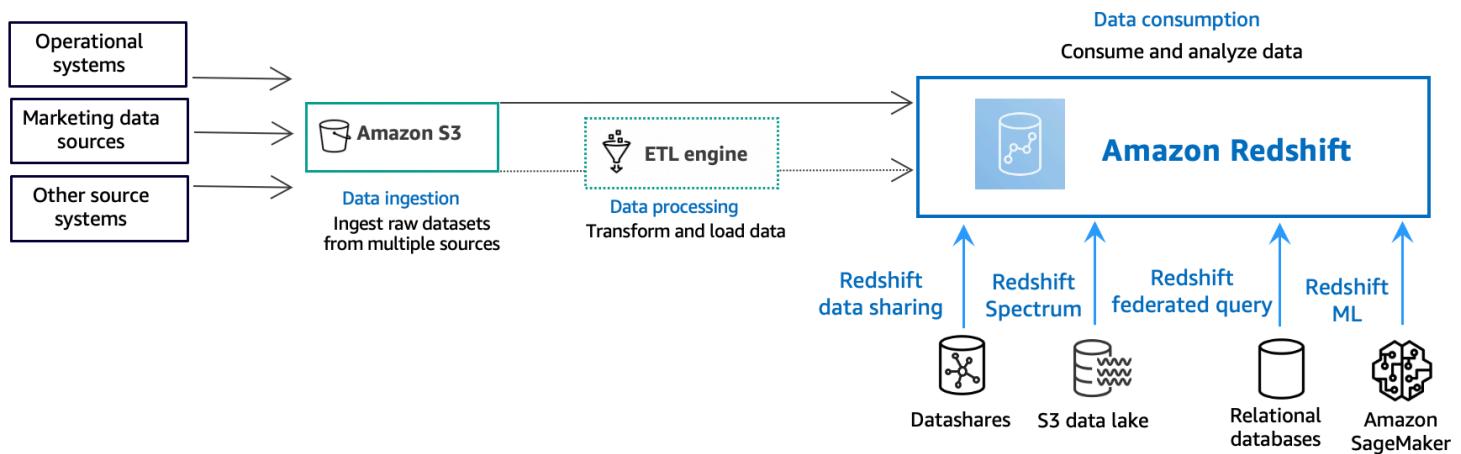
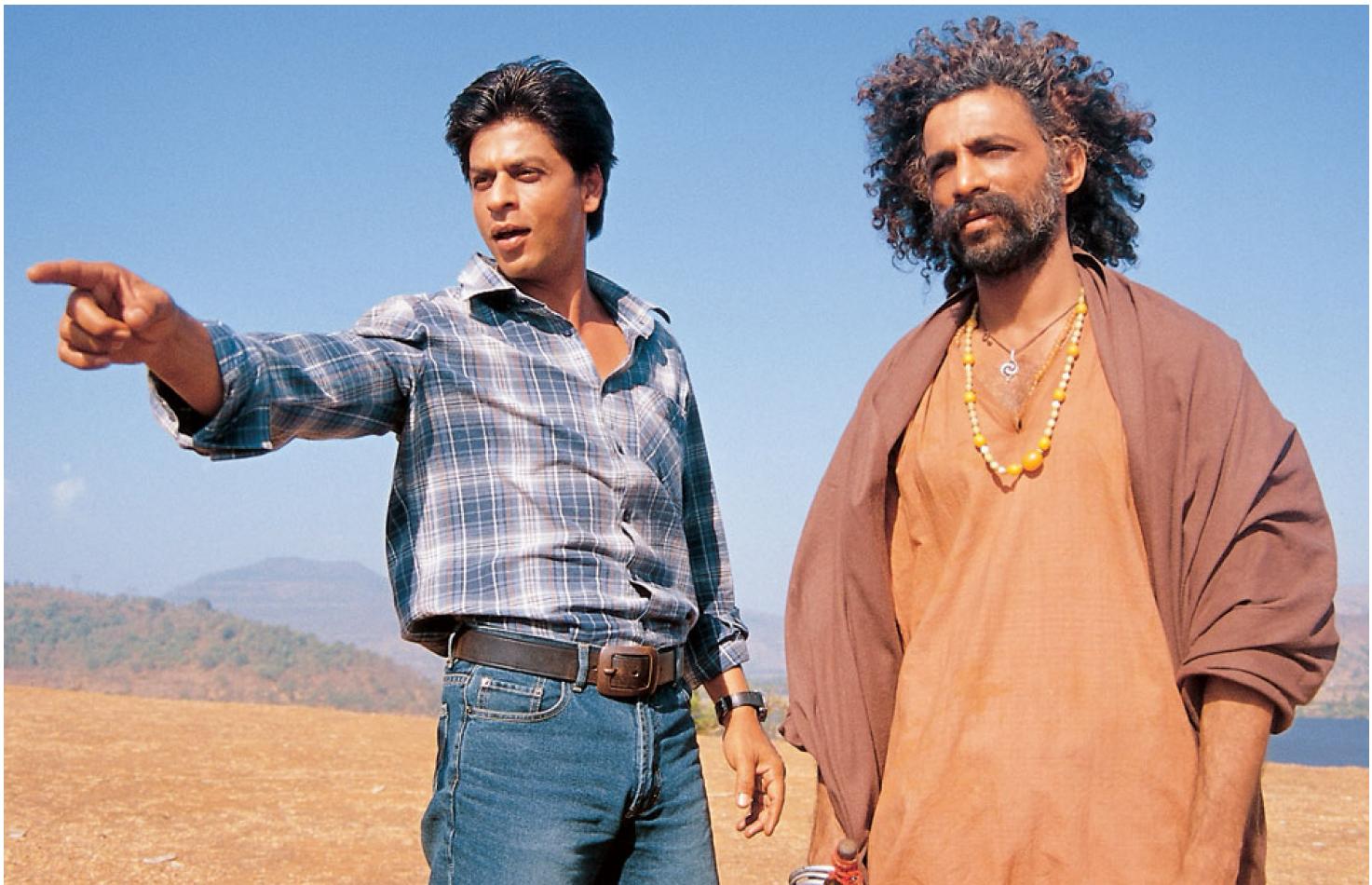
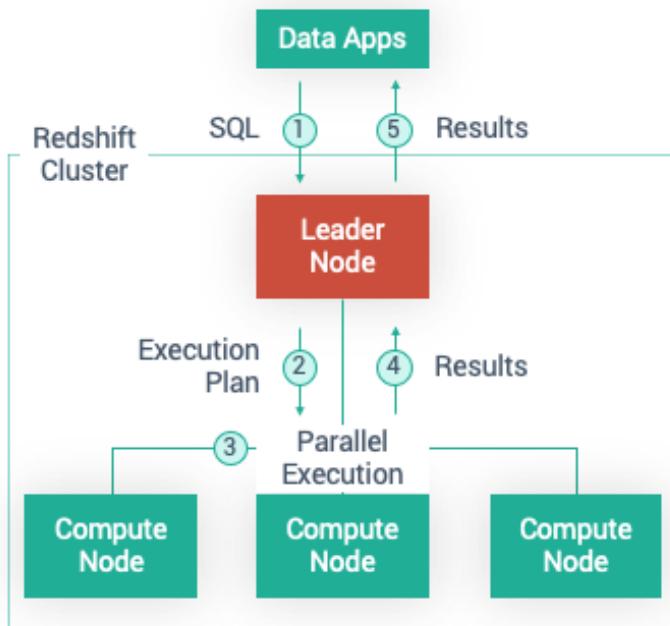


# Agenda

- a. Optimisations in Hive
- b. Essential features of Redshift
- C. Columnar Data warehouse concepts
- d. Architecture
- e. Distribution key and sort key



# Amazon Redshift Architecture: The Life of a Query



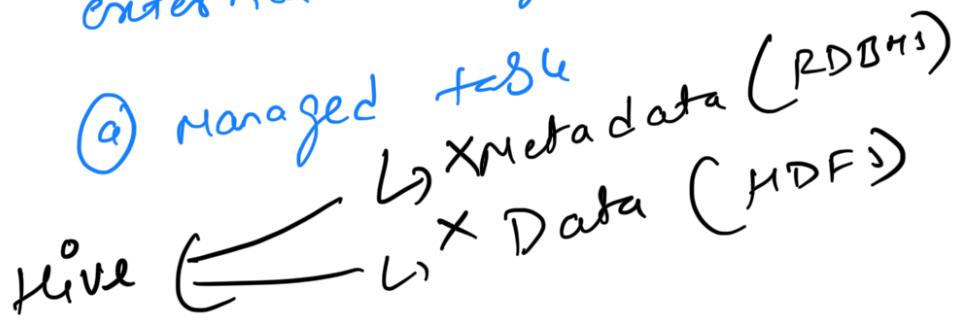
## The 5 Steps to Process a Query

- ① The cluster receives a query coming from a data app and parses the SQL in the leader node.
- ② The leader node creates an execution plan that breaks a query down into a discrete sequence of steps.
- ③ The leader node distributes the work of executing the steps in parallel across the compute nodes.
- ④ The compute nodes send the results back to the leader node to merge data into a single result.
- ⑤ The leader node addresses any final sorting or aggregation and returns the results to the data app.

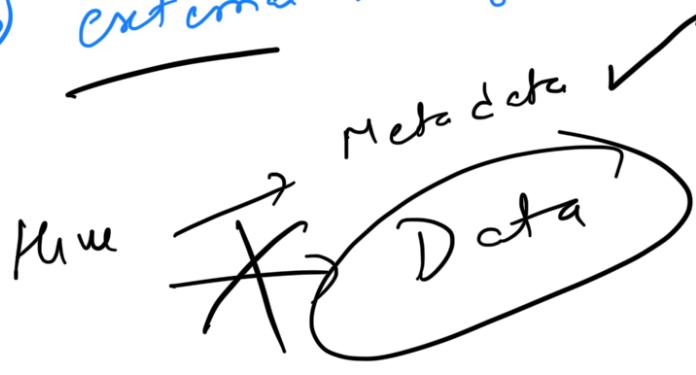
## Optimisation

① change the execution Engine  
↳ MR / Tez / Spark

② use external managed table



b) external managed table



③ File types :-

- a) Text File
  - basic human readability
  - ↳ CSV, TSV, TSV
  - Consumes more space
- b) Sequence File
  - Key - value binary format
  - More efficient
- c) Arrow
- c) Parquet

Use ORC

① Avro  $\Rightarrow$  (JSON)  
Data Serialization Framework.

Data types  $\Rightarrow$  null, boolean, int, long,  
float, double, bytes  
string

### Complex Data types

① Record (Table)  $\Rightarrow$  A UDT composed of 1 or more named fields.

② enum  $\Rightarrow$  Specific set of values

③ array  $\Rightarrow$

④ map  $\Rightarrow$

⑤ Union  $\Rightarrow$

⑥ fixed  $\Rightarrow$

Exactly one value matching a specified set of type

A fixed number of bytes

$$\therefore A = \{1, 2, 3\}$$

$$A \cup B = \{1, 2, 3, 4, 5\}$$

$$\textcircled{2} \quad A \cup B \Rightarrow : \textcircled{1} = 4, 5, 6$$

Arm Selection  
Arm De Selection

Disadvantages:-

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

- ① Arm supports now with storage will need
- ② To reach wanted scheme is needed.

## ② Parquet files

↳ Columnar storage

|   | id | name | Sal |
|---|----|------|-----|
| 1 | 1  | A    | 100 |
| 2 | 2  | B    | 200 |
| 3 | 3  | C    | 300 |
| 4 | 4  | D    | 600 |

|   |    |   |   |   |   |
|---|----|---|---|---|---|
| X | id | 1 | 2 | 3 | 4 |
|---|----|---|---|---|---|

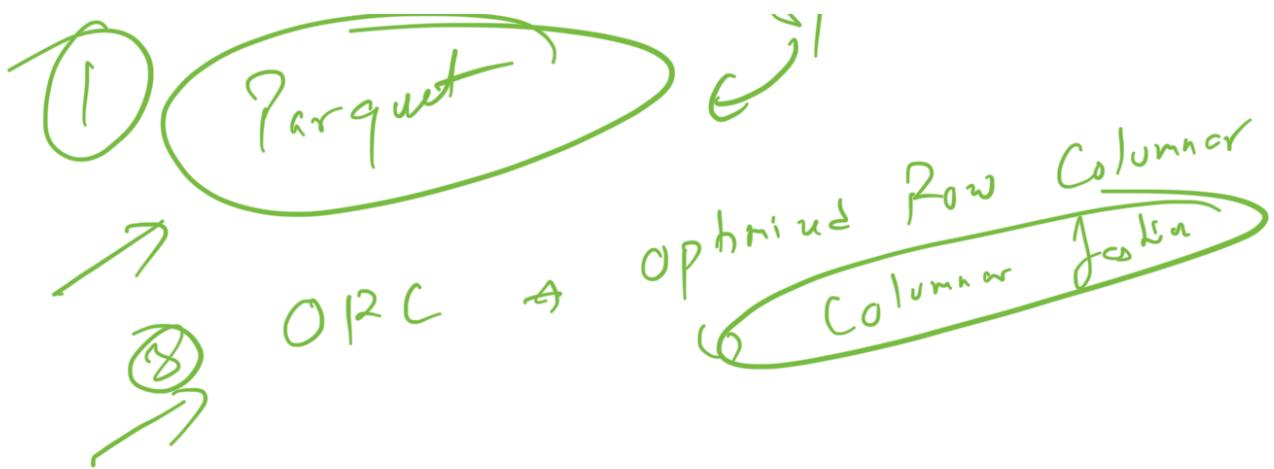
|   |      |   |   |   |   |
|---|------|---|---|---|---|
| X | name | A | B | C | D |
|---|------|---|---|---|---|

|     |     |     |     |     |
|-----|-----|-----|-----|-----|
| Sal | 100 | 200 | 300 | 600 |
|-----|-----|-----|-----|-----|

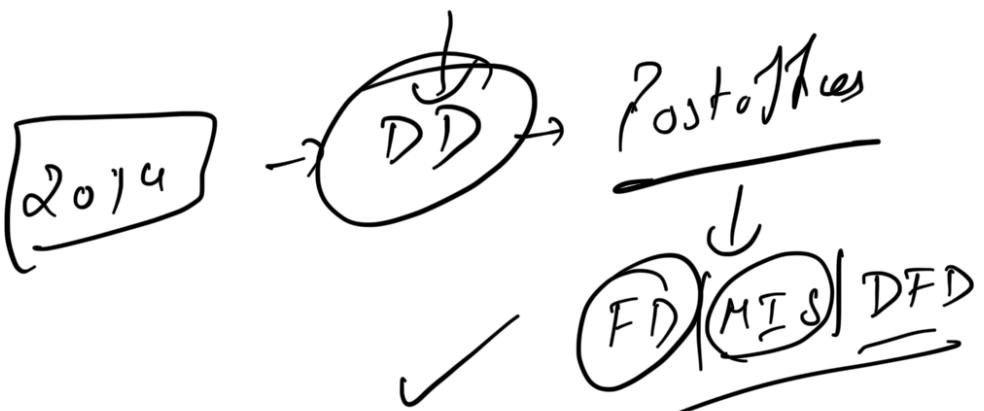
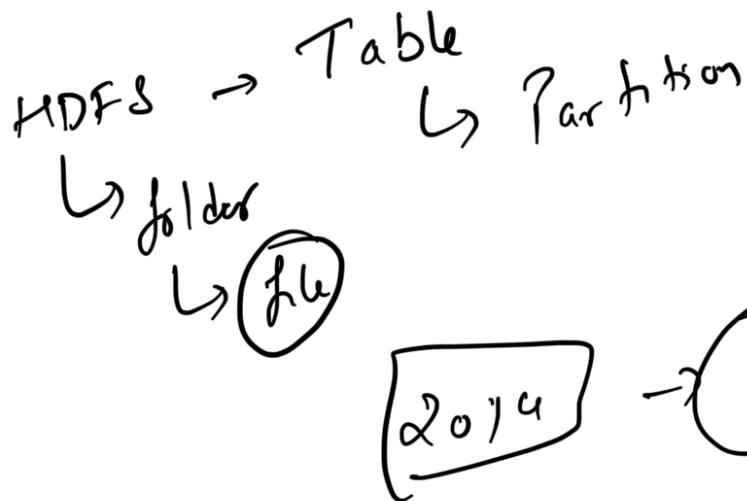
mcn (Sal)  
4 → 600

QParquet

Select \* from emp;



### ③ Partition:



✓

-,-  
;:

29