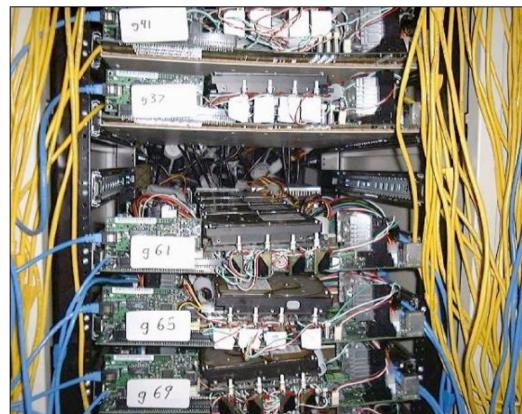
IOWA STATE UNIVERSITY

10

10

Google Circa 1999

- Multiple blade servers



IOWA STATE UNIVERSITY

11

11

Google Circa 2000

- “Data center”



IOWA STATE UNIVERSITY

12

12

Google Circa 2001

- New Data Center
 - in three days ...



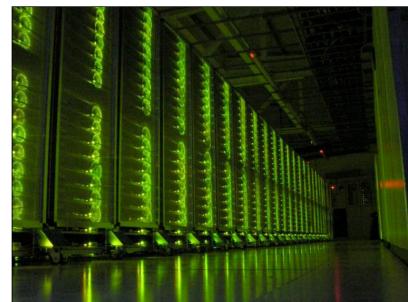
IOWA STATE UNIVERSITY

13

13

Google Today

- Multiple data centers around the world



IOWA STATE UNIVERSITY

14

14

The Joy of Real Hardware at Scale

- Typical first year for a new cluster:
 - ~0.5 overheating (power down most machines in< 5 min, ~1-2 days to recover)
 - ~20 rack failures (40-80 machines instantly disappear)
 - ~1 PDU failure (~500-1000 machines suddenly disappear)
 - ~5 racks go wonky (40-80 machines see 50% packet loss)
 - ~8 network maintenances (4 might cause connectivity losses)
 - ~1000 individual machine failures
 - ~thousands of hard drive failures
 - ~slow disks, bad memory, etc.
 - ...
- Long-haul networking breaks for unusual reasons, too
 - Wild dogs, dead horses, thieves, blasphemy, drunken hunters and sharks



IOWA STATE UNIVERSITY

15

15

Implications

- Stuff Breaks
 - If you have one server, it may stay up 1,000 days
 - If you have 10,000 servers, expect to lose ten a day
- “Ultra-reliable” hardware doesn’t really help
 - At large scales, super-fancy reliable hardware still fails, albeit less often
 - software still needs to be fault-tolerant
 - commodity machines without fancy hardware give better perf/\$

IOWA STATE UNIVERSITY

16

16

The Google File System

- Major Observations & Implications

- Component failures are normal
 - due to both quantity & quality
 - need constant monitoring, error detection, fault tolerance, recovery etc.
- Files are huge by traditional standards
 - many GBs/TBs
 - need to optimize for large I/O sizes
- Most files are modified by appending at the end
 - Random writes (and overwrites) are practically non-existent
 - need to optimize for append operations
 - Caching data at client side becomes less attractive
- Co-design applications and the file system benefits the overall system
 - increasing flexibility for both applications and the FS
 - E.g., may introduce new operation APIs

IOWA STATE UNIVERSITY

19

19

The Google File System

- Interface

- A familiar file system interface
 - Files are organized hierarchically in directories & identified by path names
 - Does not fully follow POSIX standard
- Support typical operations
 - *create, delete, open, close, read, write*
- Introduce non-traditional operations
 - *snapshot*
 - Creates a copy of a file or a directory tree at low cost
 - *record append*
 - Allows multiple clients to append data to the same file concurrently
 - Guarantees atomicity of each individual client's append

IOWA STATE UNIVERSITY

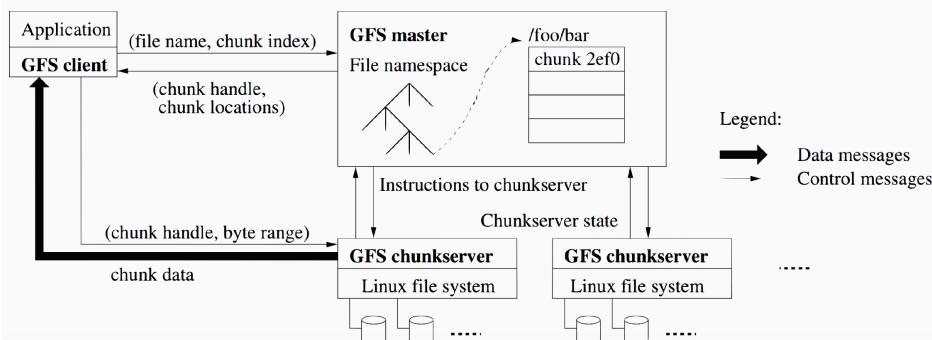
20

20

The Google File System

- Architecture

- three types of nodes
 - master, chunkserver(s), client(s)
 - each is typically a Linux machine running a user-level server process



IOWA STATE UNIVERSITY

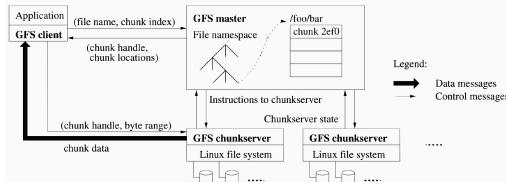
21

21

The Google File System

- Clients

- Applications invoke GFS APIs
- Ask the *master* which *chunkservers* to contact
 - i.e., metadata communication
- Interact with *chunkservers* directly for reading/writing application data
 - i.e., data communication
- Cache metadata
 - Why?
 - Why not data?



IOWA STATE UNIVERSITY

22

22

The Google File System

- Chunkservers
 - store all user files
 - Files are divided in fixed-size chunks
 - 64 MB (why?)
 - each is a Linux file on local file system (e.g., Ext4)
 - each is identified by a 64 bit unique handle in GFS
 - Do not cache file data (at GFS layer)
 - Linux buffer cache already put hot data in memory
 - replication
 - each chunk is replicated on multiple chunkservers
 - 3 by default (replication factor & locations are tunable)
 - tradeoff: space overhead vs. reliability/availability
 - erasure coding: achieve redundancy w/o full replication

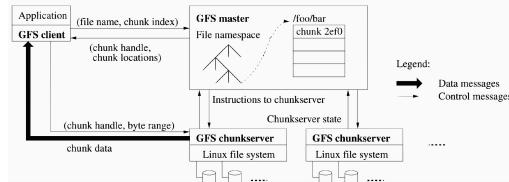
IOWA STATE UNIVERSITY

23

23

The Google File System

- Master
 - Store all GFS metadata
 - file and chunk namespace
 - mapping from files to chunks
 - chunk locations (including replicas)
 - other info
 - E.g., access-control information
 - In-memory
 - Single master
 - why?



IOWA STATE UNIVERSITY

24

24

The Google File System

- Master

- How to know chunk locations?
- Chunkservers may “come and go”
 - Fail, reboot, change name, etc.
 - Hard to sync if keep locations persistently at master



IOWA STATE UNIVERSITY

25

25

The Google File System

- Master

- How to know chunk locations?
- Chunkservers may “come and go”
 - Fail, reboot, change name, etc.
 - Hard to sync if keep locations persistently at master
- Solution
 - Poll chunkservers at startup
 - Keep up-to-date thereafter
 - Controls all chunk placement
 - Monitor chunkservers via periodic heartbeat messages



IOWA STATE UNIVERSITY

26

26

The Google File System

- Master
 - Maintain an operation log persistently
 - Central to GFS
 - A historical record of critical metadata changes
 - the only persistent record of metadata
 - Provide a logical time line that defines the order of all concurrent operations
 - Guarantee crash consistency at the GFS level
 - Store on local disk, replicate on multiple remote machines
 - Respond to a client operation only after flushing the corresponding log record to disk both locally and remotely

IOWA STATE UNIVERSITY

27

27

The Google File System

- Typical Workflow
 - Client
 - using fixed chunk size, translate filename & byte offset to chunk index, send request to master
 - Master
 - replies with chunk handle & location of chunkserver replicas (including which is “primary”)
 - Client
 - cache info using filename & chunk index as key
 - request data from nearest Chunkserver
 - no need to talk more about this 64MB chunk until cached info expires or file reopened
 - often initial request asks about sequence of chunks

IOWA STATE UNIVERSITY

28

28

The Google File System

- Lease and Mutation Order
 - Mutation
 - an operation that changes the contents or metadata of a chunk
 - E.g., write, append
 - Lease
 - A token for maintaining a consistent mutation order across replicas
 - Master grants a chunk “lease” to one replica
 - called “primary” chunkserver
 - the primary picks a serial order for all mutations to the chunk
 - serializes all mutation requests
 - communicates order to replicas

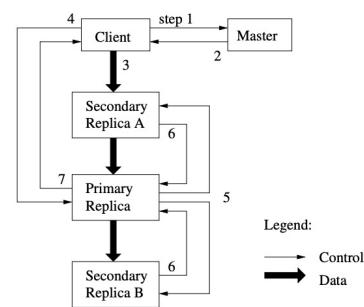
IOWA STATE UNIVERSITY

29

29

The Google File System

- Write Control and Data Flow
 - Step 1:
 - client asks the master which chunkserver holds the current lease for the chunk and the locations of the replicas
 - If no one has a lease, the master grants one to a replica



IOWA STATE UNIVERSITY

30

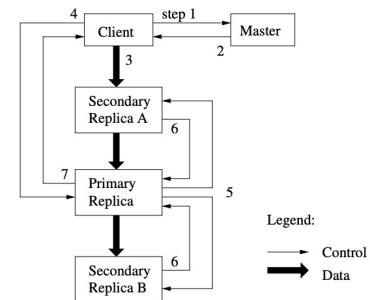
30

The Google File System

- Write Control and Data Flow

- Step 2:

- Master replies with the ID of the primary and secondary replicas
 - Client cache the info for future mutations



IOWA STATE UNIVERSITY

31

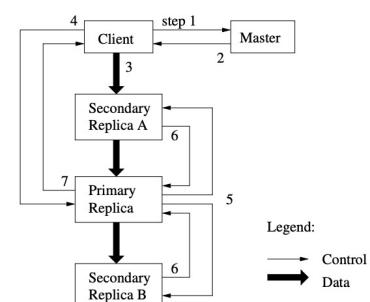
31

The Google File System

- Write Control and Data Flow

- Step 3:

- Client pushes the data to all the replicas
 - Order does not matter
 - Each chunkserver store the data in an LRU buffer in memory
 - Acknowledge to client after receiving the data



IOWA STATE UNIVERSITY

32

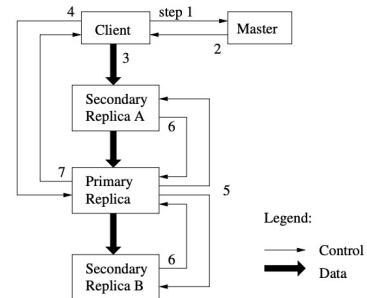
32

The Google File System

- Write Control and Data Flow

- Step 4:

- After all replicas have acknowledged receiving the data, the client sends a write request to the primary
 - The primary assigns consecutive serial numbers to all the mutations it receives, and applies the mutation to its own local state in order



IOWA STATE UNIVERSITY

33

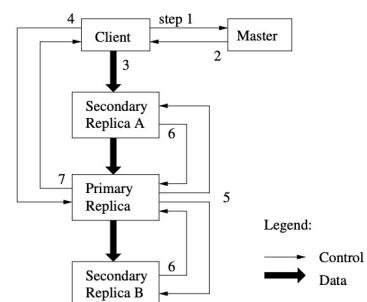
33

The Google File System

- Write Control and Data Flow

- Step 5:

- The primary forwards the write request to all other replicas
 - including the serial numbers
 - Each secondary replica applies mutations in the same serial number order assigned by the primary



IOWA STATE UNIVERSITY

34

34

The Google File System

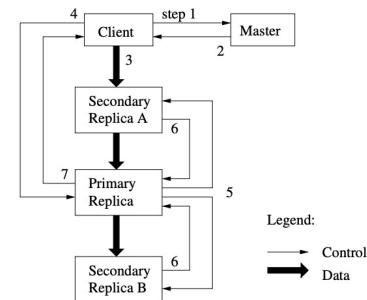
- Write Control and Data Flow

- Step 6:

- All secondary replicas reply to the primary indicating that they have completed the operation

- Step 7:

- The primary replies to the client
 - Report errors if occurred
 - Client may retry in case errors



IOWA STATE UNIVERSITY

35

35

The Google File System

- Other issues

- Atomic Appends
 - Fast snapshot
 - Master operation
 - namespace management & locking
 - replica placement & rebalancing
 - garbage collection (deleted/stale files)
 - detecting stale replicas
 - Master replication
 - master log & checkpoints replicated
 - shadow masters
 - provide read-access when primary is down

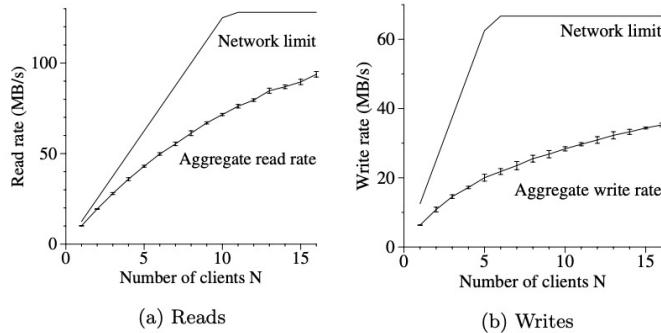
IOWA STATE UNIVERSITY

36

36

The Google File System

- Read & Write Performance
 - Both are scalable
 - i.e., read/write rates increase as the client count increases



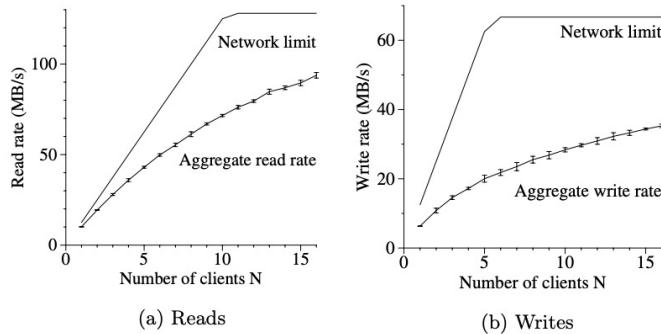
IOWA STATE UNIVERSITY

37

37

The Google File System

- Read & Write Performance
 - Why does the increasing rate drop?



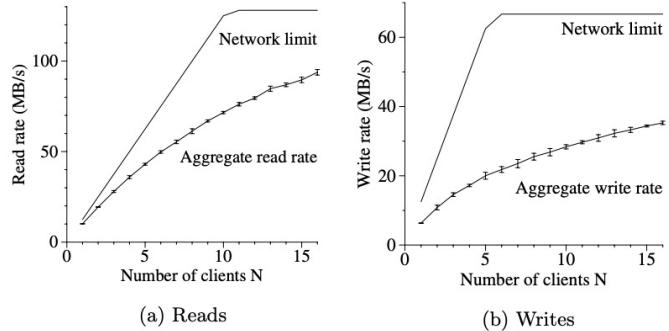
IOWA STATE UNIVERSITY

38

38

The Google File System

- Read & Write Performance
 - Why does the increasing rate of writes drop more?



IOWA STATE UNIVERSITY

39

39

Agenda

- **Recap**
- **Advanced Topics**
 - **The Google File System (GFS)**
 - **The World After GFS**

IOWA STATE UNIVERSITY

40

40

Hadoop Distributed File System (HDFS)

- An open source version of GFS

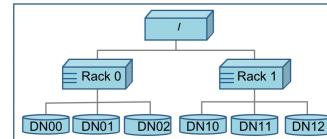
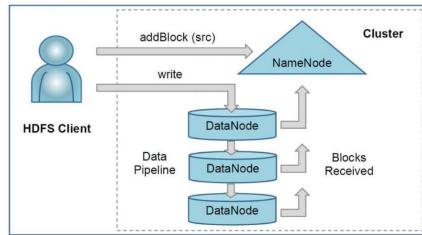
- formally published at 26th IEEE Symposium on Massive Storage Systems and Technologies (MSST'10)
 - NameNode + DataNodes
 - topology-aware replica

The Hadoop Distributed File System
Konstantin Shvachko, Hairong Kuang, Sanjay Radia, Robert Chansler
Neha Narkhede

The Hadoop Distributed File System

Konstantin Shvachko, Hairong Kuang, Sanjay Radia, Robert Chansler
Yahoo!
Sunnyvale, California USA
{Shv, Hairong, SRadia, Chansler}@Yahoo-Inc.com

Abstract—The Hadoop Distributed File System (HDFS) is designed to store very large data sets reliably and to stream them efficiently. It is built on commodity hardware in a large cluster, thousands of servers each host directly attached storage and provides a high level of fault tolerance. HDFS is a highly-computing system, many servers. Its resource can grow with demand while remaining economical at every size. We describe the design of HDFS and its implementation. We also discuss how to manage 25 petabytes of enterprise data at Yahoo!, using HDFS developed at Facebook [4], ZooKeeper [6] and Chukwa were originally developed at Yahoo! [8]. Apache Hadoop was originated at the University of California Berkeley [1] and is being used by many companies.



IOWA STATE UNIVERSITY

41

41

Hadoop Distributed File System (HDFS)

- Enable many follow-up research/innovations
 - across academia and industry

Analysis of HDFS Under HBase: A Facebook Messages Case Study

Tyler Harter, Dhruba Borthakur[†], Siying Dong[†], Amitanand Aiyer[†],
Liyin Tang[†], Andrea C. Arpacı-Dusseau, Remzi H. Arpacı-Dusseau
University of Wisconsin, Madison [†] Facebook Inc.

Abstract
We present a multilayer study of the Facebook Messages stack, which is based on HBase and HDFS. We collect and analyze HDFS traces to identify potential improvements, which we then evaluate via simulation. Measurements show that HDFS handles messages HDFS

ingly common, layered storage architecture: a distributed database (HBase, derived from BigTable [3]) atop a distributed file system (HDFS [24], derived from the Google File System [11]). Our goal is to study the interaction of these important systems, with a particular focus on the lower layer; thus, our highest-level question is: Is HDFS

Table 5. Summary statistics

Mingyuan Xia,* Mohit Saxena[†], Mario Blaum[†] and David A. Pease[†]
^{*McGill University, [†]IBM Research Almaden}

Abstract
Distributed storage systems are increasingly transitioning to the use of erasure codes since they offer higher reliability at significantly lower storage costs than data replication. However, these codes require every participant to have enough memory and CPU resources replicated across multiple machines and racks. For example, the Google File System [1] and the Hadoop Distributed File System [4] maintain three copies of each data block. Although disk storage seems inexpensive to day, replication of the entire data footprint is simply infeasible at massive scales of operation. As a result, most

HopsFS: Scaling Hierarchical File System Metadata Using NewSQL

Hopsx3: Scaling Hierarchies

Steffen Grohschmidt Mikael Ronström
Sapient AB Oracle

Sell Harald, Jim Dowling
KTH - Royal Institute of Technology
{fmunizc, mawie, harald, jdowling}@kth.se

Abstract

Recent improvements in both the performance and scalability of shared-nothing, transactional, in-memory NewSQL databases have reopened the research question of whether distributed metadata for hierarchical file systems can be managed using commodity databases. In this paper, we introduce HopFS, a next generation distribution of the Hadoop Distributed File System (HDFS) that

Spotify AB *Oracle*
 steffony@spotify.com mikael.norstrom@oracle.com

shards that change in size over time. Recent improvements in both the performance and scalability of shared-nothing, transactional, in-memory NewSQL [42] databases have reopened the possibility of storing distributed file system metadata in a commodity database. To date, the conventional wisdom has been that it is too expensive (in terms of throughput and latency) to store hierarchical file system metadata fully normalized.

FS: Hardening HDFS with Selective and Lightweight Verifications

Andrea C. Arpacı-Dusseau, Remzi H. Arpacı-Dusseau
University of Wisconsin, Madison ¹ University of Chicago

We hardened the *Hadoop Distributed File System (HDFS)* against *fail-silent* (no fail-stop) behaviors that result from memory corruption and software bugs using a new approach: selective and lightweight versioning (SLEEV). With this approach, actions performed by a single point of failure in today's systems. Growth from real systems show that these failures can lead to transient, non-deterministic errors, and make the system exhibit *fail-silent* behaviors (e.g., send corrupt messages) rather than crashing; these fail-silent errors can lead to data loss, unavailability, and prolonged debugging.

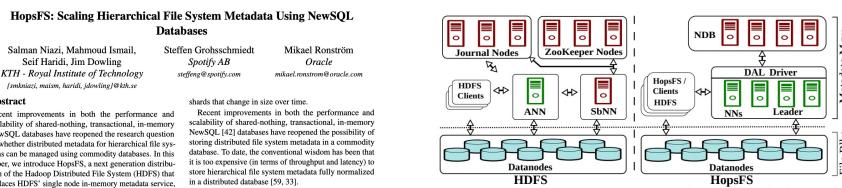
IOWA STATE UNIVERSITY

42

42

Hadoop Distributed File System (HDFS)

- Enable many follow-up research/innovations
 - E.g., HopsFS (FAST'17)
 - solve HDFS namenode's scalability issue
 - Limited metadata namespace
 - Limited concurrency
 - Leverage MySQL NDB Cluster to scale out metadata
 - 16X-37X the throughput of HDFS
 - Open source



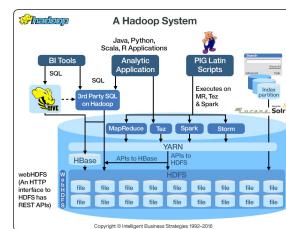
IOWA STATE UNIVERSITY

43

43

Hadoop Ecosystem

- Many software frameworks built atop HDFS
 - MapReduce, Spark, Spark Streaming, Hive, etc.
 - processing of live data streams, SQL queries on big data, etc.
 - developed by Apache Software Foundation
 - major vendors include Cloudera, Hortonworks, MapR Tech, Dell/EMC/Pivotal, IBM, Microsoft, AWS, etc.
 - market size expected to reach \$ 8240.1 million by 2026



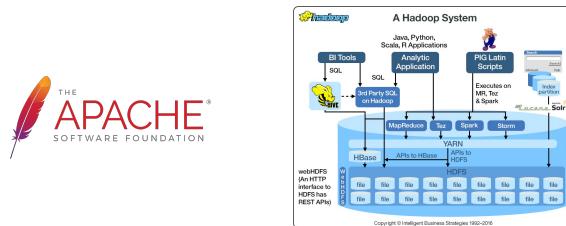
IOWA STATE UNIVERSITY

44

44

Hadoop Ecosystem

- Many software frameworks built atop HDFS
 - MapReduce, Spark, Spark Streaming, Hive, etc.
 - processing of live data streams, SQL queries on big data, etc.
 - developed by Apache Software Foundation
 - major vendors include Cloudera, Hortonworks, MapR Tech, Dell/EMC/Pivotal, IBM, Microsoft, AWS, etc.
 - market size expected to reach \$ 8240.1 million by 2026



IOWA STATE UNIVERSITY

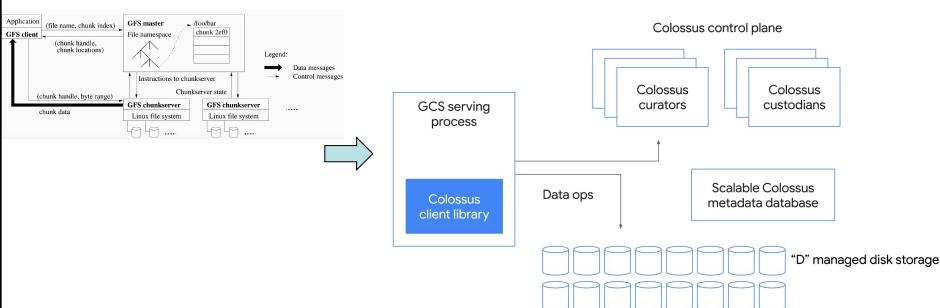
45

45

How about GFS itself?



- Colossus: successor of GFS
 - Google's cluster-level file system today
 - Enhanced storage scalability and availability
 - to handle the massive growth in data needs of an ever-growing number of applications



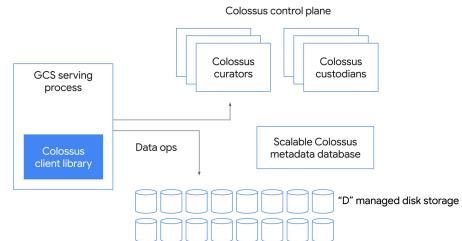
IOWA STATE UNIVERSITY

46

46

How about GFS itself?

- Colossus: successor of GFS
 - Client Library
 - Includes a lot of functionality (e.g., software RAID)
 - Tunable based on application needs
 - Control Plane
 - Curators: handle client requests (e.g., file creation)
 - Metadata database
 - Store metadata in Google's NoSQL database BigTable
 - Scale up by over 100x over GFS
 - Custodians
 - background maintenance (e.g., disk space balancing, RAID reconstruction)



IOWA STATE UNIVERSITY

47

47

Building Blocks of Google Storage

- Colossus: successor of GFS
 - Spanner: scalable relational database (OSDI'12)
 - Borg: cluster manager (EuroSys'15)



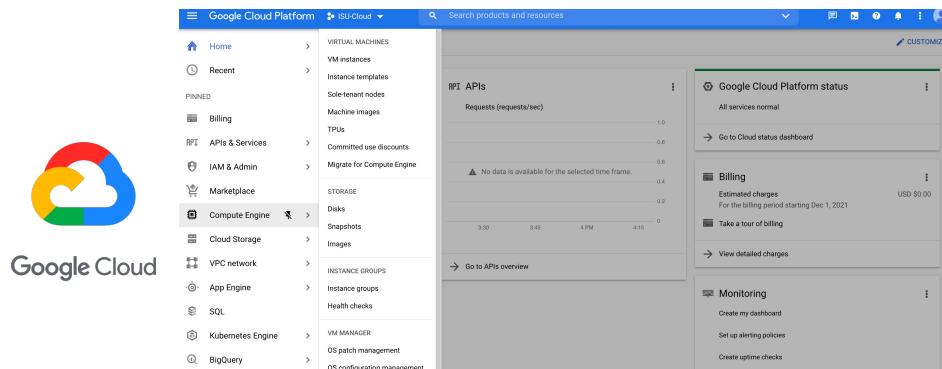
IOWA STATE UNIVERSITY

48

48

Fundamental to Google Cloud

- Various services enabled
 - E.g., VM, container, serverless functions, storage, database, ...
 - Conceptually similar to ChameleonCloud used in bonus assignment, but much more diverse
 - Free trial (\$300 credits)!



IOWA STATE UNIVERSITY

49

49

Agenda

Recap

Questions?

Advanced Topics

- The Google File System (GFS)



- The World After GFS

*acknowledgement: slides include content adapted from "Modern Operating Systems" by A. Tanenbaum, "Operating Systems Concepts" by A. Silberschatz etc., "Operating Systems: Three Easy Pieces" by R. Arpaci-Dusseau etc., lectures from Google, Stanford, UC Berkeley, and anonymous pictures from internet

IOWA STATE UNIVERSITY

50

50