
2018



SerDe (Serialization & Deserialization in Hive)

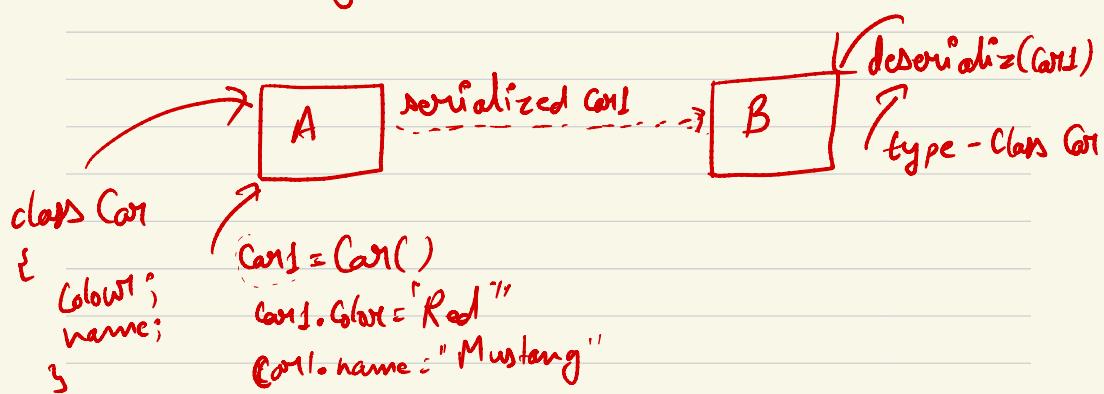
Let's break the hive query & creation.

- Serialization & Deserialization
- Hive ROW FORMAT
- Map-Reduce Input/Output

* Serialization: It is a process of converting an object into bytes that can be stored in a file or transmitted over the network

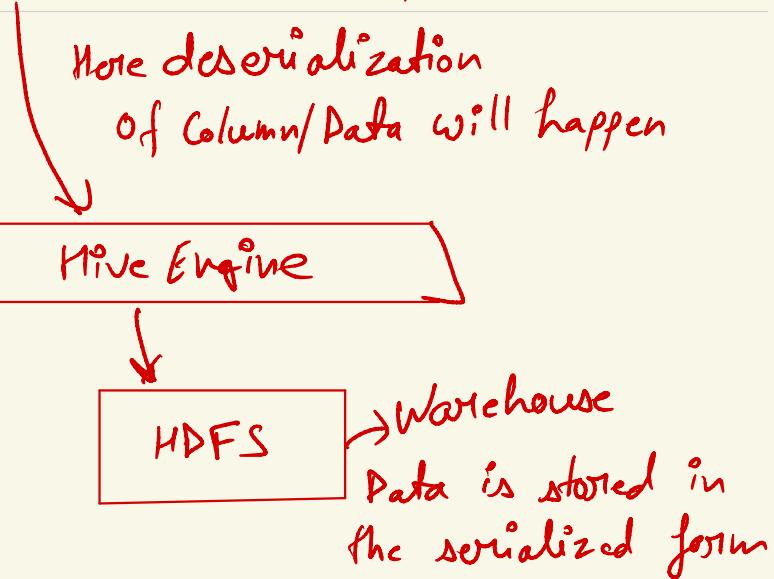
* Deserialization: Process of converting bytes into an object.

here



Hive Query

Select {name},{salary} from employee;



When we run the load command that means, we are loading serialized data in hdfs.

= > Create table mytable (a string, b string)
> **Row Format** delimited
> **Stored as** text file;

Row Format: describes the libraries used to convert a given row into column object.

Stored As: describes the input format & output format libraries used by map-reduce to read & write to HDFS.

i) Columnar file format / Databases

→ Parquet

→ ORC (Optimized Row Column)

Row Based

employee-table (name, salary)

100	105	108	1913	118
ABC	1000	DEF	2800	KLM

Read Pointer

$$1 \text{ int} = 4 \text{ Bytes}$$

alist = [1, 5, 9, 10]

100	104	108	112
1	5	9	10

$$\begin{aligned} \text{alist}[2] &= \text{a_list} + 2 \times 4 \\ &= 100 + 2 \times 4 \\ &= 100 + 8 = 108 \end{aligned}$$

Select sum(Salary) from employee;

in Row Based \rightarrow Seek Time is super high

Column Based

name	ABC	EFG	KLM
------	-----	-----	-----

Salary	1000	2000	9000
--------	------	------	------

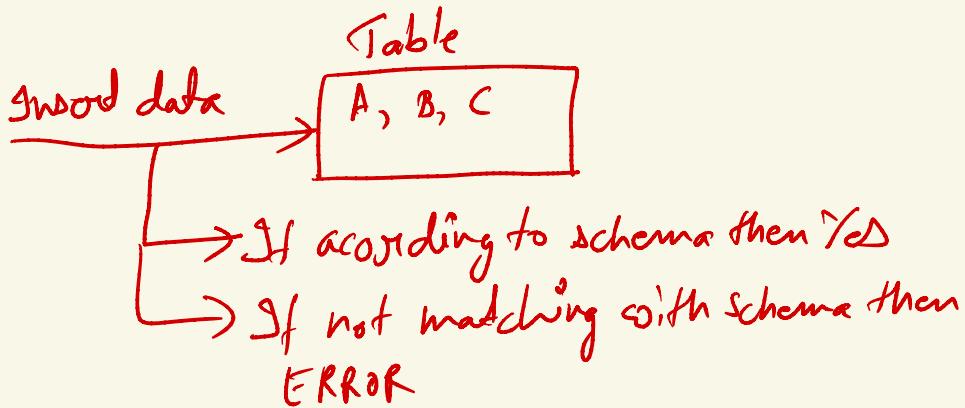
Select sum(Salary) from Employee;

ORC Comprison \rightarrow Parquet Comprison

Hive works well with ORC file format.

Spark performs well with parquet file format.

Schema on Write (RDBMS)



- x) Validation will happen during the write time or declaration time.

Schema on Read (Hive)

