


~~- CSV, JSON~~

- 1) Parquet
- 2) ORC
- 3) AVRO

CSV, JSON
[grid] { "name": Val, "col": Val2 }

"name" → string

→ L → int

→ f.d → float, double

machine can understand 0 0 1

'A' : 'B' , 'C' 1Byte = 8 bits

Char = 1 Byte for storage

int = 4 Bytes for storage

float = 4, 8, 16 bytes for storage



1000 Row → 10000

Storage for one single Row

$$= 10 + 8 = 18 \text{ bytes}$$

$$1000 \times 18 = 18000 \text{ bytes}$$

T1

name	title	salary
Shubh	DE	1000
Rahul	SDE	2000
Ankit	QA	5000

int $\underline{\underline{a}} = 4$

$\underline{\underline{x1}}$ 1000UA

$$100 + 8 = \underline{\underline{108}}$$

$$100 + 108 = \underline{\underline{118}}$$

$$118 + 8 = \underline{\underline{118}}$$

$$118 + 5 = \underline{\underline{123}}$$

$$123 + 3 = \underline{\underline{126}}$$

$$126 + 8 = \underline{\underline{134}}$$

table in MySQL DB

or any RDBMS

Row-Column format

Shashank = 8 char

XV66	108	110	118
Shashank	DE	1000	

Shashank, DE, 1000
S 100 ----- 117

* Rahul, SDE, 2000
S = 118 133

name, title, salary

Shashank, DE, 1000

Rahul, SDE, 2000

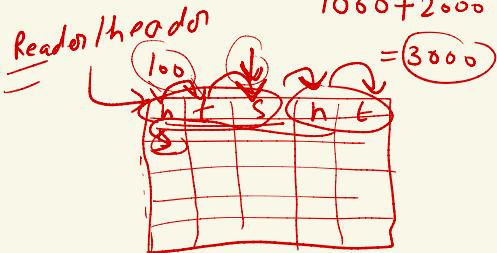
Select sum(salary) from(b);

- ⑥ Row-Columnar
⑦ Columnar Format

`int a=4;`

Print (⑦)

Disk)
Read header
Write

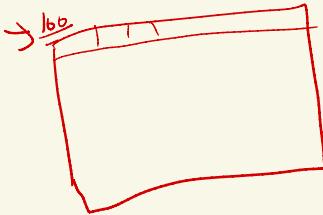


Columnar format

X

name	title	Salary
Shashank, DE	DE	1000
Rahul, SDE	SDE	2000

Columnar



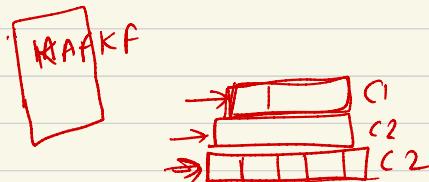
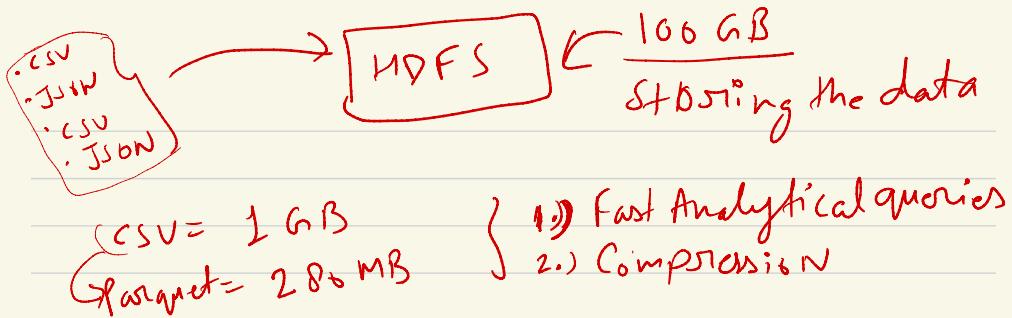
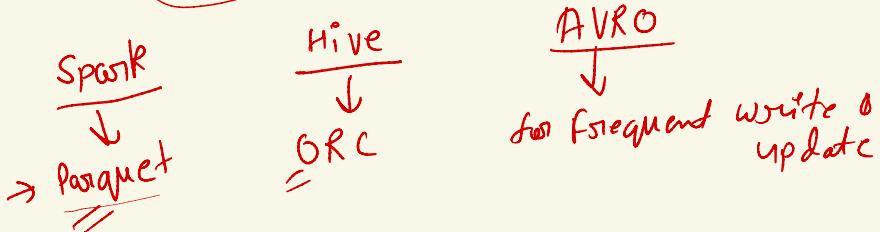
name	Shashank	Rahul	Ankuri
title	DE	SDE	QA
salary	1000	2000	3000

Disk reads → (100) (150) (200)
Select sum(Salary) from tb1;

NO SQL → Key-value, document, column

For frequent or heavy analytical queries
prefer Columnar formats or Columnar NoSQL DBs

- 1) Parquet (Columnar)
 - 2) ORC (Columnar) [Optimized Row-Columnar]
 - 3) AVRO (Row-Columnar)
- Comprehension



Comprehension

✓ ORC > Parquet > AVRO

Serilization & Deserializatior

objects

→ class, objects

↳ physical instance

```
class Car {
    String name;
    String model;
    float price;
}
```

→ class Car Car1 = new CAR();

=

Car1.name = "Hyundai"

Car1.model = "SEI"

Car1.Price = 100K

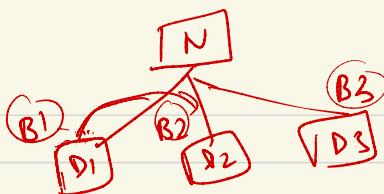
→ Car1 → memory

Car2 → ↗

Car3 → ↗

{ Car2
Car3 }

Car1.name



Data flow over Network → Serialized data

int a = 4;

4

100

Deserialize

Converting something
to its original form
or int original data types

variable
object
data

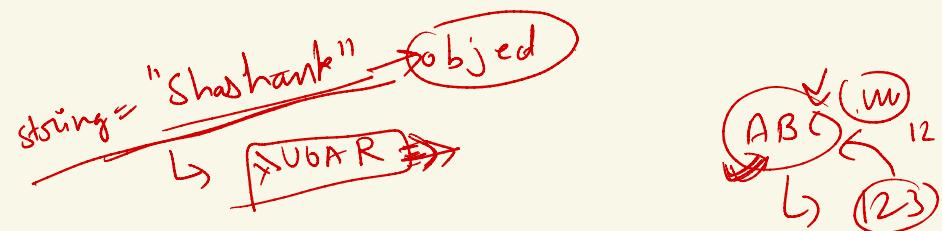
Serialized form (Converting original data
to object to object)

Data types in Byte(s)

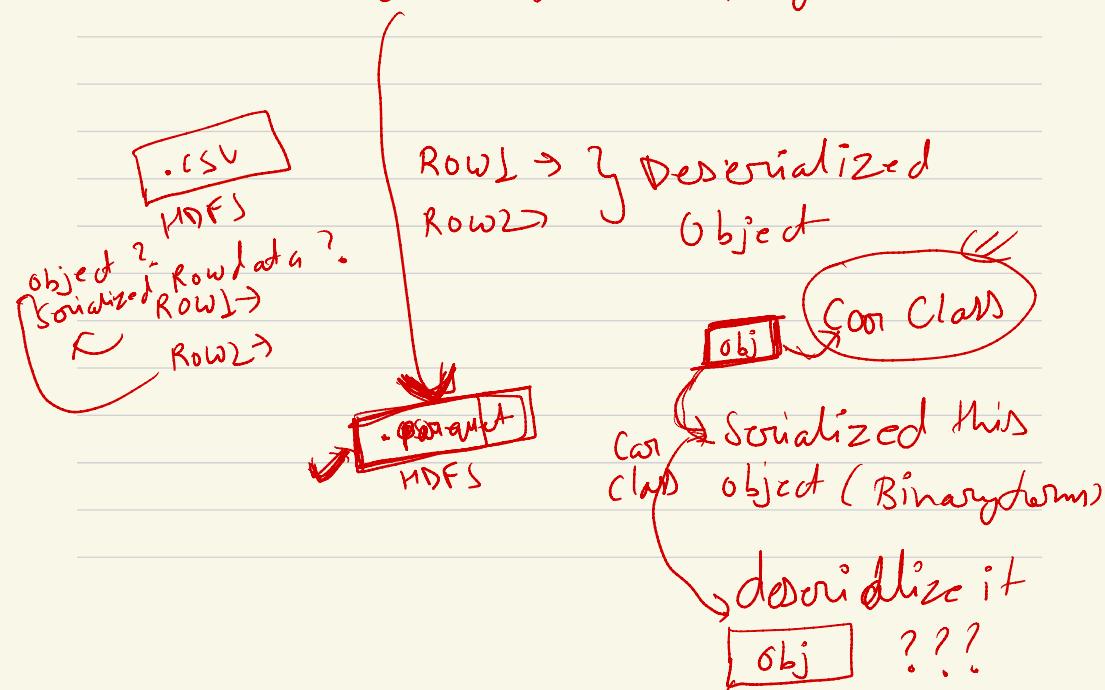
Advantage of Serialization

- Data in bytes
- Fast transmission over network
- Less data size

Select * from employee;



Select * from employee;



Class Animals

{

Car Class

???

Class Animal

= ???

String type: ...

3

Terminal

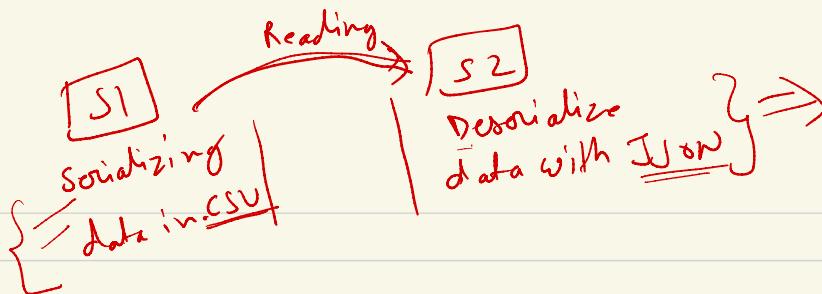
Serialized

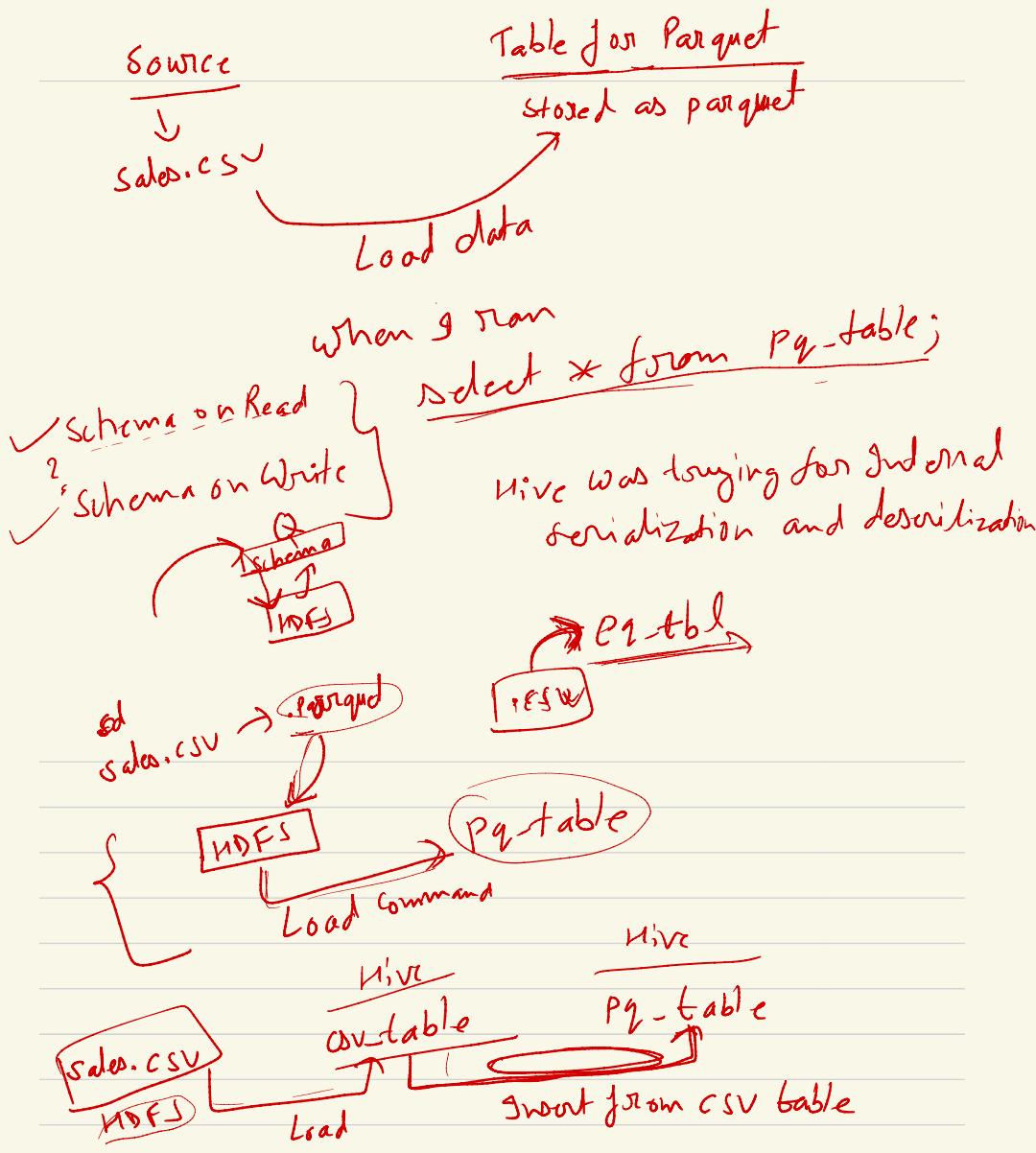
JSON

query

{ JSON }

You should have a mechanism or library which can deserialize it

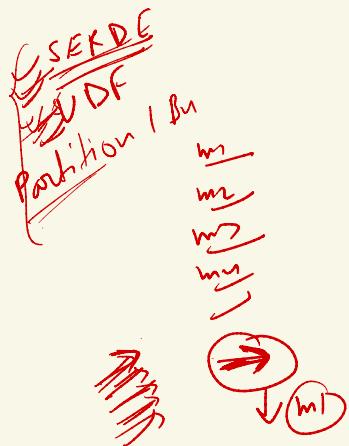




From covtable Insert overwrite Pg-table Select *;

("clothes", 1000) → object

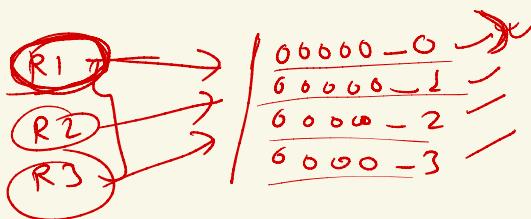
→ serialized data for parquet



→ stored as Parquet

LGB

$m_1 \leftarrow R_1$



Parquet → multi parts
ORC →