https://hive.apache.org/



Master IPS 2021-2022



Need for High-Level Languages:

- Hadoop is great for large-data processing!
 - But writing Java programs for everything is verbose and slow
 - □ Not everyone wants to (or can) write Java code
- Solution: develop higher-level data processing languages
 - □ Hive: HQL is like SQL
 - □ Pig: Pig Latin is a bit like Perl



Introduction

- Hive is an ETL and data warehouse tool on top of Hadoop ecosystem and used for processing structured and semi structured data.
 - Rich data types (structs, lists and maps)
 - □ Efficient implementations of SQL filters, joins and group-by's on top of map reduce
 - Query language is HQL, variant of SQL
 - □ Tables stored on HDFS as flat files
 - Developed by Facebook, now open source

- Hive is a database present in Hadoop ecosystem performs DDL and DML operations, and it provides flexible query language such as HQL for better querying and processing of data.
- It provides so many features compared to RDMS which has certain limitations.



Different modes of Hive

Hive can operate in two modes depending on the size of data nodes in Hadoop. These modes are :

- 1 Local mode
- 2 Map reduce mode

When to use Local mode:

- If the Hadoop installed under pseudo mode (having one data node)
- If the data size is smaller (limited to a single local machine).
- Processing will be very fast on smaller data sets present in the local machine.

When to use Map reduce mode:

- If Hadoop is having multiple data nodes and data is distributed across different node.
- It will perform on large amount of data sets and query will execute in a parallel way.
- Processing of large data sets with better performance can be achieved through this mode.

By default, Hive will work on Map Reduce mode and for local mode you can have the following setting: To work in local mode set:

SET mapred.job.tracker=local;

From the Hive version 0.7 it supports a mode to run map reduce jobs in local mode automatically.



Requirements

Note: Hive versions <u>1.2</u> onward require Java 1.7 or newer. Hive versions 0.14 to 1.1 work with Java 1.6 as well. Users are strongly advised to start moving to Java 1.8 (see <u>HIVE-8607</u>).

- Hadoop 2.x (preferred), 1.x (not supported by Hive 2.0.0
 - onward). Hive versions up to 0.13 also supported Hadoop
- □ 0.20.x, 0.23.x.

Hive is commonly used in production Linux and Windows environment. Mac is a commonly used development environment.

Installing Hive from a Stable

Release

- Start by downloading the most recent stable release of Hive from one of the Apache download mirrors (https://hive.apache.org/downloads.html).
- Next you need to unpack the tarball. This will result in the creation of a subdirectory named hive-x.y.z (where x.y.z is the release number):

```
$ tar -xzvf hive-x.y.z.tar.gz

the environment variable HIVE_HOME to point to the installation directory:

$ cd hive-x.y.z
$ export HIVE_HOME={{pwd}}}

ally, add $HIVE_HOME/bin to your PATH:

$ export PATH=$HIVE HOME/bin:$PATH
```



Hive Architecture

- Metastore: stores system catalog
- Driver: manages life cycle of HiveQL query as it moves thru'
 HIVE; also manages session handle and session statistics
- Query compiler: Compiles HiveQL into a directed acyclic graph of map/reduce tasks
- Execution engines: The component executes the tasks in proper dependency order; interacts with Hadoop
- HiveServer: provides Thrift interface and JDBC/ODBC for integrating other applications.
- □ Client components: CLI, web interface, jdbc/odbc inteface
- Extensibility interface include SerDe, User Defined Functions and User Defined Aggregate Function.



Hive Architecture

In Hive, the data is stored in HDFS and the table, database, schema, and other HQL definitions are stored in a metastore. The metastore could be any RDBMS database, such as MySQL or Oracle. Hive creates a database and a set of tables in metastore to store HiveQL definitions.

□ There are three modes of configuring a Metastore

:

- Embedded
- 2. Local
- 3. Remote



Internal Vs External Tables in Hive

- "Internal" or "Hive managed" tables are the default
 - Stored at HDFS directory/hive/warehouse by default
 - Hive assumes no other system is using
 This DATA
 - So a "drop table" statement removes the metadata AND the underlying HDFS data
- External tables can be specified
 - Point "LOCATION" to some HDFS directory path
 - "drop table" now only removes the metadata



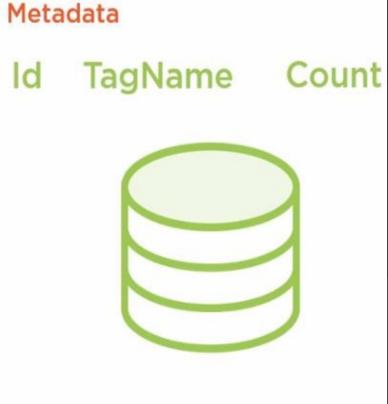
Data & Metadata

ld	TagName	Count
23 24 27 28	hadoop java solr json	2325 213123 32313 92223
4		



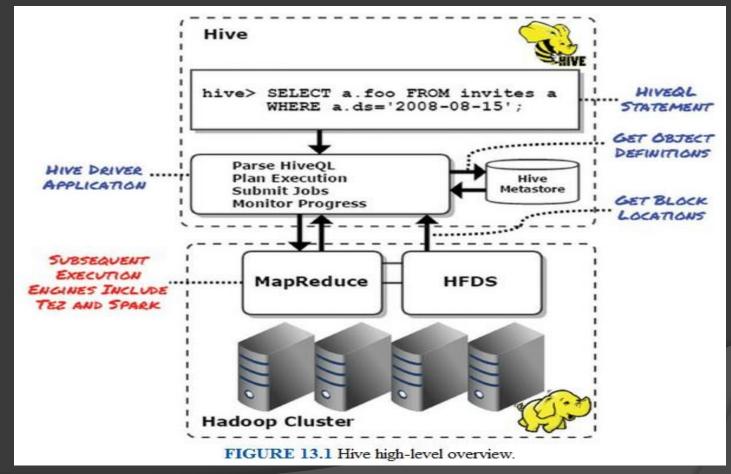
Data and MetaData

Data Metadata 2325 hadoop 23 213123 24 java 32313 27 solr 92223 28 json





Overview





- DataModel:
- Tables
 - □ Typed columns (int, float, string, boolean)
 - ☐ Also, list: map (for JSON-like data)
- Partitions
 - □ For example, range-partition tables by date
- Buckets
 - Hash partitions within ranges (useful for sampling, join optimization)



Metastore:

- Database: namespace containing a set of tables
- Holds table definitions (column types, physical layout)
- Holds partitioning information
- Can be stored in Derby, MySQL, and many other relational databases



Physical Layout:

- Warehouse directory in HDFS
 - ☐ E.g., /user/hive/warehouse
- Tables stored in subdirectories of warehouse
 - Partitions form subdirectories of tables
- Actual data stored in flat files
 - Control char-delimited text, or SequenceFiles
 - With custom SerDe, can use arbitrary format



Example:

- Hive looks similar to an SQL database
- □ Relational join on two tables:
 - □ Table of word counts from Shakespeare collection
 - Table of word counts from the bible
- SELECT s.word, s.freq, k.freq FROM shakespeare s.
- JOIN bible k ON (s.word = k.word) WHERE s.freq >= 1 AND k.freq >= 1 ORDER BY s.freq DESC LIMIT 10;
- □ the 25848 62394
- I 23031 8854
- and 19671 38985
- □ to 18038 13526
- of <u>16700 34654</u>
- □ a 14170 8057
- you 12702 2720
 - my 11297 4135
- in 10797 12445
- is 8882 6884



☐ Group By, Join, Multitable Insert ...

GROUP BY

```
hive> FROM invites a INSERT OVERWRITE TABLE events SELECT a.bar, count(*) WHERE a.foo > 0 GROUP BY a.bar; hive> INSERT OVERWRITE TABLE events SELECT a.bar, count(*) FROM invites a WHERE a.foo > 0 GROUP BY a.bar;
```

Note that for versions of Hive which don't include HIVE-287, you'll need to use COUNT (1) in place of COUNT (*).

JOIN

hive> FROM pokes t1 JOIN invites t2 ON (t1.bar = t2.bar) INSERT OVERWRITE TABLE events SELECT t1.bar, t1.foo, t2.foo;

MULTITABLE INSERT

```
FROM src

INSERT OVERWRITE TABLE dest1 SELECT src.* WHERE src.key < 100

INSERT OVERWRITE TABLE dest2 SELECT src.key, src.value WHERE src.key >= 100 and src.key < 200

INSERT OVERWRITE TABLE dest3 PARTITION(ds='2008-04-08', hr='12') SELECT src.key WHERE src.key >= 200 and src.key < 300

INSERT OVERWRITE LOCAL DIRECTORY '/tmp/dest4.out' SELECT src.value WHERE src.key >= 300;
```

Zeppelin Build and Tutorial Notebook



localhost:8080/#/notebook/2A94M5J1Z



Notebook -

Interpreter

Zeppelin Tutorial





```
%md
## Welcome to Zeppelin Tutorial.
##### This is a live tutorial, you can run the code yourself. (Shift-Enter to Run)

    bullet 2

 * Bullet 3
*emphasis*
**formatting**
_heading 1__
```

Welcome to Zeppelin Tutorial.

This is a live tutorial, you can run the code yourself. (Shift-Enter to Run)

- bullet 2
- Bullet 3

emphasis

formatting

heading 1

Took 0 seconds



METADATA

Function	Hive
Selecting a database	USE database;
Listing databases	SHOW DATABASES;
Listing tables in a database	SHOW TABLES;
Describing the format of a table	DESCRIBE (FORMATTED EXTENDED) table;
Creating a database	CREATE DATABASE db_name;
Dropping a database	DROP DATABASE db_name (CASCADE);

EZ ZAHOUT Aderrahmae



CURRENT SQL COMPATIBILITY

Hive SQL Datatypes	Hive SQL Semantics	12.027
INT	SELECT, LOAD, INSERT from query	Hive 0.10
TINYINT/SMALLINT/BIGINT	Expressions in WHERE and HAVING	Hive 0.11
BOOLEAN	GROUP BY, ORDER BY, SORT BY	
FLOAT	Sub-queries in FROM clause	Future
DOUBLE	GROUP BY, ORDER BY	
STRING	CLUSTER BY, DISTRIBUTE BY	
TIMESTAMP	ROLLUP and CUBE	
BINARY	UNION	
ARRAY, MAP, STRUCT, UNION	LEFT, RIGHT and FULL INNER/OUTER JOIN	
DECIMAL	CROSS JOIN, LEFT SEMI JOIN	
CHAR	Windowing functions (OVER, RANK, etc.)	
VARCHAR	INTERSECT, EXCEPT, UNION DISTINCT	
DATE	Sub-queries in WHERE (IN/NOT IN, EXISTS/NOT EXISTS	
	Sub-queries in HAVING	



RETRIEVING

Function	MySQL	Hive
Retrieving Information (General)	SELECT from_columns FROM table WHERE conditions;	SELECT from_columns FROM table WHERE conditions;
Retrieving All Values	SELECT * FROM table;	SELECT * FROM table;
Retrieving Some Values	<pre>SELECT * FROM table WHERE rec_name = "value";</pre>	<pre>SELECT * FROM table WHERE rec_name = "value";</pre>
Retrieving With Multiple Criteria	<pre>SELECT * FROM TABLE WHERE rec1 = "value1" AND rec2 = "value2";</pre>	<pre>SELECT * FROM TABLE WHERE rec1 = "value1" AND rec2 = "value2";</pre>
Retrieving Specific Columns	SELECT column_name FROM table;	SELECT column_name FROM table;
Retrieving Unique Output	SELECT DISTINCT column_name FROM table;	SELECT DISTINCT column_name FROM table;

EZ ZAHOUT Aderrahmae



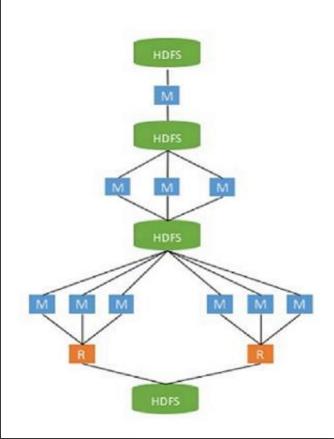
RETRIEVING

Sorting	SELECT col1, col2 FROM table ORDER BY col2;	SELECT col1, col2 FROM table ORDER BY col2;
Sorting Reverse	SELECT col1, col2 FROM table ORDER BY col2 DESC;	SELECT col1, col2 FROM table ORDER BY col2 DESC;
Counting Rows	SELECT COUNT(*) FROM table;	SELECT COUNT(*) FROM table;
Grouping With Counting	SELECT owner, COUNT(*) FROM table GROUP BY owner;	SELECT owner, COUNT(*) FROM table GROUP BY owner;
Maximum Value	SELECT MAX(col_name) AS label FROM table;	SELECT MAX(col_name) AS label FROM table;
Selecting from multiple tables (Join same table using alias w/"AS")	SELECT pet.name, comment FROM pet, event WHERE pet.name = event.name;	SELECT pet.name, comment FROM pet JOIN event ON (pet.name = event.name)



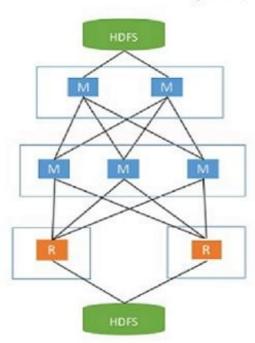
MapReduce and DAG engine

MapReduce



DAG engine

(Tez / Spark)



No intermediate DFS reads/writes!



HiveQL

DDL:

CREATE

DATABASE

CREATE TABLE

ALTER TABLE

SHOW TABLE

DESCRIBE

DML:

LOAD

TABLE

INSERT

QUERY:

SELECT

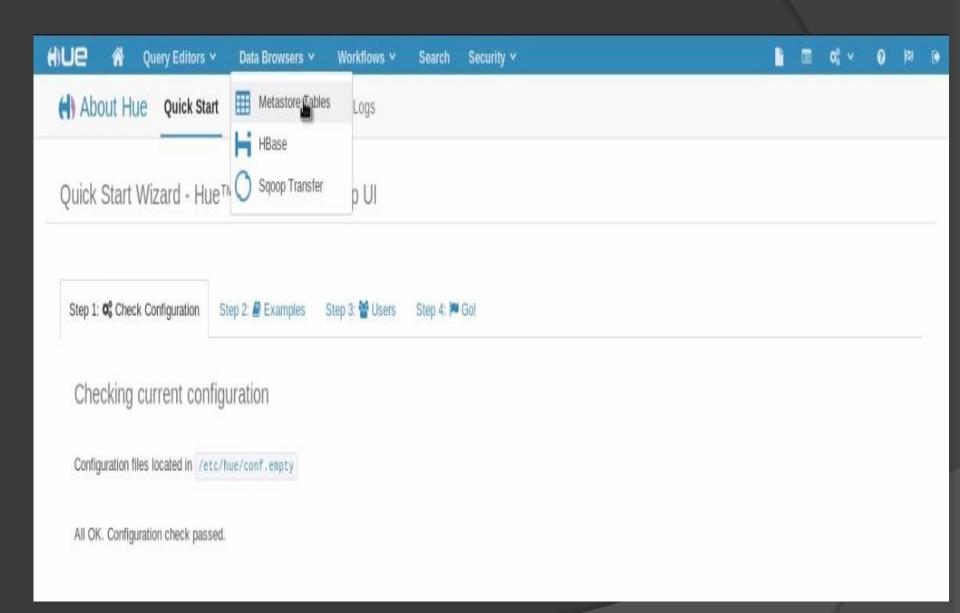
GROUP BY

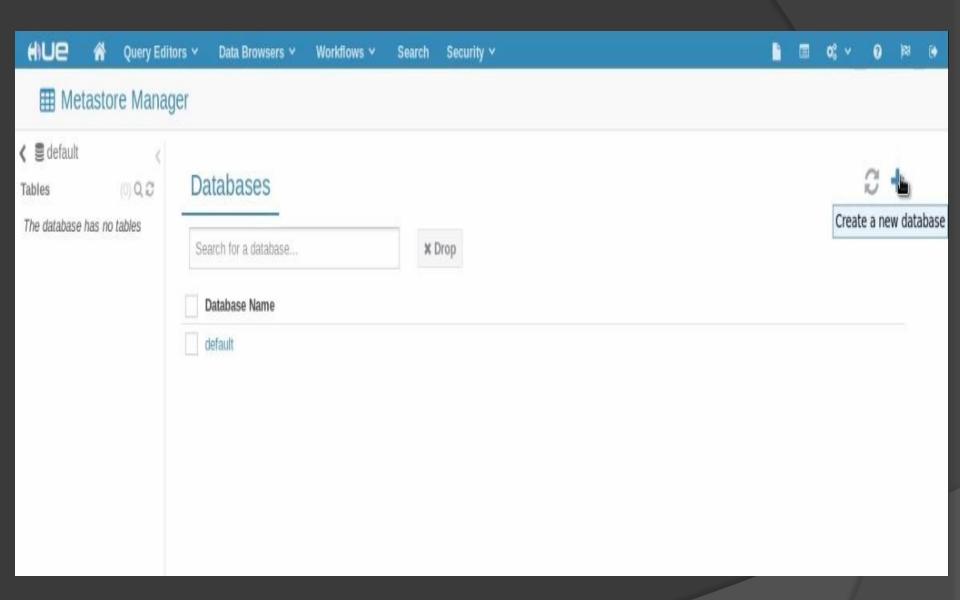
JOIN

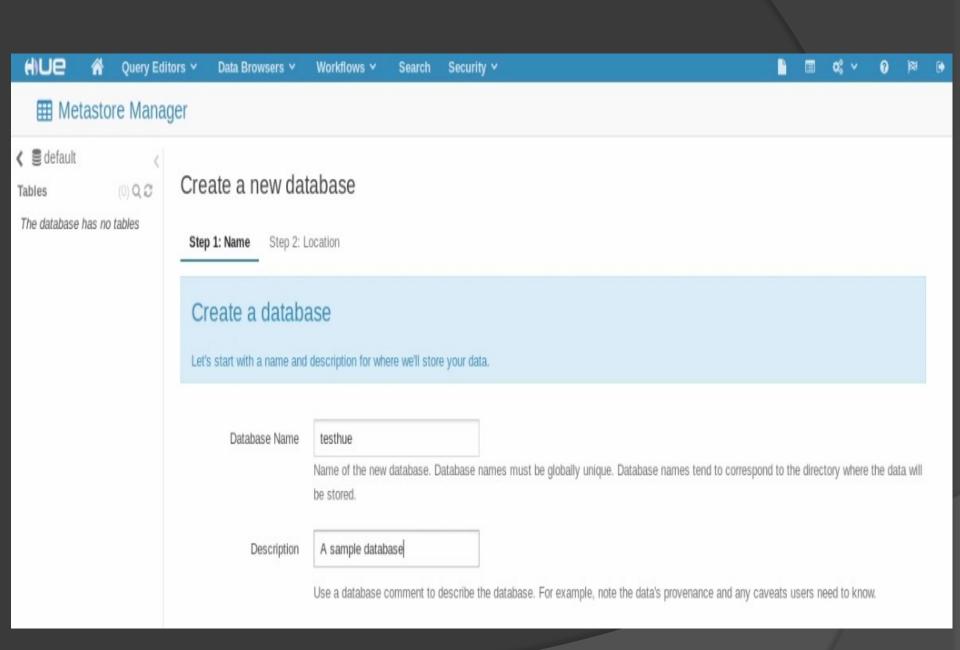
MULTI TABLE

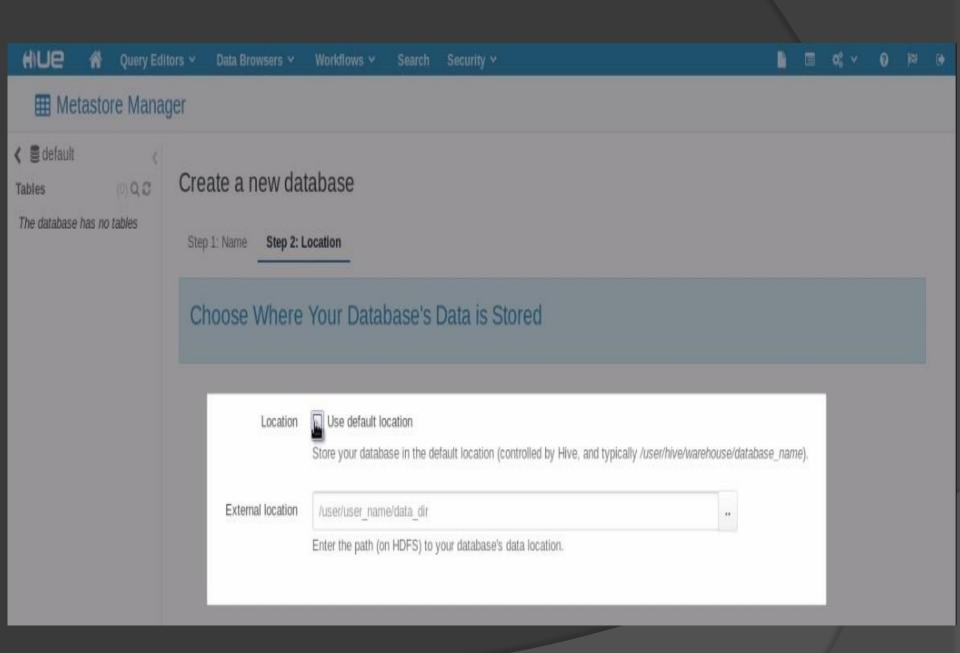
INSFRT

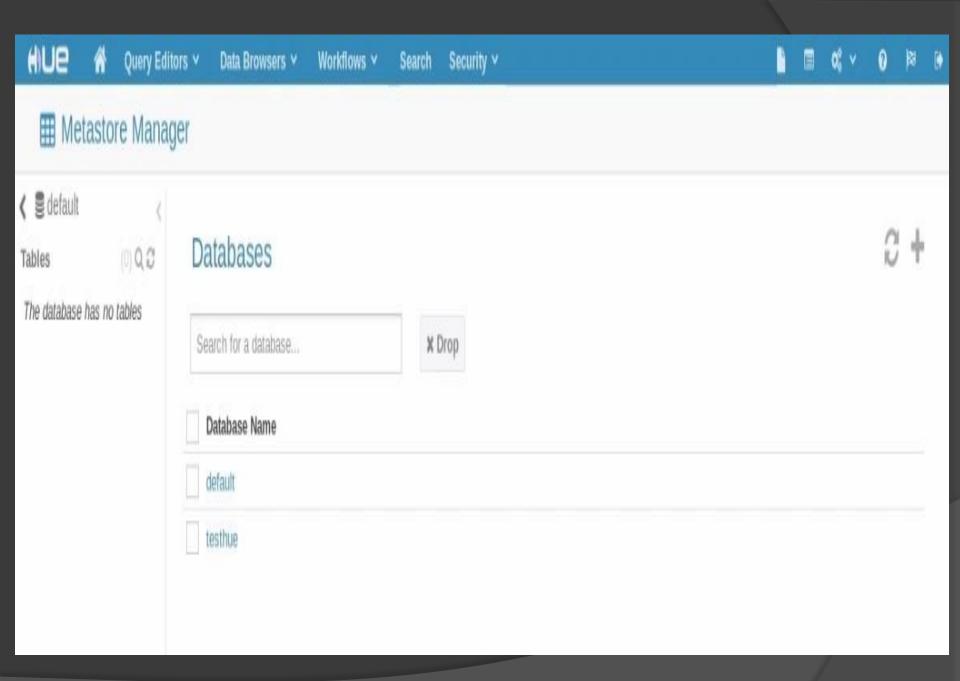
EZ ZAHOUT Aderrahmae

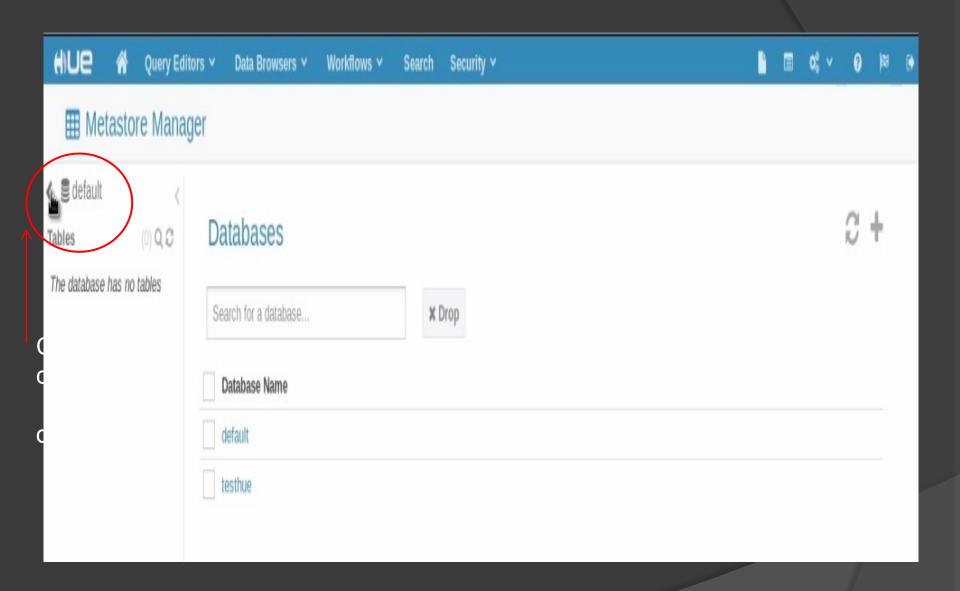


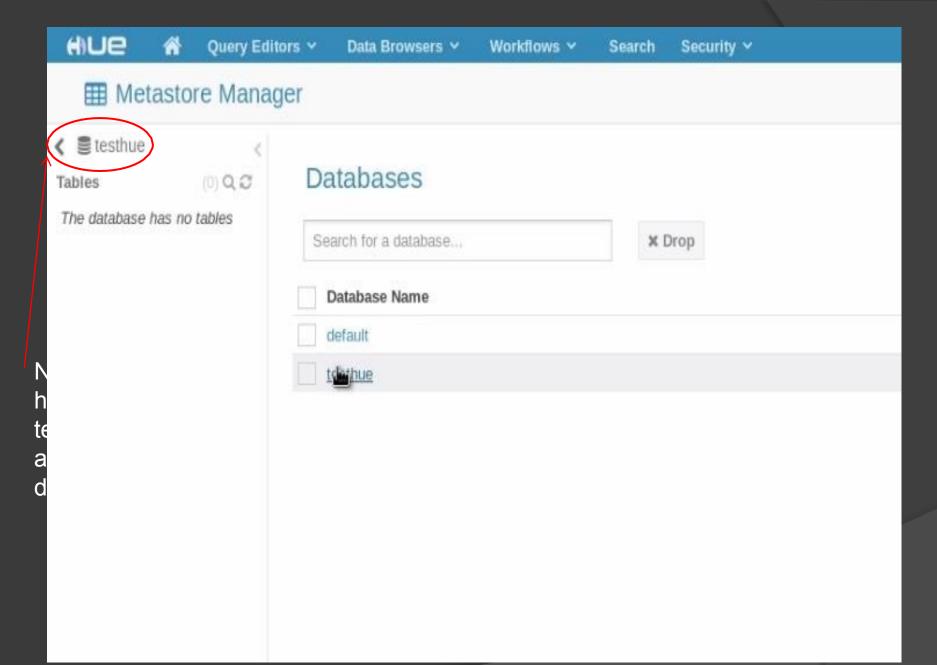




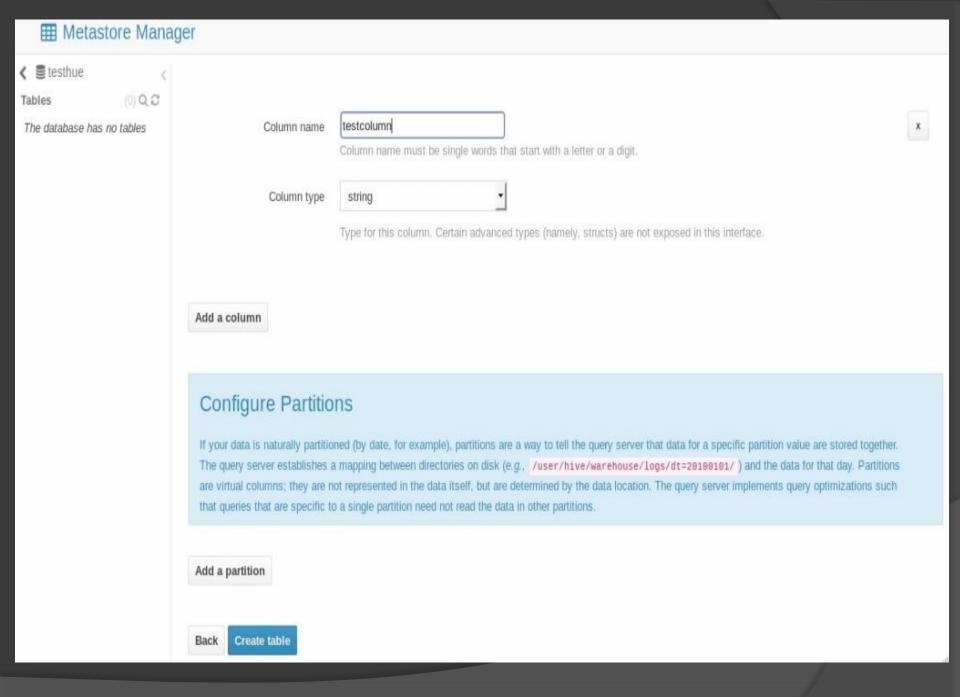








Suite : See Hive Presentation 2



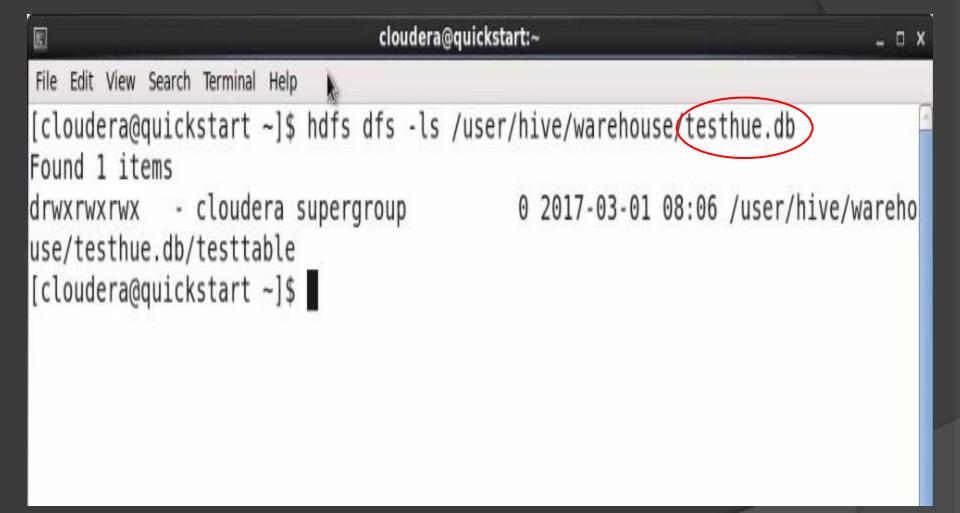
Examples

Complex operators are used to access the elements of complex type.

The following table describes complex operators available in Hive:

Operator	Operand type	Description
A[i]	A is an array object, and i is the index of an element in the array.	It will return the element at the i index of array.
M[key]	м is the map object, that is, key-value pair.	It will return the value corresponding to the specified key in the map.
S.a	S is the Struct object.	Returns the a field of S.

Suite : See Hive Presentation 2



Suite : See Hive Presentation 2