# *Hive Commands*

## Overview of Hive and Cloudera

**Apache Hive:**

- **Apache Hive** is a data warehousing and SQL-like query language interface built on top of Hadoop.
- It facilitates the reading, writing, and managing of large datasets residing in distributed storage using .
- Hive converts -like queries into MapReduce jobs for execution, allowing data analysts to perform queries and analyses without needing to understand complex Hadoop code.

**Cloudera:**

- **Cloudera**  provides a comprehensive platform for data management and analytics, integrating various components of the Hadoop ecosystem.
- Unlike Hadoop, which is primarily a framework for processing large data sets in a distributed computing environment, cloudera offers an enterprise-level solution that combines multiple data processing tools and services.
- Cloudera includes tools for data ingestion, storage, processing, and visualization, simplifying big data management.
- One notable feature within Cloudera is **Hue**, a web-based interface that streamlines interaction with Hadoop and Hive.

## Advantages of Hadoop

1. **Scalability:** Hadoop is designed to scale horizontally, allowing organizations to add more nodes to their cluster as data volumes grow. This makes it suitable for big data applications.

2. **Cost-Effective Storage:** Hadoop uses commodity hardware for storage, significantly reducing the costs associated with storing large datasets compared to traditional storage solutions.

3. **Flexibility:** Hadoop can process various types of data-structured, semi-structured, and unstructured—allowing organizations to derive insights from diverse data sources.

4. **Fault Tolerance:** Hadoop is built to handle failures. It automatically replicates data across multiple nodes, ensuring data integrity and availability even in the event of hardware failures.

5. **Large Ecosystem:** The Hadoop ecosystem includes various tools (like Pig, Hive, and HBase) that extend its functionality, providing different capabilities for data processing, analysis, and storage.

## Disadvantages of Hadoop

1. **Complex Setup:** Setting up a Hadoop cluster can be complex and time-consuming, requiring expertise in configuration and management.

2. **Performance:** While Hadoop is efficient for batch processing, it may not perform well for low-latency data processing needs, which can be a limitation for certain applications.

3. **Resource Management:** Hadoop relies on the YARN resource management framework, which can become challenging to configure and manage, particularly in large clusters.

4. **Learning Curve:** Users must familiarize themselves with Hadoop's architecture and tools, which can be daunting for those new to big data technologies.

## Advantages of Cloudera over Hadoop

1. **Ease of Use:** Cloudera provides a user-friendly interface (Hue) that allows users to interact with the Hadoop ecosystem without extensive command-line knowledge. This is particularly beneficial for those who are new to big data technologies.

2. **Integrated Tools:** Cloudera integrates various tools and services, such as Hive, Impala, HBase, and Spark, offering a comprehensive solution for data processing and analytics. This integration simplifies workflows and reduces the need for third-party tools.

3. **Enterprise Features:** Cloudera offers enterprise-level features such as enhanced security, data governance, and comprehensive support. This is especially important for organizations dealing with sensitive data or requiring regulatory compliance.

4. **Scalability:** Cloudera is designed to scale seamlessly, making it suitable for organizations that anticipate growth in their data volume or complexity.

5. **Support and Documentation:** Cloudera provides robust support and extensive documentation, which can significantly reduce the learning curve and help organizations troubleshoot issues effectively.

## Disadvantages of Cloudera

1. **Cost:** Cloudera's enterprise solutions can be expensive, which may not be feasible for smaller organizations or startups. While there is a free version, it may lack certain features available in the paid versions.

2. **Complexity:** The comprehensive nature of Cloudera can introduce complexity, particularly for users who only require basic functionalities. The wide array of tools and features may overwhelm beginners.

3. **Resource Intensive:** Running Cloudera can be resource-intensive, requiring significant hardware capabilities, especially for larger datasets and workloads.

## Steps to Create a Hive Database and Load Data

1. **Show Existing Databases:**

   *SHOW DATABASES;*

2. **Create a New Database:**

   *CREATE DATABASE madhurima_database;*

3. **Use the Newly Created Database:**

   *USE madhurima_database;*

4. **Create a Table for Student Records:**

   *CREATE TABLE Student_Records (*

   *Roll_No INT,*

   *Percentage FLOAT,*

   *Name STRING*

   *) ROW FORMAT DELIMITED*

   *FIELDS TERMINATED BY ',';*

5. **Show Existing Tables:**

   *SHOW TABLES;*

6. **Load Data into the Table:**

   *LOAD DATA LOCAL INPATH 'file:///home/cloudera/Downloads/sample_student_data.txt'*

   *INTO TABLE Student_Records;*

7. **Select Data from the Table:**

   *SELECT * FROM Student_Records;*

## Additional Tips:

- **Using GitHub for Code Management:**
  If we're unable to copy and paste commands directly into the terminal, a workaround is to upload all necessary files and code to a private GitHub repository using the Firefox browser. We can create a new issue to type the content and copy it from there. Although downloading files may not be supported, we can select and copy the desired data from the browser.

- **Hue Login Credentials:**
  If we have not set a custom password, use the default username and password:

  - **Username:** cloudera

  - **Password:** cloudera

*MADHURIMA RAWAT*

- **File                      Creation                      in                      Cloudera:**
  To create a new text file, navigate to Cloudera:

  o   Go to Applications > System Tools > File Browser > Downloads

  o   Right-click, select Create Document > Create Empty File > Create a text file
      with a .txt extension (the default format for text files).

  o   Open the file and add the contents.

By following these steps, we can efficiently create and manage warehouse Hive database within
Cloudera, leveraging the advantages of both Hive and Cloudera's integrated platform for
handling big data workloads, while also understanding the broader context of Hadoop's
capabilities and challenges.