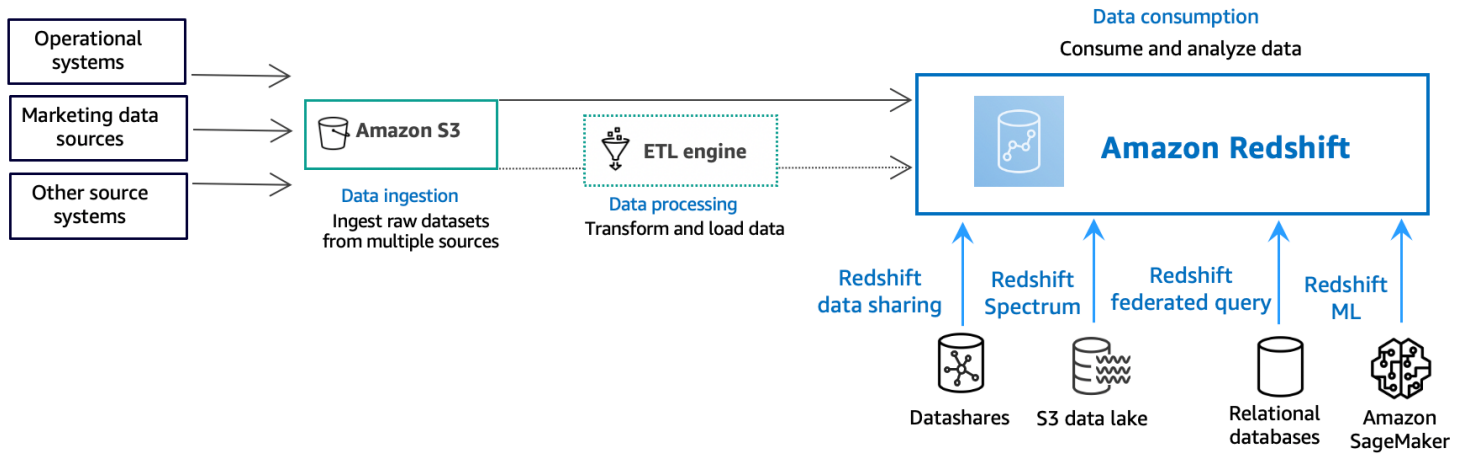


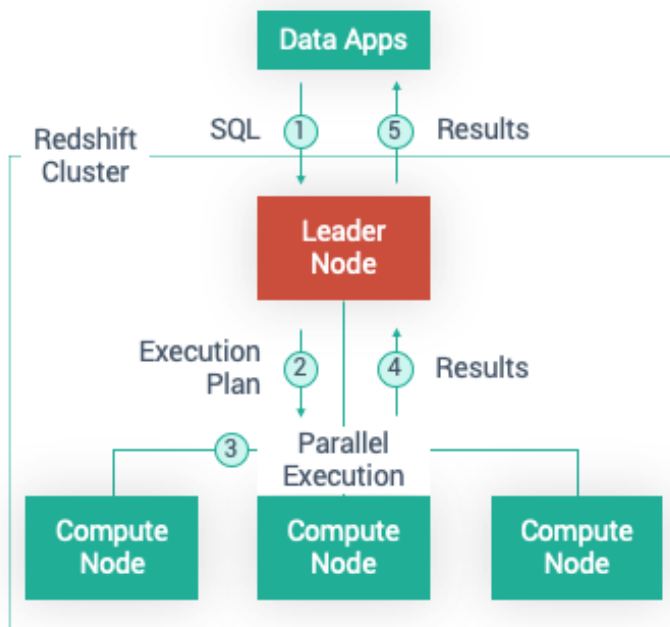
# Agenda

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



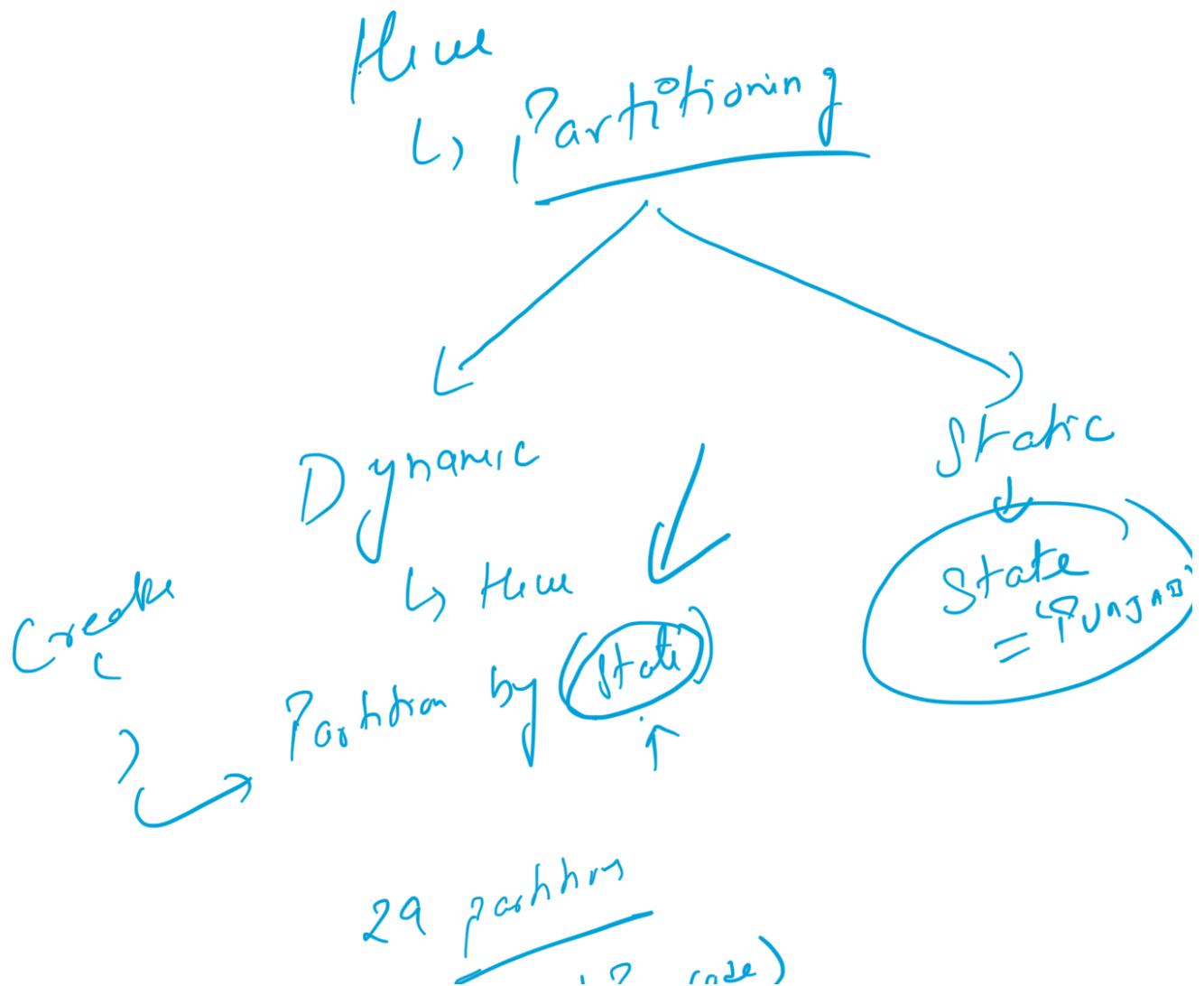
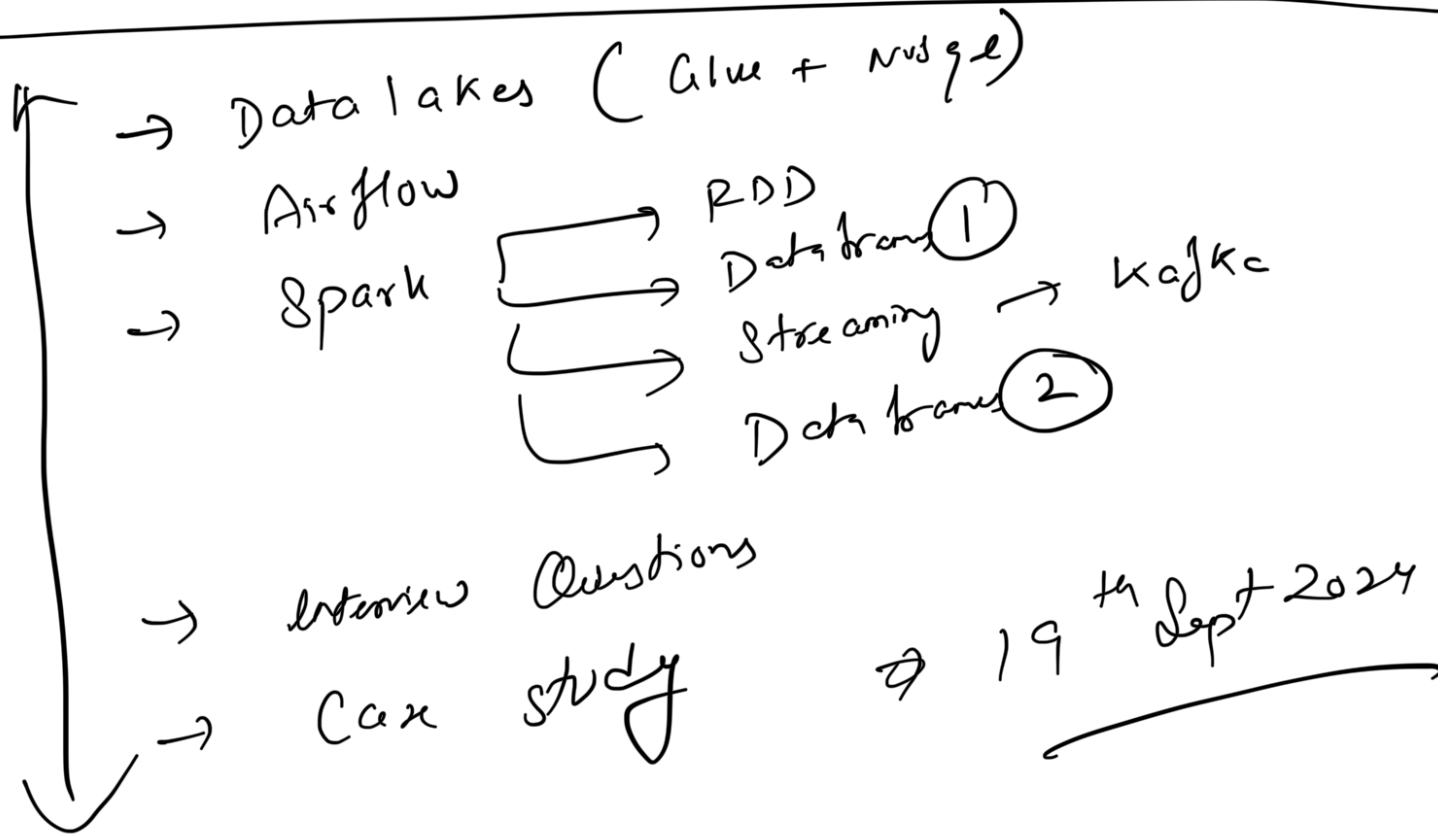


## Amazon Redshift Architecture: The Life of a Query



### 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.



↳ Partition by time

Cloud

↳ Service based Model

Tenant

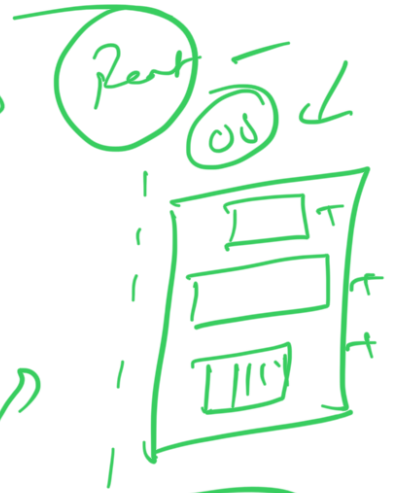
↳ Bang

fully hosted  
run

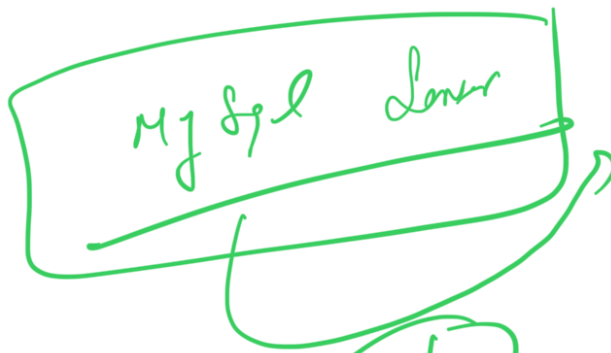


01<sup>st</sup> Sept

30<sup>th</sup> Sept



2 weeks



N/W

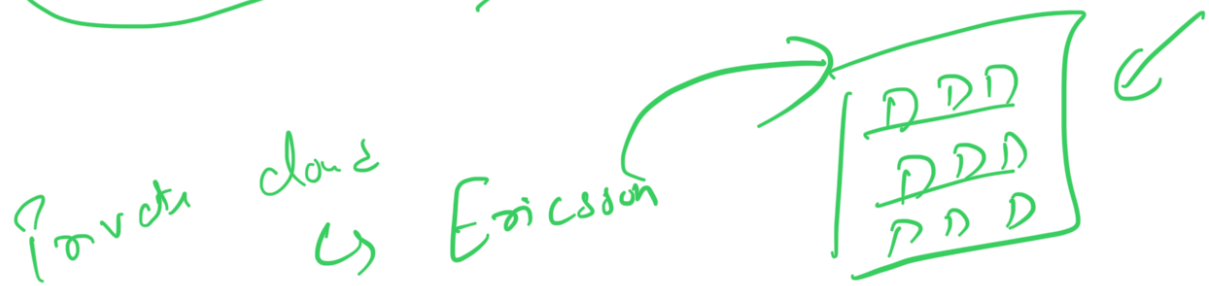
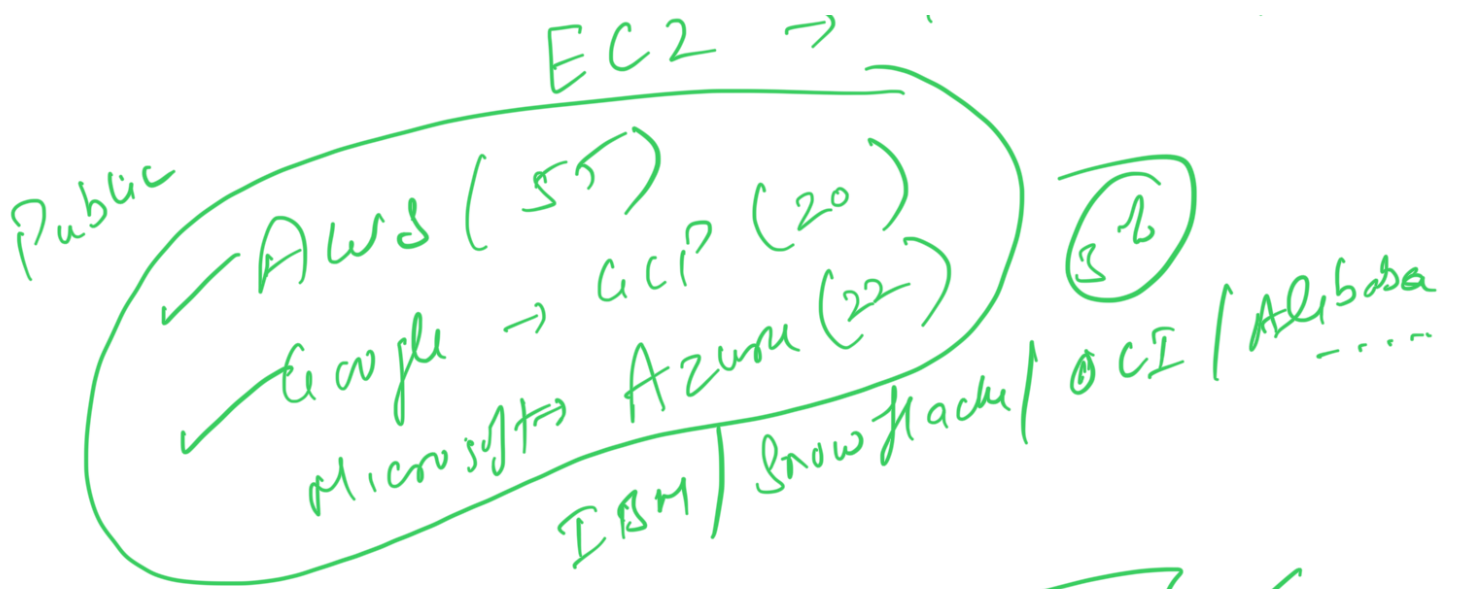
1.5L

M9

updates  
patches



1 minute  $\Rightarrow$  0.0005



AWS Redshift

↳ enterprise relational DBs that store up to PB/EB of data.

(client)

① Data Apps :-

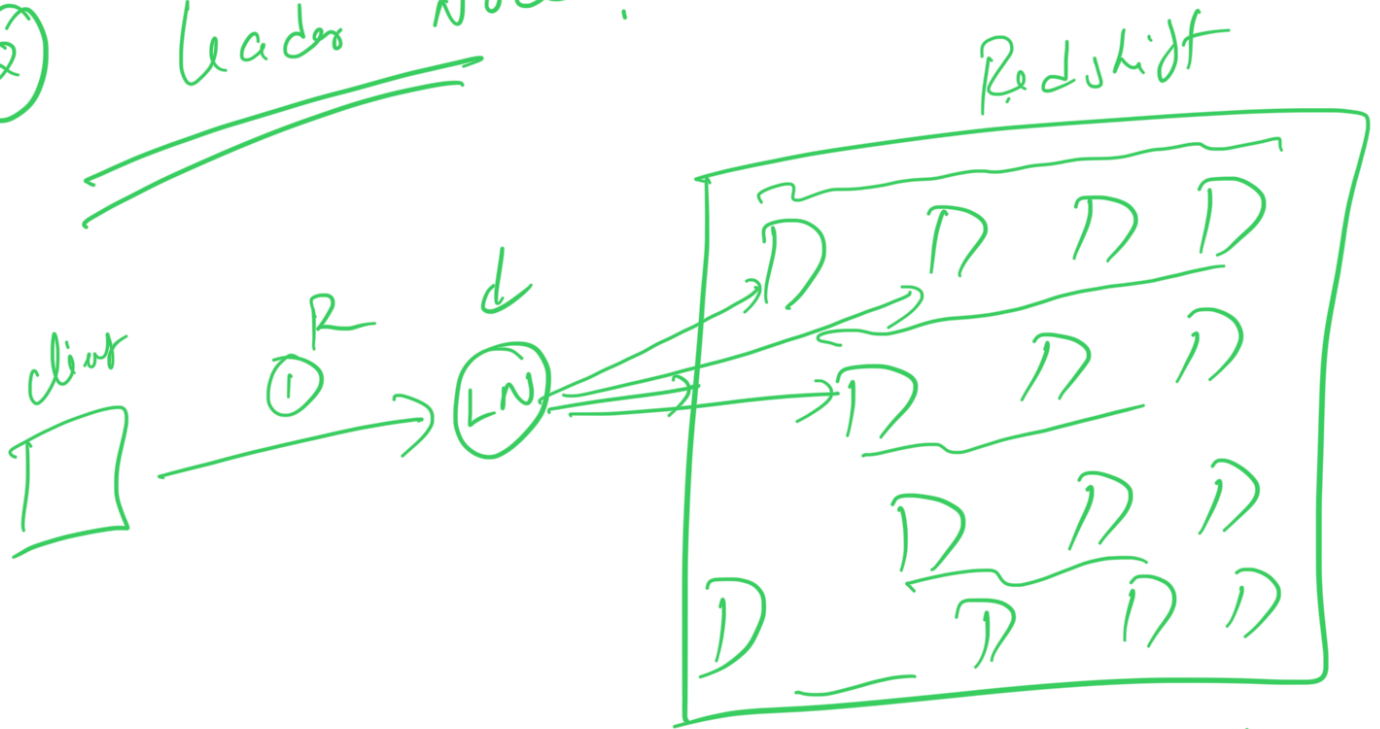
① Data Integration (Load)

② workflow orchestration  
 ↳ these are the apps that run batch jobs on a schedule.

... Analysis :-

## ③ Adv h-oc

② Leader node :- ④ Communicate with Data apps.

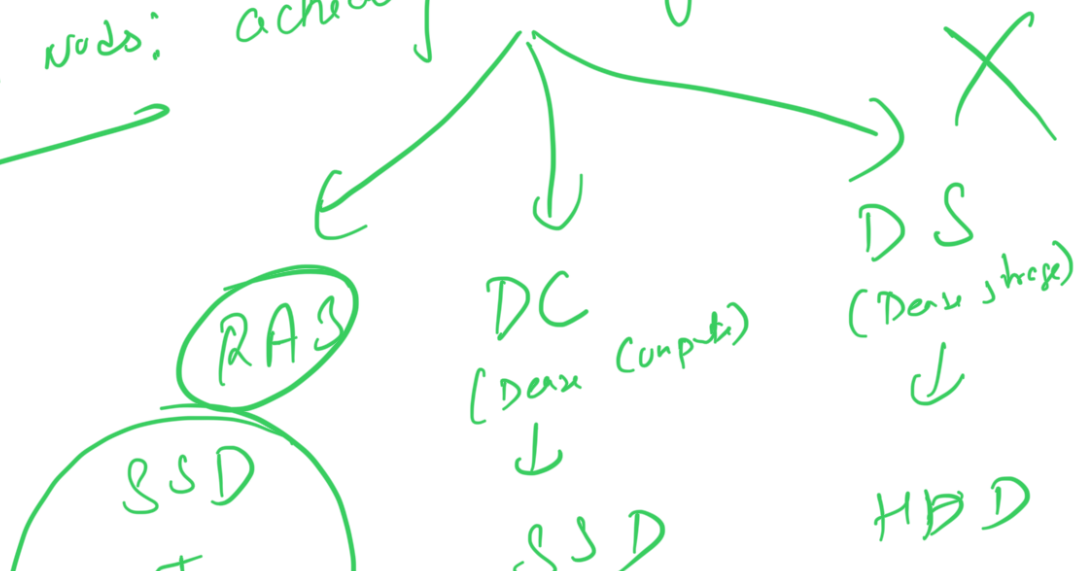


⑤ Distribution of workloads.

⑥ Caching of query results.

⑦ Maintain metadata.

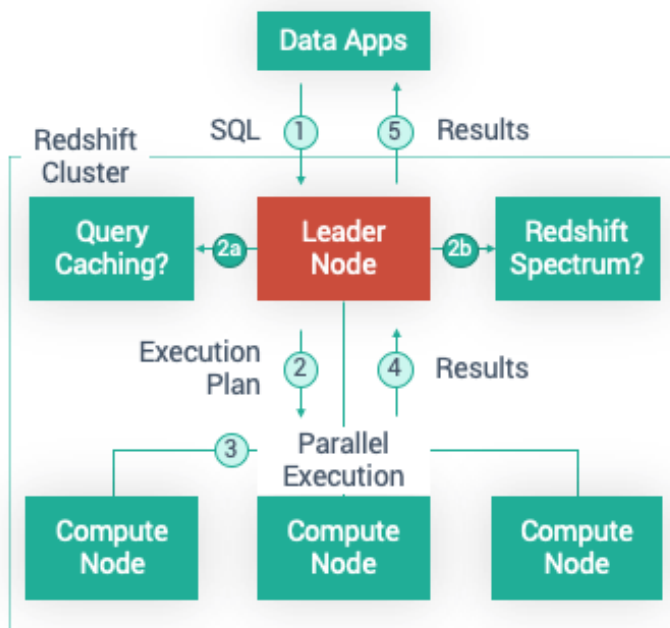
③ Compute nodes: actually store your data



SS

④ Spectrum :-

## Amazon Redshift Architecture: Query Caching & Redshift Spectrum



- 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.
  - 2a Determines if cached result is available
  - 2b Determines if query goes to Redshift Spectrum
- 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.