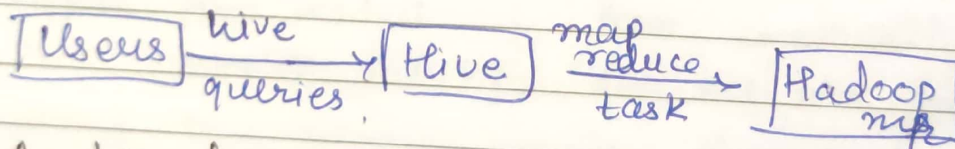# Hive

* **Why Hive?**

  **Problem** → For processing & analyzing data users found it difficult to code as not all of them were well versed with coding lang.
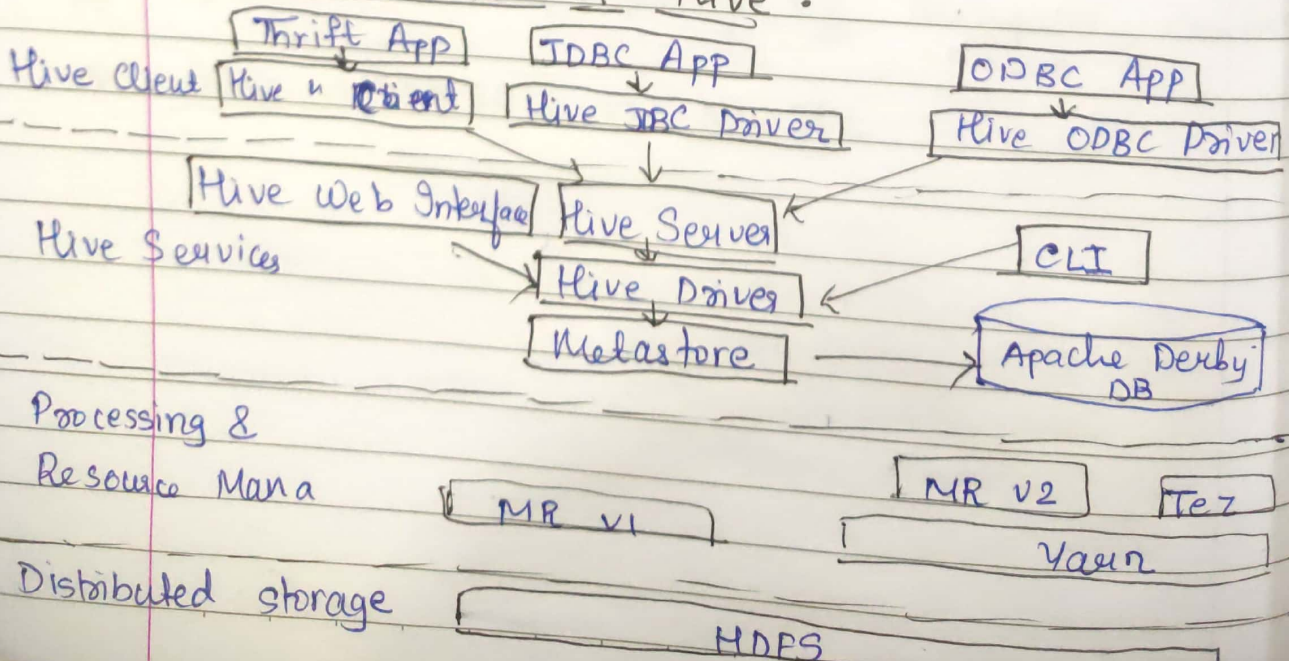
  **Solution** → Users required a lang similar to SQL which was well known to all
  → Hive QL.

* **What is Hive?**

  → Hive is a data warehouse sys which is used for quering & analyzing large datasets stored in HDFS.
  → Hive uses a query language call Hiv QL similar to SQL.

  $\boxed{Users} \xrightarrow[queries]{hive} \boxed{Hive} \xrightarrow[task]{map\ reduce} \boxed{Hadoop\ mr}$

* **Architecture of Hive:**

| | | |
|---|---|---|
| Hive Client | Thrift App → Hive Client | JDBC App → Hive JBC Driver | ODBC App → Hive ODBC Driver |

Hive Web Interface | Hive Server

Hive Services

Hive Driver

Metastore → Apache Derby DB

CLI

Processing & Resource Mana

MR V1     MR V2   Tez

Yarn

Distributed storage    HDFS

→ Hive Client supports different types of client app in diff lang for performing queries

→ Thrift is a s/f framework. Hive server is based on thrift, so it can serve the req from all the programming lang that supports thrift.

→ JDBC Java DB Connectivity app is connected through JDBC driver.

→ ODBC Open DB Connectivity app is connected through ODBC driver.

→ All the client requests are submitted to the Hive Server.

→ GUI is provided to execute Hive queries

→ Commands are executed directly on CLI.

→ Hive driver is responsible for all the queries submitted

Steps:
1) Compiler: Hive driver passes query to compiler where it is checked & analyzed
2) Optimizer: Optimized logical plan in the form of a graph of MR & HDFS tasks is^ obtained
3) Executor: the tasks are executed.

→ Metastore is a repository for Hive metadata. Store metadata for hive tables.
→ Hive uses MR framework to process queries of Hive.
→ default loc for Hive metastore db
→ to keep metadata db files isolated for each instance ^
→ can be useful while running multiple instances of Hive on same mac, each working with its own tables & metadata
→ default db for metastore is Derby.

## * Dataflow in Hive :

```
                                      Hive                     9. send result ↗  ┌─────────┐
                                                                                  │Execution│
                                                                  6. execute plan│ Engine  │──→ Hadoop
                                                                                  └─────────┘   ⇠- -⇢
                                                                         ↑              ↑        ┌──────┐
                   ┌───┐   1. execute query →  ┌────────┐                                        │ MR + │
                   │   │                       │ Driver │                                        │ HDFS │
                   │   │   8. fetch results →                                                    └──────┘
                   │UI │                              │      ↑ 5. sendplan
                   │   │   2. get plan                │      │                      7. Metdata
                   │   │   ( plan refers            │      │                           ops
                   │   │    to query                 ↓      │   3. get Metadata →    ┌──────────┐
                   └───┘    execution)        ┌──────────┐  ───────────────────────→│ Metastore│
                                              │ Compiler │ ←──────────              └──────────┘
                                              └──────────┘    4. send md.
```
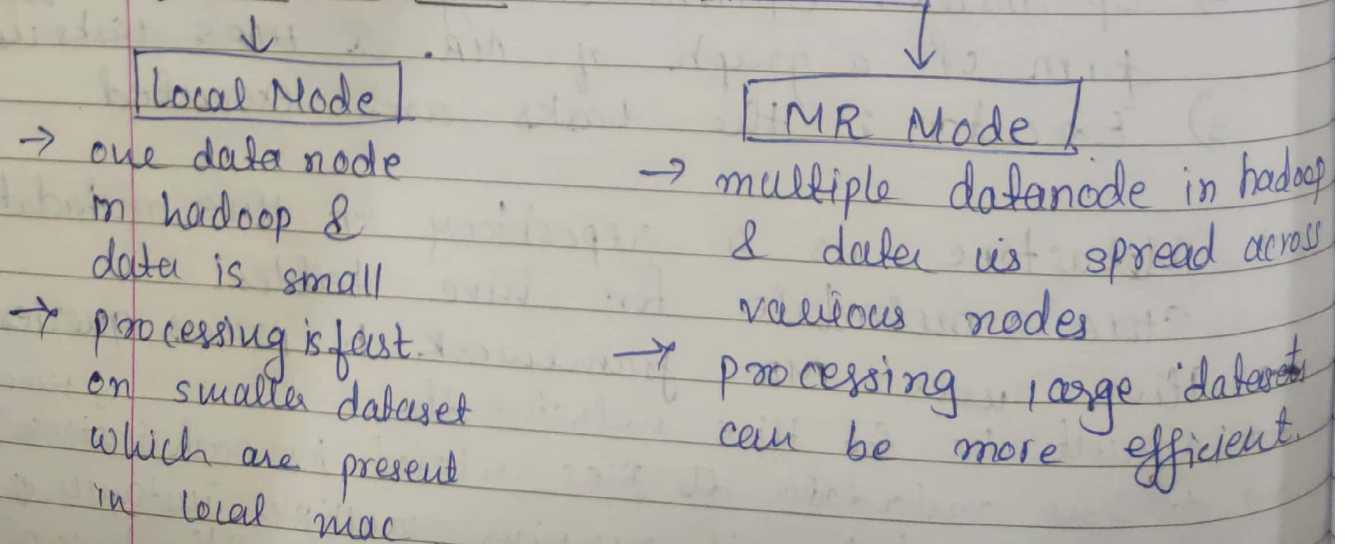
## * Hive data Modeling :

```
                                            ──────→ ┌─────────┐  data prese
                                                    │ Buckets │  in partiho
                          ┌────────┐                └─────────┘  can be further divid
                  ↓       │Partitions│                           into buckets for
          ┌────────┐      └────────┘       tables are organized    efficient
          │ Tables │                        into partitions for      query
          └────────┘                        grouping same type
       are in Hive are created              of data based on
       the same way it is                   partition key
       done in RDMS
```

## * Modes of Hive :

```
              ↓                                      ↓
      ┌────────────┐                          ┌──────────┐
      │ Local Node │                          │ MR Mode  │
      └────────────┘                          └──────────┘
```

→ one data node                    → multiple datanode in hadoop
   in hadoop &                        & data is spread across
   data is small                      various nodes
→ processing is fast.               → processing large datasets
   on smaller dataset                  can be more efficient.
   which are present
   in local mac

| Hive | RDMS |
|---|---|
| → schema on read | → on write |
| → data size in PB | → TB |
| → write once read many | → write & read many time |
| → data ware house | → database based on relational model |
| ↗ easily scalable at low cost | → not scalable at low cost. |

* **Features :-**

→ Tables are used similar to RDMS so easy to understand.
→ using hive QL, ~~query~~ multiple users can simultaneously query data
→ Hive QL - easier than long codes
→ supports variety of data formats

* **Appli :-**

→ Batch processing rather than real-time proce
→ Data warehousing
   organize & structure large volumes of data stored in HDFs
→ Large - scale Data Analysis when data is distributed ~~over~~ across Hadoop. cluster
→ Ad Hoc Quries :