

The role of Hadoop Distributed Files System(HDFS) in Technological era

Navjot Jyoti

Assistant Professor, Northwest Group of Institutions, Dhudike (Moga), Punjab, India

ABSTRACT

The popularity of cloud computing has been growing steadily for several years. It has been used more and more widely in many works of life. Cloud computing provides massive clusters for efficient large scale computation and data analysis. The MapReduce structure and its open source execution Hadoop have set up themselves as one of the most popular large informational collections analyzers. They are generally utilized by numerous cloud specialist organizations, for example, Amazon EC2 Cloud.

Keywords : MapReduce, HDFS, GFS

I. INTRODUCTION

Hadoop is an open source framework, facilitated by the Apache Software Foundation that gives a solid, fault tolerant, distributed document system and application programming interfaces. These empower its map-reduce framework to analyze substantial volumes of information in parallel. Straightforwardness of the Hadoop programming model takes into consideration clear usage of numerous applications. Java applications have the most direct access, yet Hadoop likewise has spilling capacities that take into account usage in any language.

Several organizations that need to handle large amounts of data are using map-reduce implementations to manage that data. Google started using a map-reduce system internally before 2004 [Dean 2004][1]. Yahoo runs the largest Hadoop cluster, running over a Linux cluster of over 10,000 cores [Yahoo 2008]. Vendors, such as Amazon, utilize Hadoop as part of their cloud computing service.

II. Hadoop Distributed File System

The Hadoop Distributed Files System (HDFS) keeps running over a local document framework and is just open through the Hadoop Application Interfaces Programming (APIs). **HDFS** appropriate arrangements information in similarly estimated pieces over the accessible information hubs. This division of information works best for huge records that can be put away as products of the lumping size designed for the HDFS.

In the event that the documents are littler than the lumping size, the HDFS will squander nearby record framework assets with exhaust, dispensed bytes. Excess and adaptation to internal failure are accomplished by recreating these pieces on different hubs. Hadoop endeavors to run the guide activities on duplicates of the information neighborhood the mapping undertaking. This decreases the measure of information that should be moved around. A HDFS may comprise of hundreds or thousands of server machines, each consisted of storing piece of the file information. Along these lines, recognition of flaws and speedy, programmed recuperation from them is a center design objective of HDFS. Applications that keep running on HDFS require streaming access to their data collections. They are not universally useful applications that regularly keep running on broadly useful document frameworks. HDFS is composed more for cluster preparing as opposed to intelligent use by clients [2].

III. Map-Reduce API

Hadoop uncovered three activities for executing the map reduce algorithm, mapping, joining, and decreasing. The framework is executed in Java; in any case, Hadoop additionally uncovered a gushing interface that permits programs written in any language to process every activity.

The need for real-time execution of MapReduce applications is increasingly arising and the effective processing strategy for deadline jobs has drawn attention of lots of researchers[3].

The process of MapReduce includes two major parts, Map function and Reduce function. The input files will be automatically split and copied to different computing nodes. After that the inputs will be sent to Map function in key-value pair format. The Map function will process the input pairs and also generate intermediate keyvalue pairs as inputs for Reduce function. The Reduce function will combine the inputs who have the same key and generate the final result. The final result will be written into the distributed file system[4]. The users only need to implement Map and Reduce functions. They do not need to concern about how to partition the input which is automatically done by the model.

The model will make tasks assigned evenly to machine. Sometimes, every to reduce communication optional overhead, а part Combine function can be employed. The Combine function can be regarded as a local Reduce, combines the results from Map function together locally, and returns a single value for Reduce function.

IV. HDFS Architecture

HDFS has a master/slave architecture. HDFS has an ace/slave engineering.[5]A HDFS cluster comprises of a solitary NameNode, a master server that deals with the file system namespace and manages access to files by users. Also there are various DataNodes, typically one for every cluster in the group, which manage capacity appended to the hubs that they keep running on. HDFS document framework uncovered а namespace and enables client information to be put away in files. Inside, a file is part into at least one squares and these pieces are put away in an arrangement of DataNodes. The NameNode executes file framework namespace operations like opening, closing, and renaming files and directories. It additionally decides the mapping of pieces to DataNodes. The DataNodes are in charge of serving read and write requests from file the system's clients.The DataNodes additionally perform square creation, deletion, replication upon and direction from the NameNode.



Fig 1 : HDFS Architecture

The NameNode and DataNode are bits of programming intended to keep running on item machines. These machines commonly run a GNU/Linux working framework (OS). HDFS is assembled utilizing the Java language; any machine that backings Java can run the NameNode or the DataNode programming. Use of the exceptionally versatile Java language implies that HDFS can be sent on an extensive variety of machines. A normal sending has a committed machine that just the NameNode runs programming. Each of alternate machines in the of the bunch runs one case DataNode programming. The architecture does not block running different DataNodes on a similar machine however in a genuine arrangement that is once in a while the case.

The presence of a single NameNode in a group incredibly streamlines the design of the framework. The NameNode is the judge and repository for all HDFS metadata. The framework is planned such that client information never flows through the NameNode.

Functions of NameNode:

• It is the master daemon that maintains and manages the DataNodes (slave nodes)

- It records the metadata of all the files stored in the cluster, e.g. The location of blocks stored,the size of the files, permissions, hierarchy, etc. There are two files associated with the metadata:
- **FsImage:** It contains the complete state of the file system namespace since the start of the NameNode.
- **EditLogs:** It contains all the recent modifications made to the file system with respect to the most recent FsImage.
- It records each change that takes place to the file system metadata. For example, if a file is deleted in HDFS, the NameNode will immediately record this in the EditLog.
- It regularly receives a Heartbeat and a block report from all the DataNodes in the cluster to ensure that the DataNodes are live.
- It keeps a record of all the blocks in HDFS and in which nodes these blocks are located.
- The NameNode is also responsible to take care of the replication factor of all the blocks which we will discuss in detail later in this HDFS tutorial blog.
- In case of the DataNode failure, the NameNode chooses new DataNodes for new replicas,balance disk usage and manages the communication traffic to the DataNodes.

Functions of DataNode:

- These are slave daemons or process which runs on each slave machine.
- The actual data is stored on DataNodes.
- The DataNodes perform the low-level read and write requests from the file system's clients.
- They send heartbeats to the NameNode periodically to report the overall health of

697

seconds.

The File System Namespace

HDFS supports a customary hierarchical file association. A user or an application can create directories and store files inside these directories. The file framework namespace chain of command is like most other existing document frameworks; one can make and evacuate documents, move a document starting with one directory then onto the next, or rename a file. HDFS does not yet actualize client standards. HDFS does not support hard connections or delicate connections. In any case, the HDFS architecture does not block executing these highlights.

The NameNode keeps up the record framework namespace[6]. Any change to the file framework namespace or its properties is recorded by the NameNode. An application can indicate the quantity of reproductions of a document that ought to be kept up by HDFS. The quantity of duplicates of a record is known as the replication factor of that document. This data is put away by the NameNode.

Data Replication

HDFS is intended to dependably store substantial records crosswise over machines in an expansive bunch. It stores each record as a grouping of blocks; all blocks in a file aside from the last block are a similar size. The pieces of a file are replicated for adaptation to internal failure. The block size and replication factor are configurable per document. An application can determine the quantity of copies of a document. The replication factor can be indicated at files creation time and can be changed later. Files in

HDFS, by default, this frequency is set to 3 HDFS are write once and have entirely one writer at a time.

> The NameNode settles on all choices with respect to replication of pieces. It intermittently gets a Heartbeat and a Blockreport from each of the DataNodes in the bunch. Receipt of a Heartbeat infers that the DataNode is working properly. A Blockreport contains a list of all blocks on a DataNode.



Fig 2 : Data replication

Balancer

HDFS block placement strategy does not consider DataNode plate space usage. This is to abstain from putting new information at a little subset of the DataNodes with a great deal of free storage. Along these lines data may not generally be put consistently crosswise over DataNodes. Imbalance likewise happens when new hubs are added to the clusters.

The balancer is a tool that adjusts hard disk space usage on a HDFS cluster[7]. The threshold value is taken between 0 and 1. A cluster is changed if, for each DataNode, the use of the node varies from the use of the entire cluster by near to the limit value.

The tool is conveyed as an application program that can be controlled by the cluster administration. It iteratively moves reproductions from DataNodes with higher usage to DataNodes with bring down use. One key necessity for the balancer is to keep up

information accessibility. While picking a reproduction to move and choosing its goal, the balancer ensures that the choice does not decrease either the quantity of replicas or the number of racks.

The balancer optimizes the balancing process by minimizing the inter-rack data copying. If the balancer decides that a replica A needs to be moved to a different rack and the destination rack happens to have a replica B of the same block, the data will be copied from replica B instead of replica A. A configuration parameter limits the bandwidth consumed by rebalancing operations. The higher the allowed bandwidth, the faster a cluster can reach the balanced state, but with greater competition with application processes.

The Communication Protocols

All HDFS communication protocols are layered on top of the TCP/IP protocol. A client establishes a connection to a configurable TCP port on the Namenode machine. It talks the ClientProtocol with the Namenode. The Datanodes talk to the Namenode using the DatanodeProtocol. A Remote Procedure Call (RPC) abstraction wraps both the ClientProtocol and the DatanodeProtocol. By design, the Namenode never initiates any RPCs. Instead, it only responds to RPC requests issued by Datanodes or clients.

Robustness

The primary objective of HDFS is to store data reliably even in the presence of failures. The three common types of failures are Namenode failures, Datanode failures and network partitions.

Data Disk Failure, Heartbeats and Re-Replication

Each Datanode sends a Heartbeat message to the Namenode periodically. A network partition can cause a subset of Datanodes to lose connectivity with the Namenode. The Namenode detects this condition by the absence of a Heartbeat message. The Namenode marks Datanodes without recent Heartbeats as dead and does not forward any new IO requests to them. Any data that was registered to a dead Datanode is not available to HDFS any more. Datanode death may cause the replication factor of some blocks to fall below their specified value. The Namenode constantly tracks which blocks need to be replicated and initiates replication may arise due to many reasons: a Datanode may become unavailable, a replica may become corrupted, a hard disk on a Datanode may fail, or the replication factor of a file may be increased.

Cluster Rebalancing

The HDFS architecture is compatible with data rebalancing schemes. A scheme might automatically move data from one Datanode to another if the free space on a Datanode falls below a certain threshold. In the event of a sudden high demand for a particular file, a scheme might dynamically create additional replicas and rebalance other data in the cluster. These types of data rebalancing schemes are not yet implemented.

Data Integrity

It is possible that a block of data fetched from a Datanode arrives corrupted. This corruption can occur because of faults in a storage device, network faults, or The HDFS client software. software buggy implements checksum checking on the contents of HDFS files. When a client creates an HDFS file, it computes a checksum of each block of the file and stores these checksums in a separate hidden file in the same HDFS namespace. When a client retrieves file contents it verifies that the data it received from each Datanode matches the checksum stored in the associated checksum file. If not, then the client can opt to retrieve that block from another Datanode that has a replica of that block.

Metadata Disk Failure

The FsImage and the EditLog are central data structures of HDFS. A corruption of these files can cause the HDFS instance to be non-functional. For this reason, the Namenode can be configured to support maintaining multiple copies of the FsImage and EditLog. Any update to either the FsImage or EditLog causes each of the FsImages and EditLogs to get updated synchronously. This synchronous updating of multiple copies of the FsImage and EditLog may degrade the rate of namespace transactions per second that a Namenode can support.

However, this degradation is acceptable because even though HDFS applications are very data intensive in nature, they are not metadata intensive. When a Namenode restarts, it selects the latest consistent FsImage and EditLog to use. The Namenode machine is a single point of failure for an HDFS cluster. If the Namenode machine fails, manual intervention is necessary. Currently, automatic restart and failover of the Namenode software to another machine is not supported.

COMPARATIVE ANALYSIS OF HDFS AND GFS

A GFS cluster consists of a single master and Multiple servers and is accessed by multiple clients[8] and HDFS has a master/slave architecture. HDFS cluster comprises of a solitary NameNode, a master server that deals with the file system namespace and manages access to files by users.

Properties	GFS		HDF	S
Design Goals	•	Goal is	•	One of
	to sup	port	the n	nain goals
	large files		of HDFS is to	
	•	Built	supp	ort large
	based	on the	files.	
	assum	ption	•	DBuilt
	that terabyte		based	l on the

	data sets will	assumption that		
	be distributed	terabyte data		
	across	sets will be		
	thousands of	distributed		
	disks attached	across thousands		
	to commodity	of disks attached		
	computer	to commodity		
	nodes.	compute nodes.		
	• Used	• Used		
	for data	for data		
	intensive	intensive		
	computing.	computing.		
	• Store	• Store		
	data reliably,	data reliably,		
	even when	even when		
	failures occur	failures occur		
	within chunk	within name		
	servers,	nodes, data		
	master, or	nodes, or		
	network	network		
	partitions.	partitions.		
	• GFS is	• HDFS is		
	designed more	designed more		
	for batch	for batch		
	processing	processing		
	rather than	rather than		
	interactive use	interactive use		
	by users.	by users.		
Processes	Master and	Name node and		
	chunk server	Data node		
File	• Files	• DHDFS		
Management	are organized	supports a		
	hierarchically	traditional		
	in directories	hierarchical file		
	and identified	organization		
	by path names.	• 🛛 HDFS		
	• GFS is	also supports		
	exclusively for	third-party file		
	Google only.	systems such as		
	1	CloudStore and		
		CloudStore and		
		Amazon Simple		

Scalability	• Cluster	• 🛛 🗠 🗠 🗠		Filenames	permissions for
	based	based		are random.	the user that is
	architecture	architecture		There are	the owner, for
	• The	• Hadoop		hundreds of	other users that
	file system	currently runs		thousands of	are members of
	consists of	on clusters with		files on a	the group, and
	hundreds or	thousands of		single disk,	for all other
	even	nodes.		and all the	users
	thousands of			data is	
	storage			obfuscated so	
	machines built			that	
	from			it is not	
	inexpensive			human	
	commodity			readable. The	
	parts.			algorithms	
	• The			uses for	
	largest cluster			obfuscation	
	have over 1000			changes all the	
	storage nodes,			time.	
	over 300 TB of		Security	• Google	• HDFS
	disk storage,			has dozens of	security is based
	and are			datacenters for	on the POSIX
	heavily			redundancy.	model of users
	accessed by			These	and groups.
	hundreds of			datacenters are	• Currentl
	clients on			in	y is security is
	distinct			undisclosed	limited to
	machines on a			locations and	simple
	continuous			most are	file
	basis.			unmarked for	permissions.
Protection	• Google	• The		protection.	The identity of a
	have their own	HDFS		• Access	client process is
	file system	implements a		is allowed to	just whatever
	called GFS.	permission		authorized	the host
	With GFS,	model for files		employees and	operating
	files are split	and directories		vendors only.	system says it is.
	up and	that shares		Some of the	Network
	stored in	much of the		protections in	authentication
	multiple	POSIX model.		place include:	protocols like
	pieces on	🛛 File or		24/7	Kerberos for
	multiple	directory has		guard	user
	machines.	separate		coverage,	authentication

	Electronic key	and		metadata.	cache It is a
	access, Access	encryption of		Neither the	facility provided
	logs, Closed	data transfers		sever nor the	by Mapreduce
	circuit	are yet not		client caches	framework to
	televisions,	supported.		the file data.	cache
	Alarms linked			• Chunk	application-
	to guard			s are stored as	specific, large,
	stations,			local files in a	read-only files
	Internal and			Linux system.	(text, archives,
	external			So, Linux	jars and so on)
	patrols, Dual			buffer cache	• Private
	utility power			already keeps	(belonging to
	feeds and			frequently	one user) and
	Backup power			accessed data	Public
	UPS and			in memory.	(belonging to all
	generators.			Therefore	the user of the
Database Files	Bigtable is the	HBase[15]		chunk servers	same node)
	database used	provides		need not cache	Distributed
	by GFS.	Bigtable		file data.	Cache Files.
	Bigtable is a	(Google)	Cache	• Appen	• HDFS's
	proprietary	[16]-like	Consistency	d-once-read-	write-once-
	distributed	capabilities on		many model is	read-many
	database of	top of Hadoop		adapted by	model that
	Google Inc.	Core.		Google. It	relaxes
File Serving	A file in GFS is	HDFS is divided		avoids the	concurrency
	comprised of	into large blocks		locking	control
	fixed sized	for		mechanism of	requirements,
	chunks. The	storage and		files for	simplifies data
	size of chunk	access, typically		writing in	coherency, and
	is 64MB.	64MB in		distributed	enables high
	Parts of a file	size. Portions of		environment is	throughput
	can be stored	the file can be		avoided.	access.
	on different	stored on		• Client	• Client
	nodes in a	different cluster		can append the	can only append
	cluster	nodes, balancing		data to the	to existing files
	satisfying the	storage		existing file.	
	concepts	resources and	Communicati	ТСР	RPC based
	load balancing	demand	on	connections	protocol on top
	and storage			are used for	of TCP/IP
	management.			communicatio	
Cache	• Clients	• HDFS		n. Pipelining is	
Management	do cache	uses distributed		used for data	

	transfer over	
	ТСР	
	connections	
Replication	• Chunk	• Automat
Strategy	replicas are	ic replication
	spread across	system. Rack
	the racks.	based system.
	Master	By default two
	automatically	copies of each
	replicates the	block are stored
	chunks.	by different
	• A user	Data Nodes in
	can specify the	the same rack
	number of	and a third copy
	replicas to be	is stored on a
	maintained.	Data Node
	• The	in a different
	master re-	rack (for greater
	replicates a	reliability).
	chunk replica	• An
	as soon as the	application can
	number of	specify the
	available	number of
	replicas falls	replicas of a file
	below a user-	that should be
	specified	maintained by
	number.	HDFS.
		Replication
		pipelining in
		case of write
		operations.
Available	GFS is a	Yahoo,
Implementati	proprietary	Facebook, IBM
on	distributed file	etc. are based on
	system	HDFS.
	developed by	
	Google for its	
	own	
	use.	

V. CONCLUSION

Hadoop Distributed File System and Mapreduce are the components of Hadoop project owned by Apache. Google File System is a proprietary distributed file system and is exclusive for Google Inc. Mapreduce is the programming frame work used by Google.

VI. REFERENCES

- Ke-Tha Yao,Robert F. Lucas et al. "Data anlysis for Massively distributed Simulations," Interservice/Industry Training, Simulation, and Education Conference (I/ITSEC) 2009
- [2]. Vidyasagar S. D"A Study on Role of Hadoop in Information Technology era" Volume : 2 | Issue :
 2 | Feb 2013 • ISSN No 2277 - 8160
- [3]. Jiazhen Han, Zhengheng Yuan et al. "An Adaptive Scheduling Algorithm for Heterogeneous Hadoop Systems" 2017 IEEE ICIS 2017, May 24-26, 2017, Wuhan, China
- [4]. Junbo Zhang, Dong Xiang, Tianrui Li, and Yi Pan "M2M: A Simple Matlab-to-MapReduce Translator for Cloud Computing" TSINGHUA SCIENCE AND TECHNOLOGY ISSN11007-0214 01/12 pp 1-9 Volume 18, Number 1, February 2013
- [5]. https://hadoop.apache.org/docs/r1.2.1/hdfs_desig n.html
- [6]. https://www.edureka.co/blog/apache-hadoophdfs-architecture/
- [7]. http://www.aosabook.org/en/hdfs.html
- [8]. R.Vijayakumari, R.Kirankumar, K.Gangadhara Rao "Comparative analysis of Google File System and Hadoop Distributed File System" International Journal of Advanced Trends in Computer Science and Engineering, Vol. 3, No.1, Pages :553-558 (2014)