



# Apache Hive

Apache Hive is a data warehouse system built on top of Hadoop.

It allows users to query large datasets stored in HDFS using a SQL like language called HiveQL (HQL)

Just with MR, working with hive is going to make our life a lot easy.

Why was Hive Created?

Complex Java  
based MR  
Program

SQL kind of  
interface to  
work on the  
data

developed at FB  
(2010)

SQL      no Java or  
          MR code

Analogy

Siliconian with a Mobile / iPad / smart device

PBs of  
data

# How Hive Makes Big Data Processing Easier

higher level abstraction layers  
over Hadoop

SQL like → MR, Spark or Tez  
HQL

```
Select city, count(*)  
from customers  
group by city
```

```
public class CustomerCount extends Mapper<LongWritable,  
    Text, Text, IntWritable> {  
    public void map(LongWritable key, Text value,  
        Context context) {  
        String[] fields = value.toString().split(",");  
        context.write(new Text(fields[2]), new  
            IntWritable(1)); // City as key  
    }  
  
    public class CustomerCountReducer extends Reducer<Text,  
        IntWritable, Text, IntWritable> {  
        public void reduce(Text key, Iterable<IntWritable>  
            values, Context context) {  
            int sum = 0;  
            for (IntWritable val : values) {  
                sum += val.get();  
            }  
            context.write(key, new IntWritable(sum));  
        }  
    }  
}
```

Architecture  
in-depth knowledge

Without Hive	With Hive
Writing Java-based MapReduce programs to process data	Writing SQL-like HiveQL queries
Complex and time-consuming development	Easier and faster query execution
Requires knowledge of Java and MapReduce	Anyone familiar with SQL can use Hive

# Some Common Question/Misconception about Hive

1. Hive is a Database just like MySQL or PostgreSQL.

Hive is not a database but a data warehouse tool built on top of Hadoop.

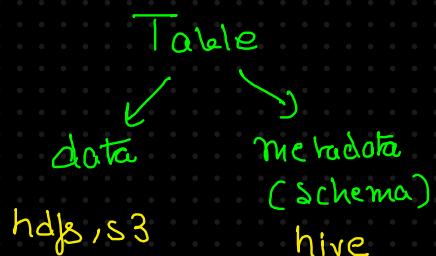
Not storing any data like MySQL

→ interface to query large datasets stored in HDFS

Storage

Query execution

hive is a metastore

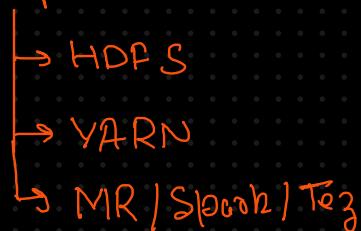


\* In big data world, any tech that supports Querying, might not be a DB

2. Hive is a replacement for Hadoop.

Hive is built on top of Hadoop & hadoop

Components



like a  
translator  
& SQL like code  
instead of MR

### 3. Hive Queries run like Normal SQL Queries in Databases

Hive does not process queries instantly  
like MySQL → MR/spark/Tez

~AI  
converts

### 4. Hive can Perform Row-level Transaction like MySQL.

↓  
dw  
olap

↓  
dB  
oltp

PB<sub>s</sub>  
↓

Hive is not meant for row-level operations



### 5. Hive only works with HDFS

Wrong. Hive can store & process data from

HDFS

S3

Amazon datalake

GCS

### 6. Hive Tables Work Like Normal Database Tables.

Hive is a metastore and not storing tables.

Feature	Hive Advantage
SQL-Like Interface	Allows non-programmers to work with big data.
Scalability	Handles <b>petabytes</b> of data.
Optimized Execution	Uses MapReduce, Tez, or Spark for distributed processing.
Storage Flexibility	Works with HDFS, S3, Azure Blob, GCS.
Schema Flexibility	Schema-on-read enables analyzing raw files.

# Hive Practical on Google Dataproc

We can run hive commands via the hive terminal or the Beeline (recommended)

Feature	Hive CLI	Beeline (Recommended)
Connection Type	Directly connects to Hive	Uses JDBC for remote execution
Security	No authentication	Supports authentication (LDAP, Kerberos)
Multiple Connections	Single session only	Supports multiple concurrent sessions
Performance	More resource-heavy	Optimized for query performance
Recommended?	✗ Deprecated	✓ Yes, for production use

## Hive Practical - 2

```
root@my-cluster-m:/# hadoop fs -ls /user/hive/warehouse/
Found 5 items
drwxr-xr-x  - root      hadoop          0 2025-02-08 03:07 /user/hive/warehouse/customers_100
drwxr-xr-x  - root      hadoop          0 2025-02-06 11:40 /user/hive/warehouse/customers_500mb
drwxr-xr-x  - anonymous hadoop          0 2025-02-16 05:13 /user/hive/warehouse/ecommerce.db
drwxr-xr-x  - root      hadoop          0 2025-02-17 11:39 /user/hive/warehouse/ecommerce_new.db
drwxr-xr-x  - root      hadoop          0 2025-02-09 05:38 /user/hive/warehouse/test_db.db
root@my-cluster-m:/#
```

## Accessing metadata in Hive

