

Introduction to Hadoop

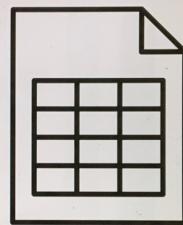
Objectives

After watching this video, you will be able to:

- Define Hadoop
- Explain the history of Hadoop
- List reasons why Hadoop was the answer to Big Data processing

What is Hadoop?

Hadoop is an open-source framework used to process enormous data sets



Why Hadoop?



Why Hadoop?

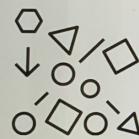


Why Hadoop?

Hadoop was designed to help organizations manage terabytes of data



Why Hadoop?



Unstructured Data



Relational Data

What is Hadoop?



used as the framework for Big Data operations

- Set of open-source programs and procedures
- Used for processing large amounts of data
- Servers run applications on clusters
- Handles parallel jobs or processes

A Hadoop Cluster is a collection of computers working together to perform tasks

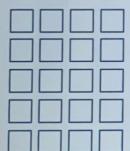
Hadoop is not a database but an ecosystem that can

A little bit of history

- 1999 • Apache software foundation established
- 2002 • Nutch web search engine created
- 2006 • Nutch was divided into Web crawler, distributed systems
- 2008 • Yahoo released Hadoop as an open-source project

Types of Big Data

Structured Data



Semi-Structured Data



Unstructured Data



How does Hadoop work?

Hadoop is now in petabytes and exabytes

core components



Hadoop Common is an essential part of the Apache Hadoop Framework that refers to the collection of common utilities and libraries that support other Hadoop modules

How does Hadoop work?

storage component

HDFS

- Hadoop stands for Distributed File System
- Handles and stores large data
- Scales a single Hadoop cluster into as much as thousand clusters

MapReduce

- Known as Hadoop's processing unit
- Processes Big Data by splitting the data into smaller units and processes them simultaneously
- The first method used to query data stored in HDFS

processes the RAM and CPU

YARN

- Yet Another Resource Negotiator acronym
- Prepares Hadoop for batch, stream, interactive and graph processing

A commodity hardware is low-specifications industry-grade hardware

The challenges of Hadoop

not good for Hadoop

- Processing transactions (random access)
- When work cannot be parallelized
- When there are dependencies in the data
- Low latency data access
- Processing lots of small files
- Intensive calculations with little data

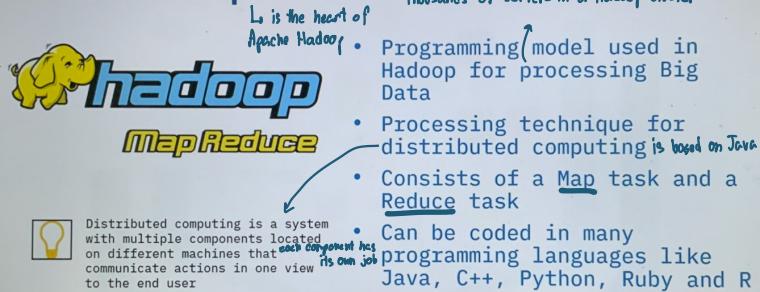
Intro to MapReduce

Objectives

After watching this video, you will be able to:

- Explain the terms Map and Reduce in MapReduce
- Describe why we use MapReduce
- List the components of MapReduce
- Outline examples of common use cases of MapReduce

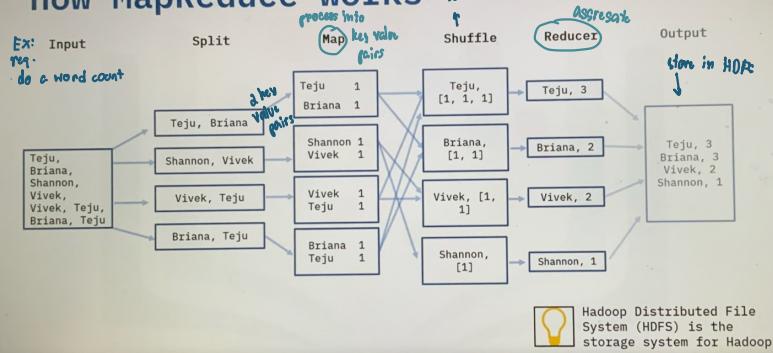
What is MapReduce?



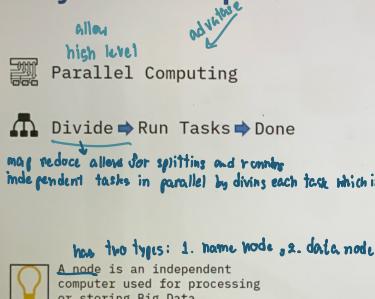
Map and reduce

- takes in*
- Input file
 - Map: Processes data into key value pairs
 - Further data sorting and organizing before the preliminary output is sent to
 - Reducer: Aggregates and computes a set of result and produces a final output
 - MapReduce keeps track of its task by creating a unique key
- these are preliminary output list*

How MapReduce works



Why use MapReduce?



Flexible

- Process data in tabular and
- Non tabular forms, such as videos
- Support for multiple languages
- Platform for analysis and data warehousing

Common use cases

Social media

Limited in, long to grow

Social media platforms can use MapReduce to analyze who visited your profile and who viewed your posts

Netflix

Recommendations

Create a recommender's system for users and provide suggestions for them based on their interest

banks

Financial industries

Can be used for fraud detection by analyzing behaviors of buyers and tracking down anomalies

Google

Advertisement

Can be used to analyze and understand the interaction with ads and the engagement levels

Hadoop Ecosystem

Objectives

After watching this video, you will be able to:

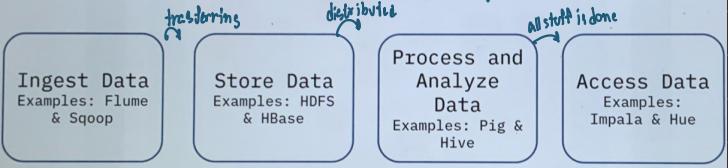
- List the stages of the Hadoop Ecosystem
- Differentiate between the core components and the extended components
- List some examples of tools used in each stage



Core components of Hadoop:
Hadoop common, HDFS, MapReduce, and YARN
stores the data collected from the ingestion and distributes the data across multiple nodes
used for making Big Data manageable by processing them in Cluster
the resource manager across cluster

Extended Hadoop ecosystem

consists of libraries software packages.
The Hadoop ecosystem is made up of components that support one another for big data processing.



Ingest data

Ingest Data

- Flume distributed service
 - Collects, aggregates, and transfers big data
 - Has a simple and flexible architecture based on streaming data flows
 - Uses a simple extensible data model that allows for online analytic application
- Sqoop
 - Designed to transfer data between relational database systems and Hadoop
 - Accesses the database to understand the schema of the data
 - Generates a MapReduce application to import or export the data

Analyze data

Analyze Data

- Pig
 - Analyzes large amounts of data
 - Operates on the client side of a cluster
 - A procedural data flow language
- Hive
 - Used for creating reports
 - Operates on the server side of a cluster
 - A declarative programming language

A declarative programming language allows users to express which data they wish to receive

Store data

Store Data

HBase

- A non-relational database that runs on top of HDFS
- Provides real time wrangling on data
- Stores data as indexes to allow for random and faster access to data

Cassandra

- A scalable, NoSQL database designed to have no single point of failure

Access data

Access Data

Impala

- Scalable and easy to use platform for everyone
- No programming skills required

Hue

- Stands for Hadoop user experience
- Allows you to upload, browse, and query data
- Runs Pig jobs and workflow
- Provides editors for several SQL query languages like Hive and MySQL

HDFS

HDFS

- allows programmers to access or store files from any network or computer
- HDFS is the acronym for Hadoop Distributed File System
- It is the storage layer of Hadoop
- Splits the files into blocks, creates replicas of the blocks, and stores them on different machines
- Provides access to streaming data
- HDFS uses a command line interface to interact with Hadoop



Streaming means that HDFS provides a constant bitrate when transferring data rather than having the data being transferred in waves

Key features

Cost efficient

The storage hardware is not expensive

Large amounts of data

it splits these large amount of data into small chunks called blocks
HDFS can store up to petabytes of data in any format

Replication

Makes copies of the data on multiple machines

Fault tolerant

If one machine crashes, a copy of the data can be found somewhere else and work continues

Scalable

One cluster can be scaled into hundreds of nodes

Portable

Can easily move across multiple platforms

HDFS concepts

Blocks

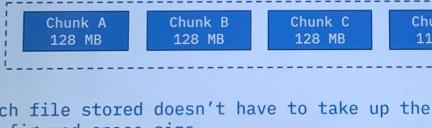
- Minimum amount of data that can be read or written
- Provides fault tolerance
- Default size is 64 MB or 128 MB

HDFS concepts

Blocks



File size: 500 MB



- Each file stored doesn't have to take up the configured space size

HDFS concepts

Nodes

This node regulates file access to the clients and maintains, manages, and assigns tasks to the secondary node **AAA data node**

Nodes

A node is a single system which is responsible to store and process data

can be hundred nodes

Primary Node

These nodes are the actual workers in the HDFS system and take instructions from the primary node

Secondary Node

HDFS concepts

Replication

- Creating a copy of the data block
- Copies are created for backup purposes
- Replication factor: Number of times the data block was copied

HDFS concepts

Rack awareness in HDFS

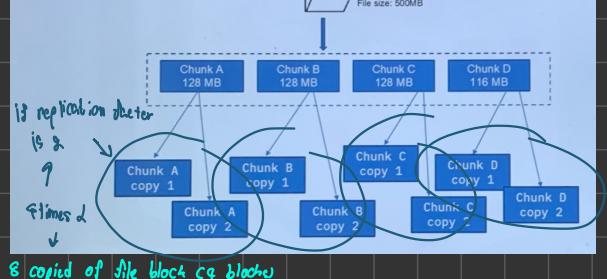
- Choosing data node racks that are closest to each other
- Improves cluster performance by reducing network traffic
- Name node keeps the rack ID information
- Replication can be done through rack awareness

*create replicas of data
node are in different rack*
A rack is the collection of about 40 to 50 data nodes using the same network switch

HDFS concepts

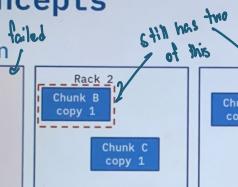
Replication

File size: 500MB



HDFS concepts

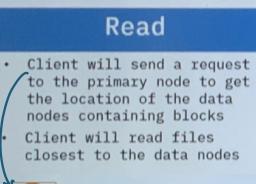
Replication



HDFS concepts

Read and Write Operations

- HDFS allows write once read many operations

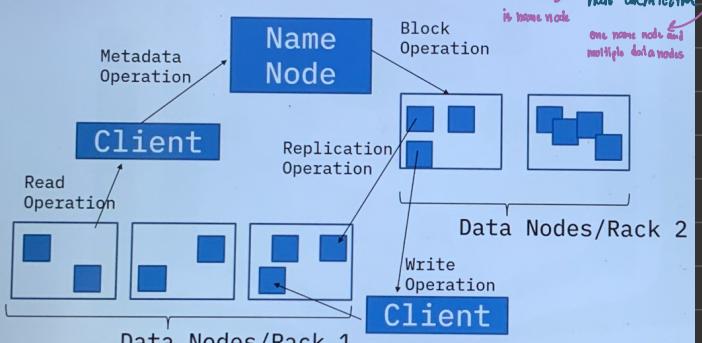


A client fulfills a user's request by interacting with the Name node and Data nodes



client receives a write message from the data node when client is done, the data nodes start creating replicas and sends confirmation

HDFS architecture



HIVE

Objectives

After watching this video, you will be able to:

- Articulate the reasons why Hive is used
- List Hive features
- Differentiate between Hive vs. traditional RDBMS
- Identify components of the Hive architecture
- Discuss Hive concepts



RDBMS - Relational Database Management System

A relational database is a database that stores data in a structured format tabular form

Hive

- store historical data from many sources
to see its insight*
- Hive is data warehouse software within Hadoop that is designed for reading, writing and managing large tabular-type datasets and data analysis:
 - It is scalable, fast, and easy to use
 - Hive Query Language (HiveQL) is inspired by SQL, making it easier for users to grasp concepts
 - It supports data cleansing and filtering depending on users' requirements



A data warehouse stores historical data from many different sources so that you can analyze and extract insights from it

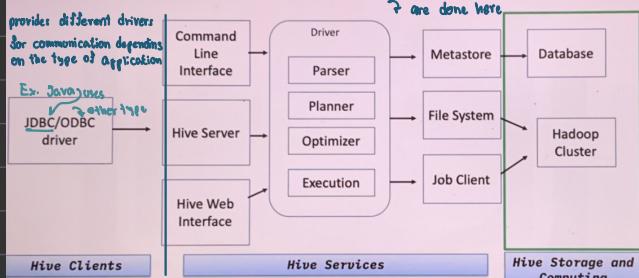
Hive and traditional RDBMS compared

Traditional RDBMS	Hive
Used to <u>maintain a Database</u> and uses SQL	Used to <u>maintain a data warehouse</u> using Hive query language
Suited for real-time/dynamic data analysis like data from <u>sensors</u>	Suited for static data analysis like a text file containing names
Designed to read and write <u>as many times</u> as it needs	Designed on the methodology of <u>write once, read many</u>
Maximum data size it can handle is <u>terabytes</u>	Maximum data size it can handle is <u>petabytes</u>
Enforced that the schema must <u>verify</u> loading data before it can proceed	Doesn't enforce the schema to <u>verify</u> loading data
May <u>not always</u> have built-in for <u>support</u> data partitioning	Supports partitioning



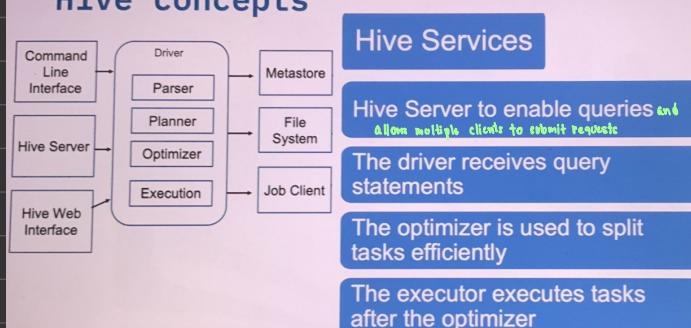
Partitioning means dividing the table into parts based on the values of a particular column, such as date or city

Hive architecture



Hive concepts

Hive concepts



Metastore stores the metadata information about the tables

Hive Clients

JDBC clients allows Java applications to connect to Hive

ODBC client allows applications based on ODBC to connect to Hive

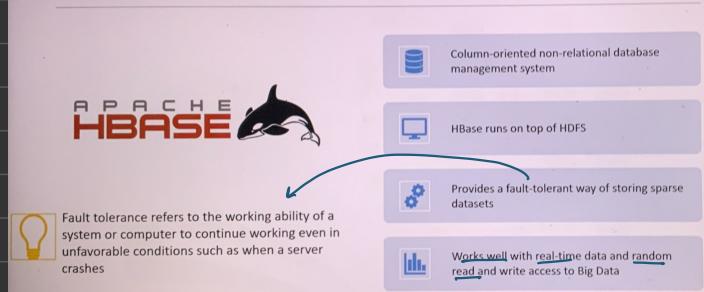
HBASE

Objectives

After watching this video, you will be able to:

- Define HBase
- Describe HBase as a columnar database
- List the HBase features and usage
- Outline the differences between HBase and HDFS
- Describe the HBase architecture
- Discuss HBase concepts

HBase



HBase features

HBase is used for write-heavy applications

HBase is linearly and modularly scalable

It is a backup support for MapReduce jobs

It provides consistent reads and writes

It has no fixed column schema

It is an easy-to-use Java API for client access

It provides data replication across clusters

An example

Patient Details		Heart Rate	Time Stamp
Patient Name	Patient Age	Heart Rate	Time Stamp
Patient A	28	120 BPM	8:50 AM
Patient A	28	110 BPM	10:10 AM
Patient B	77	95 BPM	11:00 AM
Patient C	45	150 BPM	11:30 AM
:			
Patient B	77	115 BPM	12:30 AM

each row is a log record

A few things to note:

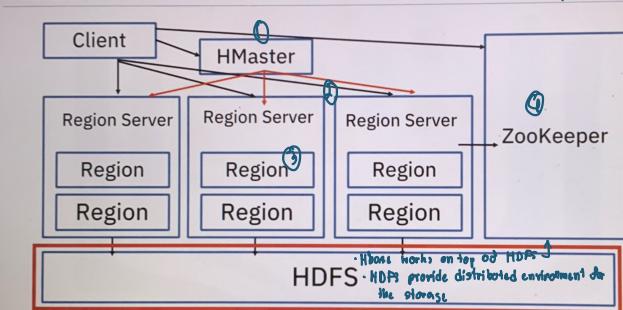
- Predefine the table schema and specify column families
- New columns can be added to column families at anytime
- HBase schema is very flexible
- HBase has master nodes to manage the cluster and region servers to perform the work

Differences between HBase and HDFS

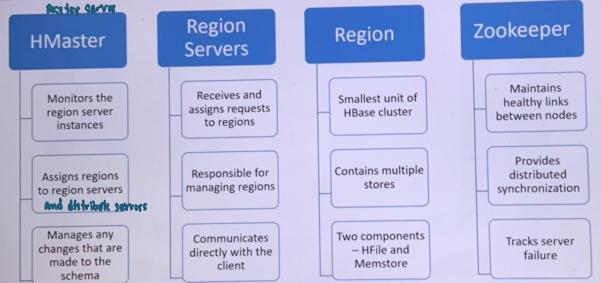
both used to store massive amounts of data

HBase	HDFS
HBase stores data in the form of <u>columns</u> and <u>rows</u> in a table	HDFS stores data in a <u>distributed</u> manner across different nodes on that network
HBase allows <u>dynamic</u> changes	HDFS has a <u>rigid architecture</u> that doesn't allow changes
HBase is suitable for <u>random writes</u> and <u>reads</u> of data stored in HDFS	HDFS is suited for <u>write once</u> and <u>read many times</u>
HBase allows for storing and processing of Big Data	HDFS is for <u>storing only</u>

HBase architecture



HBase concepts



CHEAT SHEET

Module 2 Cheat Sheet: Introduction to the Hadoop Ecosystem

Package/Method	Description	Code Example
bin/hadoop	All Hadoop commands are invoked by the bin/hadoop script. Running the Hadoop script without any arguments prints the description for all commands.	<p>Running Hadoop script without arguments:</p> <pre>1 bin/hadoop</pre>
cat	Reads each file parameter in sequence and writes it to standard output. If you do not specify a file name, the cat command reads from standard input. You can also specify a file name of - (dash) for standard input.	<p>Create two sample files.</p> <pre>1 echo "This is file 1" > file1.txt 2 echo "This is file 2" > file2.txt</pre> <p>Use the cat command to read and display the contents of both files</p> <pre>1 cat file1.txt file2.txt</pre> <p>Sample output (Contents of file1.txt and file2.txt):</p> <pre>1 This is file 1 2 This is file 2</pre>
cd	Used to move efficiently from the existing working directory to different directories on your system.	<p>Basic syntax of cd command:</p> <pre>1 cd [options]... [directory]</pre> <p>Example 1: Change directory location to "folder1"</p> <pre>1 cd /usr/local/folder1</pre> <p>Example 2: Get back to the previous working directory</p> <pre>1 cd -</pre> <p>Example 3: Move up one level from the present working directory tree</p> <pre>1 cd ..</pre>
create table	Used to create a new table in a database	<p>Create a new database (if not already created).</p> <pre>1 CREATE DATABASE your_database;</pre> <p>Use the newly created database.</p> <pre>1 USE your_database;</pre> <p>Create a new table named "employees" in Hive.</p> <pre>1 CREATE TABLE employees (2 id INT, 3 first_name STRING, 4 last_name STRING, 5 email STRING, 6 hire_date DATE 7) 8 ROW FORMAT DELIMITED 9 FIELDS TERMINATED BY ',' 10 STORED AS TEXTFILE;</pre> <p>Show the list of tables in the database.</p> <pre>1 SHOW TABLES;</pre> <p>Sample Output (List of Tables):</p> <pre>1 OK 2 employees</pre>
curl	A command-line tool (pronounced "curl") that allows data to be exchanged between a device and a server through a terminal. The user specifies the server URL, the location where they want to send the request, and the data they want to send to the server URL using this command-line interface (CLI).	<p>Example 1: Sending a GET request and displaying the response</p> <p>Send a GET request to a server and display the response.</p> <pre>1 curl https://www.example.com</pre> <p>In this example, we use the curl command to send a GET request to https://www.example.com and display the HTML response from the server.</p> <p>-----</p> <p>Example 2: Sending data to a server using POST Request</p> <p>Send a POST request with data to a server and display the response.</p> <pre>1 curl -X POST -d "name=John&age=30" https://www.example.com/api</pre> <p>In this example, we use the curl command to send a POST request to https://www.example.com/api with data name=John&age=30 and display the JSON response from the server.</p>

docker exec	<p>Runs a new command in a running container. It only runs when the container's primary process is running, and it is not restarted if the container is restarted.</p>	<p>Running a command in a running Docker container: Run a new command inside a running Docker container.</p> <pre>1 docker exec -it container_name_or_id ls /app</pre> <p>Sample Output (List of files in the '/app' Directory inside the container):</p> <pre>1 file1.txt 2 file2.txt 3 subdirectory</pre> <p>In this example:</p> <ul style="list-style-type: none"> • <code>docker exec</code> is used to run a new command (<code>ls /app</code>) inside a running Docker container. • <code>-it</code> enables an interactive terminal session, which allows you to see the output of the command. • <code>container_name_or_id</code> is the name or ID of the running Docker container you want to execute the command in. • <code>ls /app</code> is the command that lists the files and directories in the '/app' directory inside the container.
docker-compose	<p>Compose is a tool for defining and running multi-container Docker applications. It uses the YAML file to configure the services and enables us to create and start all the services from just one configuration file.</p>	<p>Starting Docker containers using docker-compose: Suppose you have a <code>docker-compose.yml</code> file like this:</p> <pre>1 version: '3' 2 services: 3 web: 4 image: nginx:latest 5 ports: 6 - "80:80" 7 db: 8 image: postgres:latest 9 environment: 10 POSTGRES_PASSWORD: example_password</pre> <p>You can use <code>docker-compose</code> to start the services defined in the <code>docker-compose.yml</code> file as follows:</p> <p>Navigate to the directory containing the <code>docker-compose.yml</code> file.</p> <pre>1 cd /path/to/your/docker-compose-project</pre> <p>Start the Docker containers defined in the <code>docker-compose.yml</code> file</p> <pre>1 docker-compose up</pre>
docker pull	<p>You can download Docker images from the internet.</p>	<pre>1 docker pull [OPTIONS] IMAGE_NAME[:TAG]</pre>
docker run	<p>It runs a command in a new container, getting the image and starting the container if needed.</p>	<pre>1 docker run [OPTIONS] IMAGE [COMMAND] [ARG...]</pre>
git clone	<p>You can create a copy of a specific repository or branch within a repository.</p>	<pre>1 git clone REPOSITORY_URL [DESTINATION_DIRECTORY]</pre>
hdfs dfs	<p>Apache Hadoop hadoop fs or hdfs dfs are file system commands to interact with HDFS. These commands are very similar to Unix commands. Hadoop provides two types of commands to interact with the file system: hadoop fs or hdfs dfs. The major difference is that Hadoop commands are supported with multiple file systems like S3, Azure, and many more.</p>	<p>Example-1: Listing files and directories in HDFS: List files and directories in the root directory of HDFS.</p> <pre>1 hdfs dfs -ls /</pre> <p>Example-2: In this example, we use the <code>hdfs dfs -ls</code> command to list files and directories in the root directory of HDFS.</p> <pre>1 hdfs dfs -ls /</pre> <p>Sample output:</p> <pre>1 drwxr-xr-x - hdfs hduser 0 2023-09-13 10:00 /user 2 drwxr-xr-x - hdfs hduser 0 2023-09-13 10:05 /tmp 3 drwxr-xr-x - mapred hduser 0 2023-09-13 10:10 /mapred</pre> <p>Create a new directory named "mydata" in HDFS.</p> <pre>1 hdfs dfs -mkdir /user/your_username/mydata</pre>
hdfs dfs -cat	<p>Display the contents for a file.</p>	<p>Display the contents of a file in HDFS.</p> <pre>1 hdfs dfs -cat /path/to/file.txt</pre>
hdfs dfs -mkdir	<p>Creates a directory named path in HDFS</p>	<p>Create a directory in HDFS.</p> <pre>1 hdfs dfs -mkdir /user/username/mydirectory</pre>
hdfs dfs -put	<p>Upload a file or folder from the local disk to HDFS.</p>	<p>Upload a file from the local file system to HDFS.</p> <pre>1 hdfs dfs -put localfile.txt /user/username/hdfsfile.txt</pre>
LOAD DATA INPATH	<p>Hive provides the functionality to load precreated table entities either from the local file system or from HDFS. This command is used to load data into the hive table.</p>	<p>Load data from HDFS into a Hive table.</p> <pre>1 LOAD DATA INPATH '/user/username/hdfsfile.txt' INTO TABLE 2 mytable;</pre>

		Basic command syntax
ls	Writes to standard output the contents of each specified Directory parameter or the name of each specified file parameter, along with any other information you ask for with the flags. If you do not specify a file or directory parameter, the ls command displays the contents of the current directory.	<pre>1 ls [options] [file/directory]</pre> <p>Example 1: Sorts the file names displayed in the order of last modification time. 't' is for displaying in reverse order</p> <pre>1 ls -lt 2 ls -ltr</pre> <p>Example 2: Displays hidden files</p> <pre>1 ls -a</pre>
mkdir	Used to create one or more directories specified by the Directory parameter. Each new directory contains the standard entries dot () and dot dot (..). You can specify the permissions for the new directories with the -m Mode flag.	Create a new directory named "myfolder."
SELECT * FROM	Lists all the rows from the table to check if the data has been loaded from the file.	Select all rows from a table.
show tables	Used to see all the tables in the database that have been selected.	Show all tables in the selected database.
tar	Looks for archives on the default device (usually tape) unless you specify another device. When writing to an archive, the tar command uses a temporary file (the /tmp/tar* file) and maintains in memory a table of files with several links.	Create a tar archive of a directory.