

What is Hive?

Apache Hive is a data warehousing and SQL-like query system built on top of Hadoop. It allows users to write queries in a language called HiveQL (similar to SQL) to interact with data stored in Hadoop Distributed File System (HDFS). Hive abstracts the complexity of Hadoop's underlying MapReduce framework, allowing users to query and manage large datasets without needing to write low-level MapReduce code.

Basic Components of Hive:

1. **HiveQL (Hive Query Language):**
 - A SQL-like query language used to perform operations like data retrieval, insertion, and analysis.
2. **Metastore:**
 - Stores metadata about the databases, tables, columns, data types, and table locations in HDFS. The metastore allows Hive to retrieve information quickly when executing queries.
3. **Driver:**
 - Acts as a controller for managing the execution of the queries, compiling the queries into MapReduce jobs, and executing them.
4. **Compiler:**
 - Converts HiveQL queries into MapReduce tasks. The query is parsed and optimized before sending it to the execution engine.
5. **Execution Engine:**
 - Executes the MapReduce jobs generated by the compiler and returns the results back to the user.
6. **Hive Server:**
 - It allows clients to connect to Hive and submit queries. This server can handle multiple requests from various clients simultaneously.

7. CLI (Command Line Interface):

- Provides an interface where users can execute HiveQL queries and interact with the Hive system.

Architecture of Hive:

The architecture of Hive is composed of several key components that interact with Hadoop and provide a user-friendly interface for managing large datasets.

1. User Interface:

- Users interact with Hive via a web UI, command line interface (CLI), or JDBC/ODBC interface. This is where queries are written and submitted to Hive.

2. HiveQL Process Engine:

- Once a query is submitted, the HiveQL process engine parses the HiveQL, checks for errors, and compiles it into a directed acyclic graph (DAG) of MapReduce tasks.

3. Metastore:

- The metastore holds all metadata about the tables in Hive, such as table schemas, column types, and table locations. This metadata is crucial for query execution, as it helps in locating the data and understanding the structure of the datasets.

4. Execution Engine:

- The execution engine is responsible for the actual execution of the queries. It takes the compiled MapReduce jobs and runs them in the Hadoop cluster.

5. HDFS (Hadoop Distributed File System):

- Hive stores its data on HDFS, Hadoop's underlying file system. All data queried or inserted via Hive is stored in HDFS.

6. Driver:

- The driver is responsible for managing the query lifecycle. It coordinates between different stages of the query, including the compilation, optimization, and execution phases. It also handles the communication between Hive and the metastore.
-

Usage of Hive:

1. Data Warehousing:

- Hive is primarily used for performing data warehousing tasks on large datasets. It provides the ability to run complex analytical queries over large amounts of structured and semi-structured data.

2. Handling Structured Data:

- Hive is ideal for handling structured data, with support for multiple formats such as text files, ORC, Parquet, etc. It allows users to organize data into tables, partitions, and buckets.

3. Batch Processing:

- Hive is designed for batch processing of data. It is not meant for real-time querying but excels in analyzing historical data stored in HDFS.

4. SQL-Like Interface for Hadoop:

- Hive makes Hadoop accessible to analysts and data scientists who are familiar with SQL, offering a more user-friendly way to interact with large datasets without needing to learn MapReduce programming.

5. Integrating with BI Tools:

- Hive can be integrated with business intelligence (BI) tools that use JDBC or ODBC for reporting and data analysis, making it a suitable choice for organizations looking to use Hadoop for business analytics.

Installation of HIVE:

- Downloading hive files

```
lnaad@lnaad:~$ wget https://downloads.apache.org/hive/hive-3.1.3/apache-hive-3.1.3-bin.tar.gz
--2024-08-23 17:54:20-- https://downloads.apache.org/hive/hive-3.1.3/apache-hive-3.1.3-bin.tar.gz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.208.237, 2a01:4f8:10a:39da::2, ...
Connecting to downloads.apache.org (downloads.apache.org)[135.181.214.104]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 326940667 (312M) [application/x-gzip]
Saving to: 'apache-hive-3.1.3-bin.tar.gz'

apache-hive-3.1.3-bin.tar.gz      100%[=====] 311.79M  5.09MB/s   in 70s
2024-08-23 17:55:30 (4.46 MB/s) - 'apache-hive-3.1.3-bin.tar.gz' saved [326940667/326940667]

lnaad@lnaad:~$
```

- Locating hive folder after extraction

```
lnaad@lnaad:~$ tar xzf apache-hive-3.1.3-bin.tar.gz
lnaad@lnaad:~$ ls
apache-hive-3.1.3-bin  apache-hive-3.1.3-bin.tar.gz  'Assignments (TV) 05_LAB'  Desktop  Documents  Downloads  Music  Pictures  Public  snap  Templates  Videos
lnaad@lnaad:~$
```

- Setting path for hive in Hadoop

```
GNU nano 6.2
export HADOOP_HOME=/home/hadoop/hadoop
```

- Creating hive and warehouse folder in hdfs

```
hadoop@lnaad:~$ hdfs dfs -mkdir -p /user/hive/warehouse
hadoop@lnaad:~$ hdfs dfs -chmod g+w /user/hive/warehouse
hadoop@lnaad:~$ hdfs dfs -ls /user/hive
Found 1 items
drwxrwxr-x - hadoop supergroup          0 2024-08-23 18:30 /user/hive/warehouse
hadoop@lnaad:~$
```

- Updating warehouse part with desired values

```
</property>
<property>
  <name>hive.metastore.db.type</name>
  <value>DERBY</value>
  <description>
    Expects one of [derby, oracle, mysql, mssql, postgres].
    Type of database used by the metastore. Information schema &amp;ap; JDBCStorageHandler depend on it.
  </description>
</property>
</property>
<property>
  <name>hive.metastore.warehouse.dir</name>
  <value>/user/hive/warehouse</value>
  <description>Location of default database for the warehouse</description>
</property>
<property>
  <name>hive.metastore.warehouse.external.dir</name>
  <value></value>
  <description>Default location for external tables created in the warehouse. If not set or null, then the normal warehouse location will be used as the default location.</description>
</property>
<property>
  <name>hive.metastore.uris</name>
  <value></value>
  <description>Thrift URI for the remote metastore. Used by metastore client to connect to remote metastore.</description>
</property>
<property>
  <name>hive.metastore.uri.selection</name>
```