# Comprehensive Lecture Notes: Apache Hive and Big Data Warehousing

Based on Lecture by Dr. Tariq Mahmood
Fall 2025

# Contents

# 1   Introduction: The Motivation for Hive

## 1.1   The Challenge with MapReduce

Before the advent of Apache Hive, processing data on Hadoop required writing raw MapReduce jobs. While MapReduce is a powerful paradigm for distributed computing, it presented significant barriers to entry for data analysts and businesses:

- **Complexity:** MapReduce requires writing verbose code in Java. Modern alternatives include PySpark, but the underlying complexity of distributed systems remains.

- **Imperative Programming:** Developers must define *how* to get the data (step-by-step procedures) rather than *what* data they want (declarative).

- **Lack of SQL Support:** Standard business intelligence (BI) tools and analysts are fluent in SQL, not Java.

- **Latency:** MapReduce is optimized for batch processing, resulting in high latency. Simple tasks often require a long execution stack.

- **No Schema Management:** Raw MapReduce operates on files without inherent knowledge of the data's structure (schema), increasing the maintenance burden.

## 1.2   Enter Apache Hive

Apache Hive was developed (initially by Facebook) to address these challenges by bringing data warehousing capabilities to the Hadoop ecosystem.

**Definition:** Hive is a data warehouse infrastructure built on top of Hadoop that facilitates data summarization, ad-hoc querying, and the analysis of large datasets stored in Hadoop compatible file systems (HDFS).

**Core Capabilities:**

- **SQL-Like Interface:** Hive provides a query language called **HiveQL (HQL)**, which closely resembles SQL-92.

- **Abstraction:** It translates HQL queries into distributed processing jobs (MapReduce, Apache Tez, or Apache Spark), hiding the complexity from the user.

- **Structure:** It projects structure onto data stored in HDFS, allowing users to query data as tables.

# 2   Hive Architecture

Hive is not a database in the traditional sense; it is a translation layer. Its architecture consists of several key components:

## 2.1   Major Components

1. **User Interfaces:**

    - **CLI (Command Line Interface):** The legacy shell.
    - **Beeline:** The modern, recommended CLI that connects via JDBC.
    - **Web UI:** Graphical interfaces (e.g., Hue, Ambari).

2. **Hive Driver:** The "brain" of Hive. It receives the query and manages the lifecycle:

- **Compiler:** Parses the query, performs semantic analysis, and generates an execution plan.

- **Optimizer:** Optimizes the logical plan (e.g., predicate pushdown, map-side joins).

- **Executor:** Submits the physical plan (DAG of jobs) to the cluster (YARN).

3. **Metastore:** The central repository of Hive metadata. It stores information about tables, columns, partitions, schema definitions, and data locations.

- Unlike the data itself (stored in HDFS), the metadata is stored in a relational database (RDBMS) like MySQL, PostgreSQL, or Derby.

## 2.2  Execution Engines

Hive queries must be executed by a distributed engine. Hive supports:

- **MapReduce:** The original engine. Stable but slow due to disk I/O between stages.

- **Apache Tez:** A DAG (Directed Acyclic Graph) based engine. It optimizes performance by eliminating unnecessary write-to-disk steps between map and reduce phases.

- **Apache Spark:** Uses in-memory processing for significantly faster performance.

## 2.3  Metastore Configurations

The Metastore can be deployed in three modes:

1. **Embedded:** The Metastore and Hive service run in the same JVM, using a local Derby database. *Limitation:* Single user only; used for testing.

2. **Local:** The Hive service connects directly to an external RDBMS (e.g., MySQL) on the same or remote machine.

3. **Remote (Production Standard):** A dedicated Metastore Server runs as a separate process. Hive clients connect to this server via Thrift. This allows centralized security, scalability, and multiple clients (Spark, Impala) to share metadata.

# 3  Data Model and Organization

Hive organizes data into a hierarchy that maps logical concepts to physical HDFS directories.

## 3.1  Hierarchy

1. **Databases:** Namespaces that group tables. Default location: `/user/hive/warehouse/dbname.db`.

2. **Tables:** Homogeneous units of data with the same schema. Corresponds to a directory in HDFS.

3. **Partitions:** Virtual columns that define directory hierarchies to segregate data.

4. **Buckets (Clusters):** Files within partition directories, created by hashing a column.

## 3.2   Partitioning

Partitioning is a coarse-grained method of dividing data.

- **Concept:** Creates sub-directories based on the value of a column (e.g., Country, Date).

- **Benefit:** *Partition Pruning.* When a query filters by the partition column (e.g., `WHERE country='US'`), Hive scans *only* the 'US' directory, ignoring the rest. This drastically reduces I/O.

- **Example Path:** `/warehouse/sales/country=US/date=2025-10-01/`

## 3.3   Bucketing (Clustering)

Bucketing is a fine-grained method of dividing data.

- **Concept:** Data is distributed into fixed-size files (buckets) based on the hash of a column (e.g., `hash(user_id) mod N`).

- **Benefit:**
  - **Efficient Sampling:** Allows querying a random sample of data without scanning the whole dataset.
  - **Map-Side Joins:** If two tables are bucketed by the join key into the same number of buckets, Hive can join them efficiently at the Map stage, skipping the Reduce/Shuffle phase.

# 4   HiveQL: Key Concepts and Operations

## 4.1   Schema-on-Read vs. Schema-on-Write

- **RDBMS (Schema-on-Write):** The schema is enforced when data is loaded. Loading is slower (checking constraints), but query performance is optimized.

- **Hive (Schema-on-Read):** Data is stored as raw files. The schema is applied only when the data is read by a query.

- **Implication:** Hive is very flexible (you can change the schema without rewriting data) and loads data fast, but queries may fail at runtime if data types mismatch.

## 4.2   Managed vs. External Tables

This is a critical distinction in Hive data lifecycle management.

| Managed (Internal) Table | External Table |
|---|---|
| Hive manages both Metadata and Data. | Hive manages Metadata; Data is external. |
| Data is stored in the Hive Warehouse directory. | Data can be stored anywhere in HDFS. |
| **DROP TABLE** deletes both the metadata and the actual data files. | **DROP TABLE** deletes only the metadata. The data files remain safe in HDFS. |
| Used when Hive is the sole owner of the data. | Used when data is shared with other tools (e.g., Spark, Pig) or raw logs. |

### 4.3 Data Types

- **Primitive:** TINYINT, INT, BIGINT, BOOLEAN, FLOAT, DOUBLE, STRING, TIMESTAMP, DECIMAL, BINARY.

- **Complex:**

  - ARRAY<TYPE>: Ordered collection. Accessed via index [n].
  - MAP<KEY, VALUE>: Key-value pairs. Accessed via ['key'].
  - STRUCT<col:TYPE, ...>: Collection of named fields. Accessed via dot notation .col.

## 5 Working with Hive: Syntax and Examples

### 5.1 Creating Tables

```
1  -- Creating a Partitioned and Bucketed Table using ORC format
2  CREATE TABLE ecommerce_sales (
3      order_id STRING ,
4      customer_id STRING ,
5      amount FLOAT
6  )
7  PARTITIONED BY (country STRING)
8  CLUSTERED BY (customer_id) INTO 8 BUCKETS
9  STORED AS ORC;
```

### 5.2 Loading Data

Hive does not typically use row-level inserts. It uses bulk loading.

- **LOAD DATA:** Moves or copies files into the Hive directory.

```
1      -- Load from local filesystem
2      LOAD DATA LOCAL INPATH '/home/user/data.csv'
3      INTO TABLE sales PARTITION(country='US');
4
```

- **INSERT OVERWRITE:** Loads data by querying another table or raw file.

```
1      INSERT OVERWRITE TABLE sales_summary
2      SELECT product , SUM(amount) FROM sales GROUP BY product;
3
```

### 5.3 SQL Limitations and Extensions

While HQL is SQL-like, it has historical limitations compared to standard RDBMS:

- **No Update/Delete (Historically):** Standard Hive assumes immutable data. (Note: Modern Hive 3.x supports ACID transactions for Update/Delete on ORC tables with bucketing).

- **Latency:** Not suitable for OLTP (Real-time transaction processing).

- **Extensions:**

  - **Multi-Table Insert:** Scan a table once and insert into multiple tables/partitions.
  - **CTAS (Create Table As Select):** Create and populate a table in one step.

# 6 Optimization and File Formats

## 6.1 File Formats

Choosing the right storage format impacts performance and storage cost.

- **TEXTFILE:** Default. Human-readable, bulky, slow to parse.

- **SEQUENCEFILE:** Binary, splittable, supports block compression.

- **Columnar Formats (Recommended for Analytics):**

  - **RCFile / ORC (Optimized Row Columnar):** High compression, indexes, optimized for Hive. Supports ACID.
  - **Parquet:** Industry standard columnar format, excellent interoperability with Spark.

## 6.2 Views

- **Logical Views:** Saved queries. Run every time the view is referenced.

- **Materialized Views:** Pre-computed results stored on disk.

```
1    CREATE MATERIALIZED VIEW sales_agg AS
2    SELECT id, SUM(amount) FROM sales GROUP BY id;
3
4    -- Needs rebuilding to refresh data
5    ALTER MATERIALIZED VIEW sales_agg REBUILD;
6
```

Used to speed up expensive aggregations.

# 7 Summary

Apache Hive democratized Big Data by allowing analysts to query massive HDFS datasets using familiar SQL syntax. While it is not a replacement for real-time databases (OLTP), it serves as the de-facto standard for SQL-based ETL and Data Warehousing on Hadoop.

- **Use Partitioning** to reduce the amount of data scanned.

- **Use Bucketing** for sampling and join optimization.

- **Use Columnar Formats** (ORC/Parquet) for compression and performance.

- **Use External Tables** for raw data to prevent accidental data loss.