# unit-IV

- Hive – data warehousing tool
- Query structured data on Hadoop
- uses HDFS for storage & execution – Map-reduce
- Metadata is stored in RDBMS
- Data warehouse applications – suitable
- Batch processing jobs

Eg: web logs & Application logs

## History of Hive

2009 – used by Facebook – analyse incoming log data
2008 – apache Hadoop – sub project

## Recent release of Hive

| Hive 0.10 | Hive 0.13 | Hive 0.14 |
|---|---|---|
| – Batch processing | – interactive data | – transaction with ACID |
| – Read only data | – read only data | Properties |
| – Hive QL | – substantial QL | – cost based optimiser |
| | | – SQL temporary tables |

HQL (Hive Query language)
- Process Queries → Mapreduce jobs →
- provides various data types functions & format
  summarization & analysis
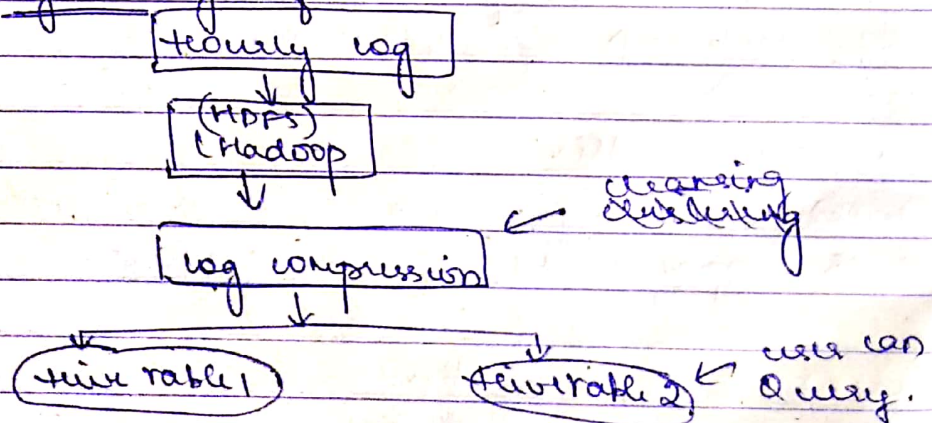                                              (continued)

## Features of Hive
i) Similar to SQL
ii) Easy to code
iii) Supports extensive data types like Structs, list
   & maps for efficient processing
iv) Supports some filters and of SQL, group by,
   orderby process calcles
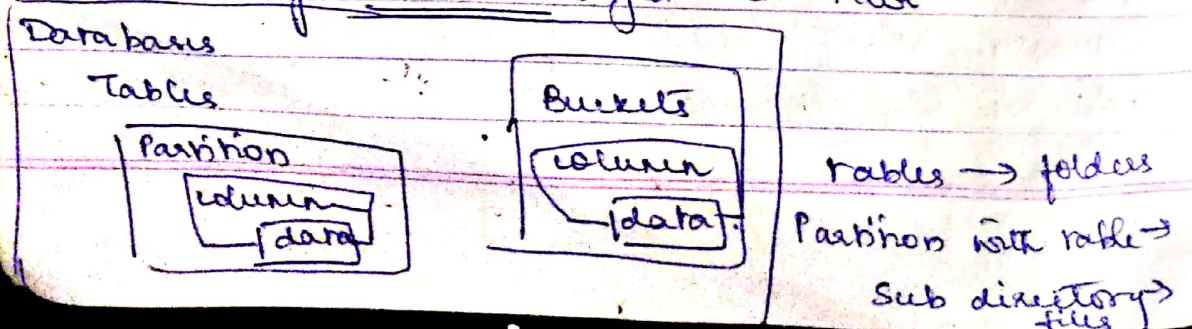v) few costume types & functions are supported


- Hive Data units
i) Databases: Namespace for different tables.
ii) Tables:, Set of records having similar scheme
iii) Partition: logical separation of i/p data-specific att
        & stored in form of folders
iv) Bucket (cluster): uses hash function for segregating
    data of i/p & decides which record goes to
    which bucket / cluster

emp
Flow of log analysis file
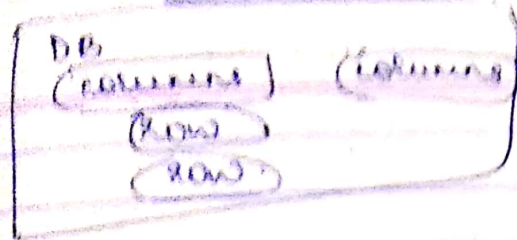
```
        ┌─────────────┐
        │ Hourly log  │
        └─────────────┘
               │
               ▼
        ┌─────────────┐
        │  (HDFS)     │
        │ (Hadoop)    │
        └─────────────┘
               │
               ▼                          cleansing
        ┌─────────────┐              ←    clustering
        │ log compression│
        └─────────────┘
         │              │
         ▼              ▼              user can
    ( Hive Table 1 )   ( Hive Table 2 ) ← Query.
```

Data units arranged in hive

```
Databases
 ┌──────────────────────────────────────┐
 │ Tables              ┌──────────────┐  │
 │  ┌──────────────┐   │  Buckets     │  │
 │  │ Partition    │   │  ┌────────┐  │  │
 │  │  ┌────────┐  │   │  │ column │  │  │
 │  │  │ column │  │   │  │  data  │  │  │
 │  │  │  data  │  │   │  └────────┘  │  │
 │  │  └────────┘  │   └──────────────┘  │
 │  └──────────────┘                     │
 └──────────────────────────────────────┘
```

tables → folders
Partition with table →
Sub directory →
files

semblance of hive with DB

```
DB
(columns)  (columns)
  (row)
  (row)
```

If directory
→ type and query/partition
→ Hdl file

## Architecture of Hive (7-8 marks)

```
Hive
[Hive command line          ]    [Hive web interface]
        interface

[Hive server]  [Driver (Query compile, Executing)]

        [Metastore]
    job tracker
        Task tracker    [HDFS]
    Hadoop
```

1) HCL) - interface used to interact with hive directly.
2) Hive web interface - used to communicate with hive & execute Query.
3) Hive server - optional server, to i/p jobs from remote client
4) JDBC/ODBC - jobs can be submitted from JDBC client we need to write java programs to connect
5) Hive driver - Query is sent to driver for orientation & execution
6) Hive metastore - defn of tables, mapping to the tables.
   — metastore service — interface to hive
   — database — mapping of data & file defn of database

Content of Metadata

i) ID of database
ii) ID of tables
iii) ID of indexes
iv) Time of creation of table
v) Input format used by table
vi) output format used by table

} update when created table & modify
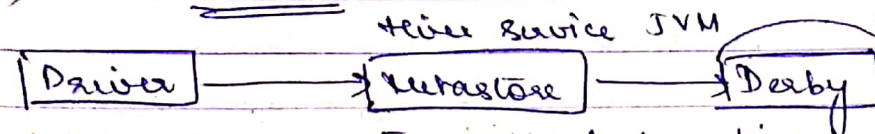
9MP

What are metastores supported by Hive?

- There are 3 meta stores
i) Embedded Metastore
ii) local "
iii) Remote "

i) Embedded Metastore :
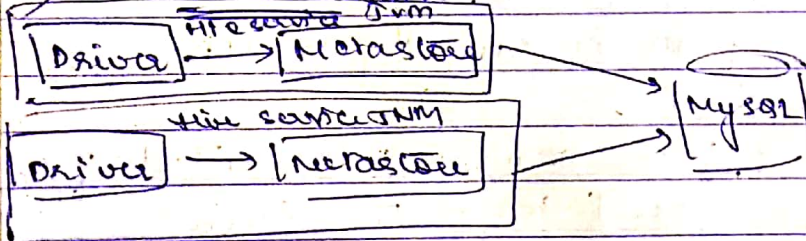
Hive Service JVM

[Driver] ⟶ [Metastore] ⟶ [Derby]

→ default metastore used by hive
→ used for performing unit tests & only one process can connect to metastore at a time
→ DB & metastore embedded in system

ii) local Metastore :

Hive service JVM
[Driver] ⟶ [Metastore]
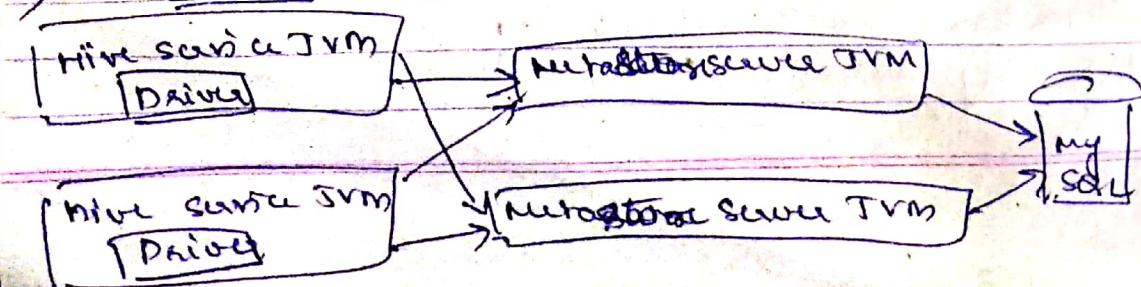
Hive service JVM
[Driver] ⟶ [Metastore]

→ [MySQL]

→ Metadata can be stored in RDBMS / MySQL
→ Multiple connection is allowed
→ metastore executed separately in server process, myHQL data DB runs in separate process.

iii) Remote metadata store

Hive service JVM
[Driver]

Metastore service JVM

→ [My SQL]

Hive service JVM
[Driver]

Metastore service JVM

→ MySQL DB is fetched from user
→ credentials of DB are completely isolated from user
→ Driver & metastore have separate JVM & therefore executed in different machines

Q.2
## Data types of Hive

| Primitive | collection | ③ Miscellaneous |
|---|---|---|
| ① Data type | Data Type | boolean |
| i) TINYINT (1 byte) | i) struct- | binary |
| ii) SMALL INT (2 byte) | 'c' structure, accessd using (.) | |
| iii) INT (4 byte) | ii) map - sequence of k value pair | |
| iv) BIG INT (8 byte) | accessed using [ ] | |
| v) FLOAT (4 byte) - Single precision | iii) array - similar data types | |
| vi) DOUBLE (8 byte). double "FP" | -accessed using index | |

② String types
  i) STRING    'abc' or "abc"
  ii) VAR CHAR     } available in hive
  iii) CHAR        } version of 0.120 & 0.130

## File formats in Hive

① text file        ② Rc file        ③ sequential file

their records are encoded in file
① → default file, separate delimiters (Ctrl+A) → separate different fields, (Ctrl+B) → separate array, struct elements, (Ctrl+C) - separate key value pair
② → content in key value pair, allows compression CPU time, I/O operations

Ans
② record column file - data stored in columns, supports aggregation Eg! c1 c2 c3 c4     row+ group

|  |  |  |  |
|---|---|---|---|
| 1 | 2 | 3 | 4 |

row wise distribution   5 6 7 8
9 10 11 12
13 14 15 16

row1 group
| c1 | c2 | c3 | c4 |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |

row 2 group
| c1 | c2 | c3 | c4 |
|---|---|---|---|
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

but in RC file format
it is distributed on column wise

| row group 1 | | row group 3 | |
|---|---|---|---|
| 1 &1 5 &1 ; | | 9 ; | 13 ; |
| &2 &2 6 ; | | 10 , | 14 ; |
| 3 , 7 ; | | 11 , | 15 ; |
| 4 , 8, ; | | 12 , | 16 ; |

HQL - features
→ creates & manages table, partition
→ relational, arithmetic, logical etc operations
are available
→ Evaluate functions
→ Download contents of Queries, tables etc to
a directory.

DDL - create table & alter parameters
updation supported - creating & managing DB
                    - drop or Alter the table/terminate
- alter |partition| column of a table
- create| prop| view
- create | prop |index
- show - contents
- describe - contents & other description

DML - retrive the data
      - alter | modify | delete | update
- O|P can be loaded to the table.

# Hive

--. CREATE DATABASE IF NOT EXISTS STUDENTS
comments 'STUDENTS details' with DBproperties ('C
                              'create' = 'abc' )

hive> show DATABASES

hive> DESCRIBE DATABASE STUDENTS
                              EXTENDED — properties of DB along
                                        with other contents
hive> ALTER DATABASE STUDENTS SET PROPERTY ('editdb'
                                            = abc')
> USE DATABASE    > DROP DATABASE STUDENTS

## Tables

**i) Manageable warehouse**
- under hive dictionary
- entire life cycle is managed by hive
- dropping internal table remove data along with metadata

**Creating Manageable table**

```
CREATE TABLE IF NOT EXISTS
STUDENT ( ROLL INT, name STRING,
id INT) ROW FORMAT DELI-
MITED BY 'lt'
```

- HIVE> DESCRIBE STUDENT
  - roll INT

*along with table data - metadata is also removed

**ii) External table (self manag d table)**
- table is dropped, data is stored in underlying location
- create → EXTERNAL - key word used.
- specify the location where data is to be stored, if table is dropped.

**CREATE EXTERNAL TABLE**

```
STUDENT (roll, INT, name STRING,
id INT) ROW FORMAT DELITE
MITED FIELDS TERMINATED
BY 'lt' LOCATION='/DOC/Sh'
;
```

---

## Loading data from file into table:

```
LOAD DATA (LOCAL) INPATH '/student/docs/aati'
OVERWRITE INTO TABLE STUDENT
```

→ loading data from local file system

loading from HDFS then no LOCAL keyword.

---

## Collection Data types

```
CREATE TABLE STUDENTINFO (ROLL INT, name STRING,
sub ARRAY<STRING>, marks MAP<STRING, INT>)
ROW FORMAT DELIMITED BY 'lt';
COLLECTION ITEMS DELIMITED BY ':';
MAP KEYS DELIMITED BY '!';
LOAD DATA LOCAL INPATH (/student/filename) INTO
TABLE STUDENT INFO;
```

---

## Querying table

```
SELECT * FROM STUDENT;
SELECT rollno, name, FROM STUDENT;
SELECT name, MARKS['marks'] FROM STUDENT INFO;
```

Partitions — 2 Type
i) static                         ii) dynamic
- different columns, extracted    - only during execution time
  during compile time               when values are known
CREATE TABLE STUDENT             - CREATE TABLE IF NOT EXISTS
  IF NOT EXISTS
STUDENT (round INT, name          DYNAMIC_PARTITION_STUDENT,
STRING) PARTITIONED BY            (round INT, id INT) PARTI-
(roll INT) ROW FORMAT             TIONED BY (gpa >'4') ROW FORM
FIELDS DELIMITED TERMI-           DELIMITED FIELDS TERMINA,
NATED '\t';                       '\t'


loading data into partitions from table,
INSERT OVERWRITE STATIC_PART_TABLE PARTITION
  (roll INT = '10');
SELECT name, rollno. FROM EXT_STUDENT WHERE
  rollno = )

ALTER TABLE STAT_PARTITION_STUDENT ADD PARTITION(
                                                  gpa=4
INSERT . . . . .
SELECT . . . . .


Bucketing — if table size is too large, No. of partitions
  made
- to limit no. of partitions — bucketing
- partitions — directory    bucket — file
- set hive.enforce.bucketing = true    # bucket is enabled
CREATE TABLE IF NOT EXISTS BUCKET_STUDENT (
roll INT, name STRING, id INT)
CLUSTERED BY id into 3 buckets    # create 3 bucket
|| load data into bucket table ||
FROM STUDENT
INSERT OVERWRITE TABLE BUCKET_STUDENT
SELECT rollno, name, id
SELECT DISTINCT GRADE FROM BUCKET_STUDENT
TABLESAMPLE (BUCKET 1 OUT OF 3 on ID);

# VIEWS

CREATE VIEW STUDENTVIEW { AS SELECT rollNO, id

FROM STUDENTS;

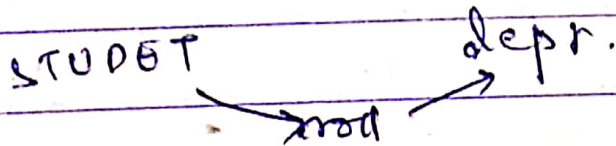SELECT rollno., id;

FROM STUDENTVIEW

LIMIT 4;           || limit to 4 rows

DROP VIEW STUDENTVIEW ;

## Nested

- o/p of one Query is i/p to one query
- uses clauses like GROUP BY, ORED BY etc

## Join

STUDET → roll → dept.

select name, roll, id FROM student a JOIN DEPT d ON

a.roll = d.roll;

## Aggregation Func

count(), avg() etc.

## Group by ... Having

slect roll, name from STUDENT having id > 20

group by roll;