

## Agenda

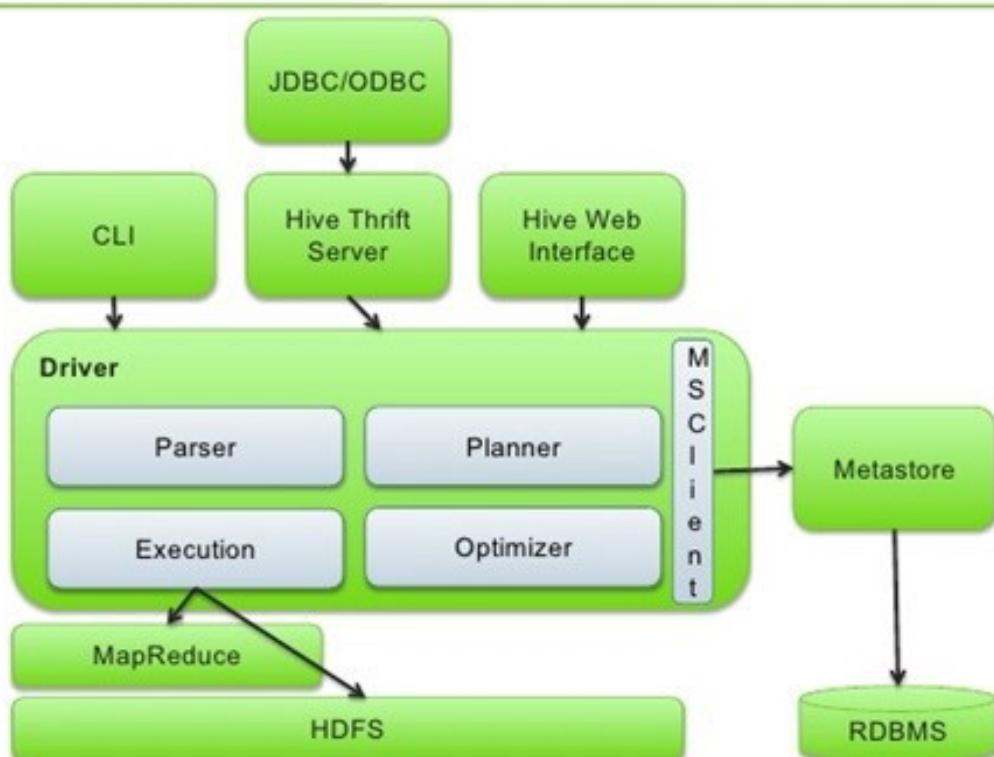
What are data warehouses and what is OLAP?  
How does Apache hive act as a data warehouse?  
What makes hive as Hive? (Architecture)  
How to interact with Hive?  
What are different HQL commands that are used commonly?  
What are the optimization techniques in Hive ?  
Use cases of hive/ where to use hive/why and why not

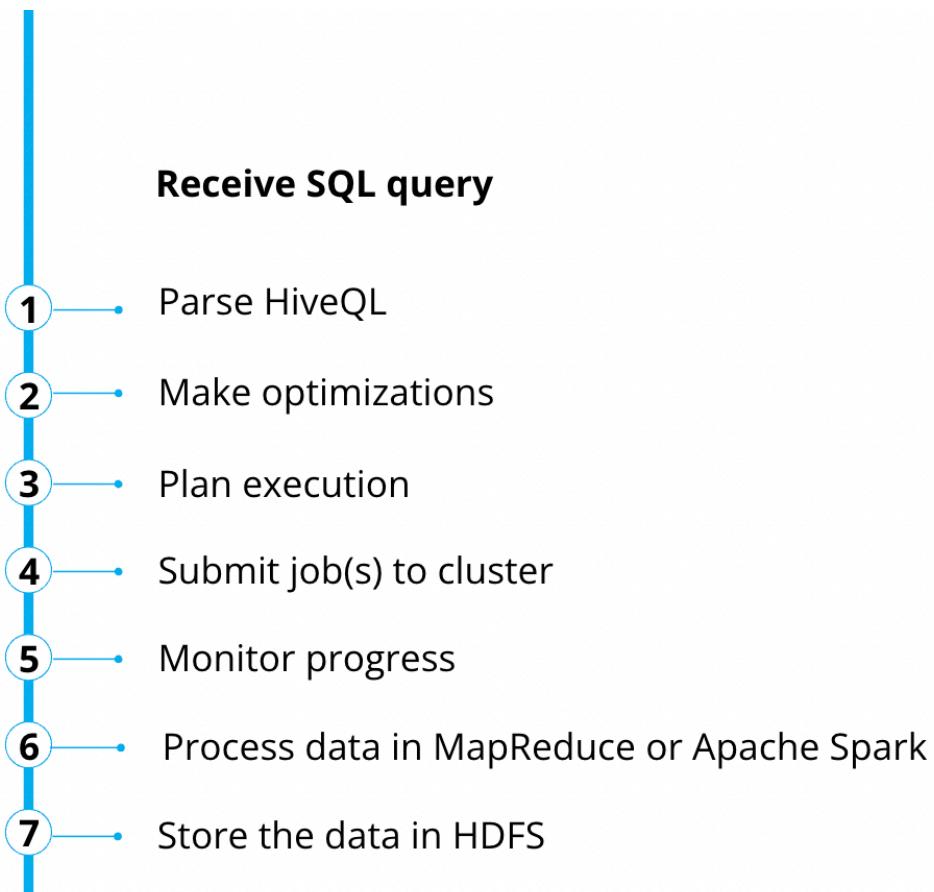


# AIRBNB DATA INFRA



## Apache Hive Architecture





+-----+
dim_hosts
+-----+
host_id (PK)
host_name
host_since
host_location
host_response_time
host_response_rate
+-----+
FK
v
+-----+
fact_listings
+-----+
listing_id (PK)

host_id (FK)
neighbourhood_id (FK)
room_type_id (FK)
price
minimum_nights
number_of_reviews
last_review
reviews_per_month
calculated_host_listings_count
availability_365
+-----+

|  
|  
| FK  
|  
v

+-----+
dim_neighbourhoods
+-----+
neighbourhood_id (PK)
neighbourhood_group
neighbourhood
latitude
longitude
+-----+

|  
|  
| FK  
|  
v

+-----+
dim_room_types
+-----+
room_type_id (PK)
room_type
+-----+

+-----+
dim_dates
+-----+

date_id (PK)	
date	
year	
quarter	
month	
day	
day_of_week	
+-----+	

Hive

↳ Datawarehouse developed by FB in 2002

↳ open source Apache

↳ High abstraction layer on top of MR.

↳ use HQL

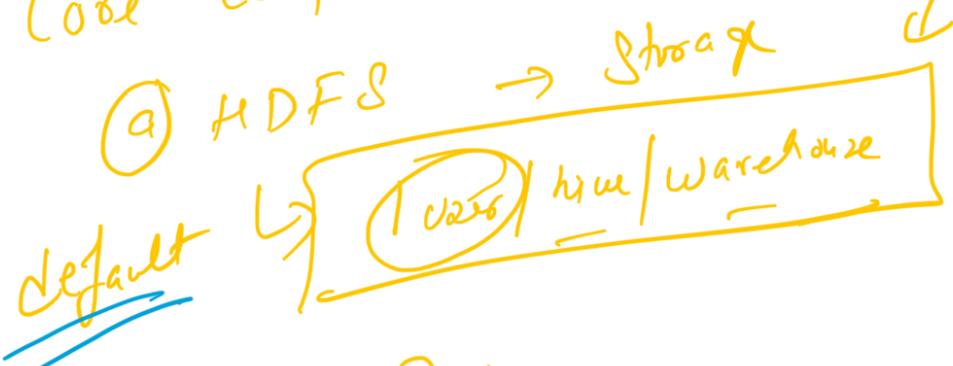
↳ Invertible for structured Data

Hive Architecture

(a) Hadoop Core Components (Hive-site.xml)

(a) HDFS → Storage

default



(b) Map Reduce

SQL → MR Jobs

(c) Metastore

... in all metadata

(b)    
↳ Consider part for now as information related to table, columns, partitions, is present in it.

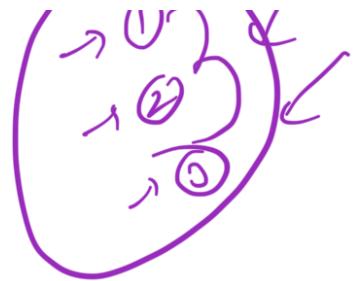
↳ Derby | MySQL | Oracle | PostgreSQL | SQLite.

(c) Hive Server :-  
↳ enables client to execute query against service.  
↳ Designed to provide better support for open API like (JDBC/ODBC)

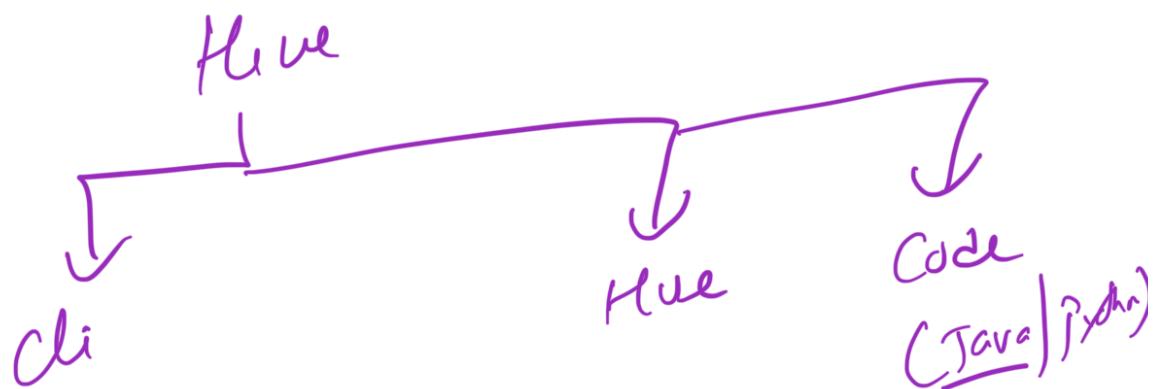
(d) Compiler  
↳ Parses the query  
↳ Does semantic analysis  
↳ Optimizes expressions

(e)    
↳ generates an execution plan (table mD)

(f) Execution Engine :-  
↳ Planner presents the execution plan to Engine.  
↳ ... in front of



- ↳ plan is in 3 stages.
- ↳ SQL/HQL will be converted to Map Reduce Jobs. (Java)



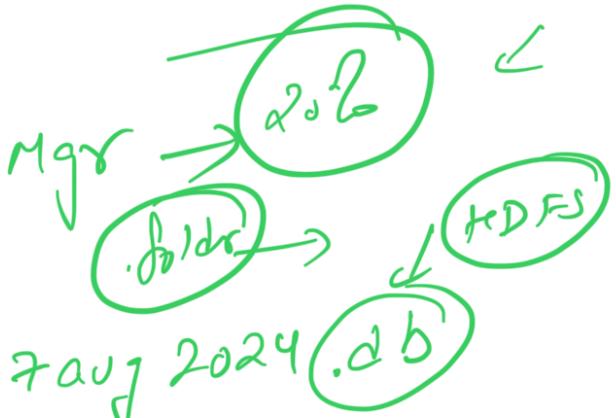
Hive *	Beeline
<ul style="list-style-type: none"> <li>→ Native cli tool</li> <li>→ No Concurrency support</li> <li>→ Fat client (TOFS)</li> <li>→ Higher resource usage as it connects directly to MS.</li> <li>→ No support to advanced</li> </ul>	<ul style="list-style-type: none"> <li>→ JDBC based CLI</li> <li>→ Concurrency support</li> <li>→ thin client (only relevant)</li> <li>→ lower resource usage as it runs on HiveServer2</li> <li>→ supports all advance features like ACID, CLAP.</li> </ul>

Question : Get the details of all travelers about their choice of stays in New York city i.e whether they book Entire home/apt or Private Room or they preferred Shared room

- An.1
- (a) Create database ;
  - (b) Create tables
  - (c) Insert Data

(d) Query

(e) Somech



✓ ① Database → 27 Aug 2024 .db

✓ ② Table

① Apache Tez

Optimizations  
2013-15

Set `hive.execution.engine = tez;`

② Apache Spark