

---

---

---

---

---

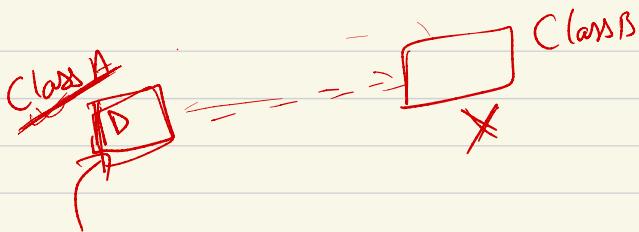


Serialization → Converting Objects into bytes  
Deserialization → Converting bytes into objects

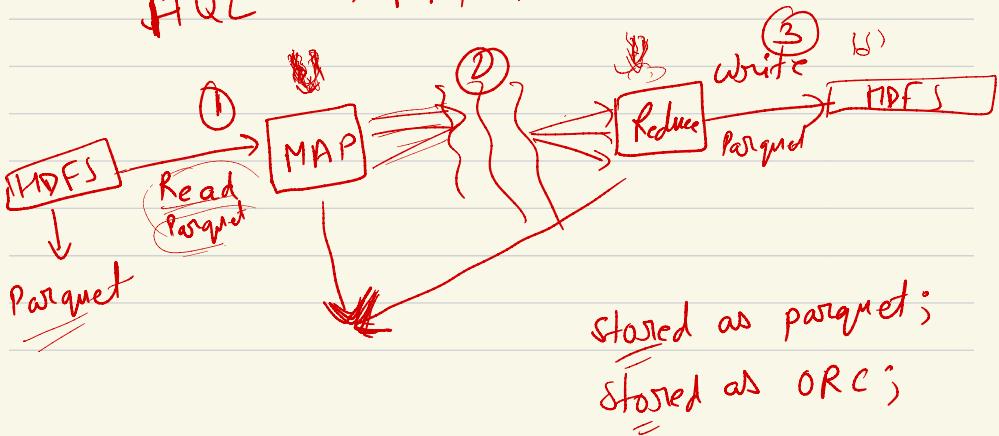
Input format → Parquet SerDe

Output format → Parquet SerDE

SerDe library → Parquet Hive SerDe



FQL → MAP-Reduce



Row format delimited ; .csv

Row format sende "OpenCSVReader";  
    {      }  
    {  
        }

.csv  
↓

quoted

Nested delimiter

Row

select \* from employee

name, salary

Serde library [ Parse the raw data ]

Deserialized object (Row)

" "

" Shashank "

" , "

" Shashank Mishra "

QuoteChar = " = "

QuoteChar = " " "

" Shashank, Mishra "

escapeChar = " \\" "

" \Shashank Shashank, Mishra \" DE "

" Shashank, Mishra " DE " "

Raw data → " Shashank, Mishra \" DE \" "

Output      ↗ shashank, Mishra " DE "

a = "Shashank"

a = "Shashank"

"Shashank, Mishra"

name, location  
shash, "Bangalore, India"

"10"

15

quartz

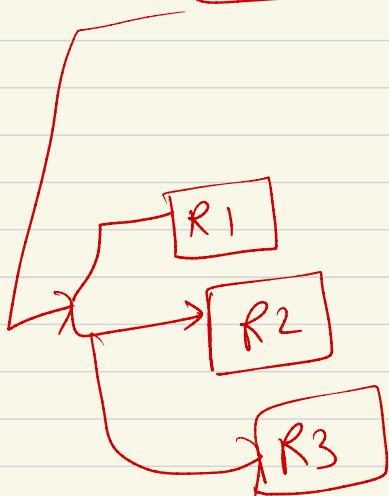
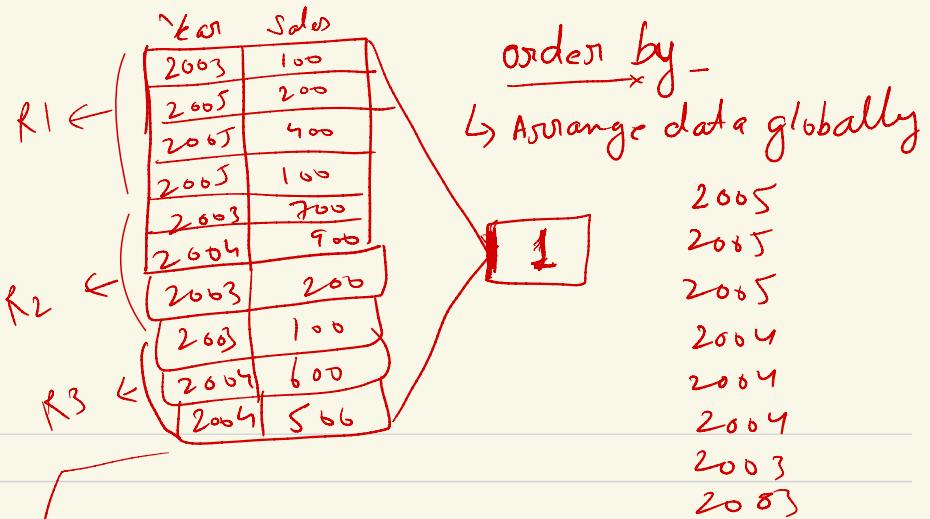
X "Shashank"

Sales	Year-ID	Qtr
200	2003	1
500	2003	1
200	2004	2
100	2004	3
300	2005	4
200	2005	5

select year\_id, sum(sales) as total\_sales

from sales\_data group by YearID

{ 2003, 900  
2004, 800  
2005, 500 }



Sort by  
 ↪ local sorting

2005  
 2004  
 2003

2005  
 2004  
 2003

2005  
 2004  
 2003

→ R1

→ R2

→ R3

## Partitioning -

R1 → Shaank, IT  
R2 → Rahul, QA  
R3 → Amit, IT

1M → Record

## Employee

	<u>Name</u>	<u>Salary</u>	<u>Country</u>
	Shaank	1	India
	Amit	1600	UK
	Amit	500	India
	Rahul	400	China
	Surbhi	800	USA
	Kapil	900	USA
	Mulk	1000	Russia

→ which column should be used for partitioning?

→ Columns used to filter the data frequently

→ Always avoid Primary key type of columns for partitioning  
↳ 1M (emp\_id)

1M (Super Bad)

 loss Uniqueness (Repeated multiple  
time or used)

any hive table with partition

## nominal Interval hierarchical table

L) /user/hive/warehouse/db-name/tb-name/

→ /user/hive/Warehouse/db-name/tb-name/Country=IND/

- Laundry - USA //

- Country = China

+ (and say = Russia)

/ Candy = UK )

**Static Partitioning** → We know in which partition data will get loaded

 **Dynamic Partitioning** → Hive will scan the partition column and will create partition inhds for each distinct value

## Q) Dynamic & Static Partitioning

which one perform slow?

create dynamic table

{

=

↳ Partitioned by (Country, Year-id)

half-path — / Country = USA / Year-id = 2005  
/ Country = USA / Year-id = 2004

