



Introduction to Hadoop Distributed File System (HDFS)

Hadoop Distributed File System (HDFS) is a distributed, scalable, and portable filesystem written in Java for the Hadoop framework. It provides high-throughput access to application data and is designed to be fault-tolerant and suitable for use on commodity hardware.

LIST FILES

<code>hdfs dfs -ls /</code>	List all the files/directories for the given hdfs destination path
<code>hdfs dfs -ls -d /hadoop</code>	Directories are listed as plain files. In this case, this command will list the details of hadoop folder
<code>hdfs dfs -ls -h /data</code>	Format file sizes in a human-readable fashion (eg 64.0m instead of 67108864)
<code>hdfs dfs -ls -R /hadoop</code>	Recursively list all files in hadoop directory and all subdirectories in hadoop directory
<code>hdfs dfs -ls /hadoop/dat*</code>	List all the files matching the pattern. In this case, it will list all the files inside hadoop directory which starts with 'dat'

READ/WRITE FILES

<code>hdfs dfs -text /hadoop/derby.log</code>	Prints the content of the file in text format. Only one file is allowed. Text is output through TextOutputStream
<code>hdfs dfs -cat /hadoop/test</code>	This command will display the content of the HDFS file test on your stdout
<code>hdfs dfs -appendToFile /home/ubuntu/test1/hadoop/test2</code>	Appends the content of a local file test1 to a hdfs file test2

UPLOAD/DOWNLOAD FILES

<code>hdfs dfs -put /home/ubuntu/sample /hadoop</code>	Copies the file from local file system to HDFS
<code>hdfs dfs -put -f /home/ubuntu/sample /hadoop</code>	
<code>hdfs dfs -put -l /home/ubuntu/sample /hadoop</code>	
<code>hdfs dfs -put -p /home/ubuntu/sample /hadoop</code>	
<code>hdfs dfs -get /newfile /home/ubuntu/</code>	
<code>hdfs dfs -get -p /newfile /home/ubuntu/</code>	modification times, ownership and the mode
<code>hdfs dfs -get /hadoop/*.txt /home/ubuntu/</code>	Copies all the files matching the pattern from local file system to HDFS
<code>hdfs dfs -copyFromLocal /home/ubuntu/sample /hadoop</code>	Works similarly to the put command, except that the source is restricted to a local file reference
<code>hdfs dfs -copyToLocal /newfile /home/ubuntu/</code>	Works similarly to the put command, except that the destination is restricted to a local file reference
<code>hdfs dfs -moveFromLocal /home/ubuntu/sample /hadoop</code>	Works similarly to the put command, except that the source is deleted after it's copied

FILE MANAGEMENT

<code>hdfs dfs -mv /hadoop/file1 /hadoop1</code>	Copies file from source to destination on HDFS. In this case, copying file1 from
--	--

<code>hdfs dfs -rm -r /hadoop</code>	Deletes the directory and any content under it recursively
<code>hdfs dfs -rm -R /hadoop</code>	
<code>hdfs dfs -rmr /hadoop</code>	
<code>hdfs dfs -rm -skipTrash /hadoop</code>	The -skipTrash option will bypass trash, if enabled, and delete the specified file(s) immediately
<code>hdfs dfs -rm -f /hadoop</code>	If the file does not exist, do not display a diagnostic message or modify the exit status to reflect an error
<code>hdfs dfs -rmdir /hadoop1</code>	Delete a directory
<code>hdfs dfs -mkdir /hadoop2</code>	Create a directory in specified HDFS location
<code>hdfs dfs -mkdir -f /hadoop2</code>	Create a directory in specified HDFS location. This command does not fail even if the directory already exists
<code>hdfs dfs -touchz /hadoop3</code>	Creates a file of zero length at <path> with current time as the timestamp of that <path>

<code>hdfs dfs -find /hadoop/files</code>	Prints the names of files or files that match the file pattern <src> to stdout
<code>hdfs dfs -chmod 755 /hadoop/file1</code>	Changes permissions of the file
<code>hdfs dfs -chmod -R 755 /hadoop</code>	Changes permissions of the files recursively
<code>hdfs dfs -chown ubuntu:ubuntu /hadoop</code>	Changes owner of the file. 1st ubuntu in the command is owner and 2nd one is group
<code>hdfs dfs -chown -R ubuntu:ubuntu /hadoop</code>	Changes owner of the files recursively
<code>hdfs dfs -chgrp ubuntu /hadoop</code>	Changes group association of the file
<code>hdfs dfs -chgrp -R ubuntu /hadoop</code>	Changes group association of the files recursively

Storage Operations

Learn basic HDFS commands for storage and data management.

<code>hdfs dfs -du -h /hadoop/file</code>	Show the amount of space, in bytes, used by the files that match the specified file pattern. Formats the sizes of files in a human-readable fashion
---	---

ADMINISTRATION

<code>hdfs balancer -threshold 30</code>	Runs a cluster balancing utility. Percentage of disk capacity. This overwrites the default threshold
<code>hadoop version</code>	To check the version of Hadoop
<code>hdfs fsck /</code>	It checks the health of the Hadoop file system
<code>hdfs dfsadmin -safemode leave</code>	The command to turn off the safemode of NameNode
<code>hdfs dfsadmin -refreshNodes</code>	Re-read the hosts and exclude files to update the list of nodes that are allowed to connect to the Namenode and decommissioned or recommissioned
<code>hdfs namenode -format</code>	Formats the NameNode

Overview of HDFS and Storage Commands

What is HDFS?

HDFS is the primary storage system used by Hadoop applications.

Creating HDFS Directory

- **Objective:** Create a directory structure in HDFS.
- **Command I:** `hdfs dfs -mkdir -p /user/your_username/your_directory`
- **Command II:** `hdfs dfs -put /path/to/local/file /user/your_username/your_directory`

```
2 2022-03-22 23:58
66736 2022-03-22 23:57
47128 2022-03-22 23:57
88969 2022-03-22 23:58
2539 2022-03-23 21:10
86359 2022-03-22 23:58
```

```
txt /user/cloudera/demo
```

```
66736 2022-03-22 23:57
47128 2022-03-22 23:57
88969 2022-03-22 23:58
86359 2022-03-22 23:58
```

```
2 2022-03-22 23:58
2539 2022-03-23 21:10
```

Moving Files to HDFS

1

Demonstrate

Explain the process of moving files from the local machine to HDFS.

2

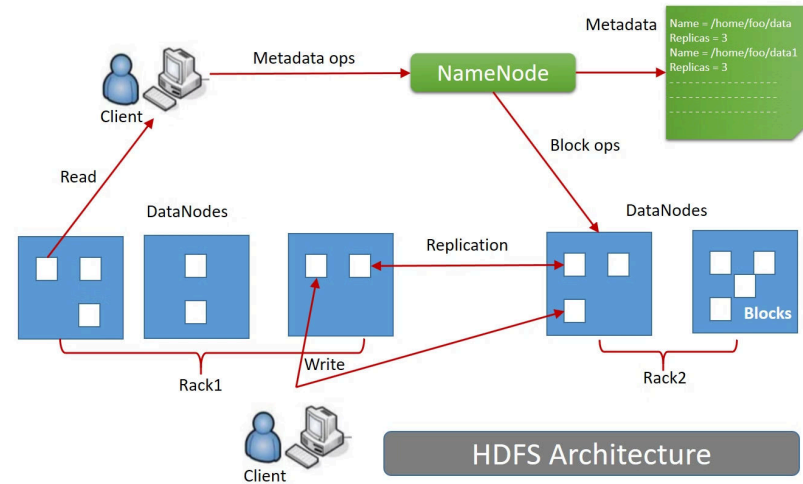
Command Usage

Illustrate the use of the command: `hdfs dfs -put /path/to/local/file /user/your_username/your_directory`

Exploring HDFS Data

Demonstrate how to explore and view data in Hadoop Distributed File System (HDFS) using the following commands:

- `hdfs dfs -cat`
`/user/your_username/your_directory/your_file`
- `hdfs dfs -ls`
`/user/your_username/your_directory`



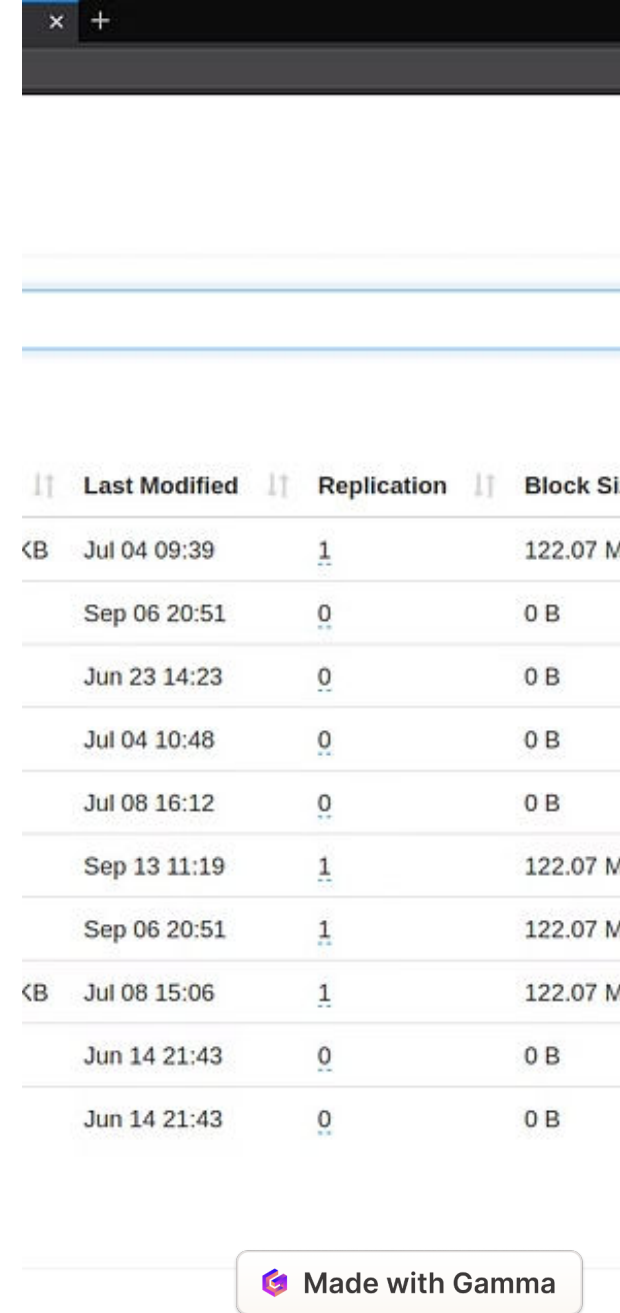
Retrieving Data from HDFS

Objective

Copy a file from HDFS to the local file system.

Command I

Use `hdfs dfs -get /user/your_username/your_directory/your_file /path/to/local/directory` to copy the file to the local system.



The screenshot shows a web browser window with a table of HDFS file metadata. The table has columns for file size, last modified date, replication factor, and block size. The data is organized into two sections, each starting with a file size indicator (e.g., <B).

	Last Modified	Replication	Block Size
<B	Jul 04 09:39	1	122.07 M
	Sep 06 20:51	0	0 B
	Jun 23 14:23	0	0 B
	Jul 04 10:48	0	0 B
	Jul 08 16:12	0	0 B
	Sep 13 11:19	1	122.07 M
	Sep 06 20:51	1	122.07 M
<B	Jul 08 15:06	1	122.07 M
	Jun 14 21:43	0	0 B
	Jun 14 21:43	0	0 B

Made with Gamma

Lab Objectives

Recap Objectives

Using HDFS commands to move data for processing.

Importance of Data Preparation

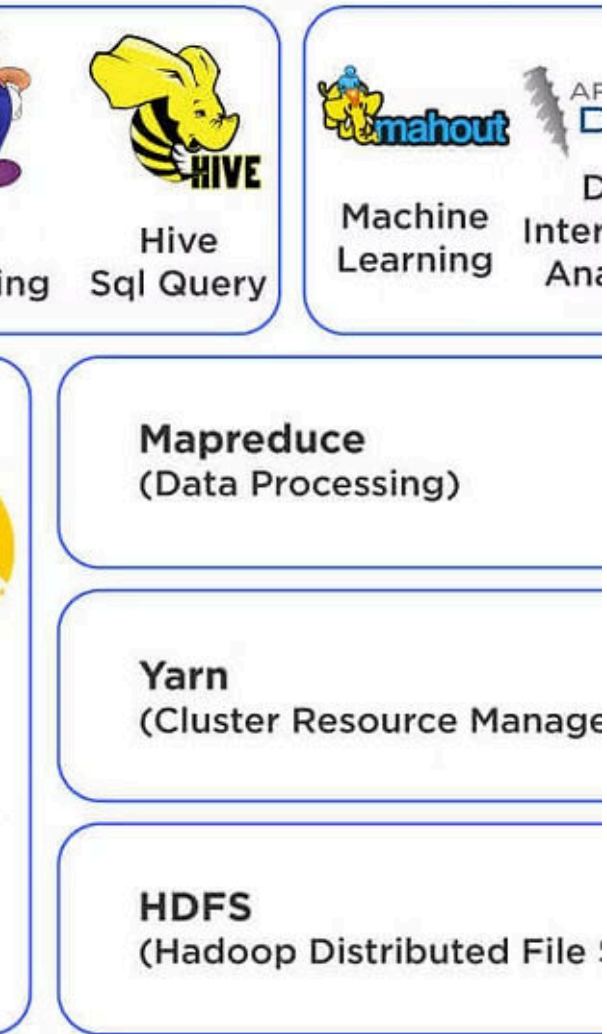
Emphasize the importance of preparing data in HDFS for subsequent processing.

Executing HDFS Commands

Reinforce the importance of utilizing HDFS commands effectively for managing data.

Ensure that all files stored in HDFS can be easily accessed using the appropriate commands.

Hadoop Ecosystem



Source: Data Flair

Summary and Conclusion

1

Recap Key Points

Review the main concepts of HDFS.

2

Importance of Mastery

Highlight the significance of mastering HDFS commands.

3

Efficient Data Management

Emphasize the role of Hadoop commands in efficient data management.



Summary and Conclusion

Recap the key points covered in the presentation, emphasizing the importance of mastering HDFS commands for efficient data management in Hadoop. Understanding HDFS and its storage commands is fundamental for anyone working with big data. Efficient management and manipulation of data in HDFS are essential skills for data engineers and analysts.