

Agenda

- a. Parquet Files demo in hive
- b. Data warehousing in cloud Redshift
- c. Distributed system and core components in Hadoop.
- d. YARN architecture in detail

Parquet File Format

→ AVRO → JSON File

→ text
→ Binary
→ Avro

→ Parquet File Format

Columns					Format	Location
id	name	Sal	Dept			
1	Amit	7000	IT			Delhi
2	Amit	9000	IT			Ahmedabad
3	Rakesh	10000	IT			Kochi
4	Abhishek	10000	IT			Panaji
5	Ayush	9000	AD			Delhi

Ayush → location
→ Sal
→ Dept

100,000

25 cols

25 cols

Delhi, 9000, AD

1, 2, 3, 4, 5
Amit, Amit, Rakesh, Abhishek, Ayush
7000, 9000, 11000, 10000, 9000, 11000

5, Ayush, -, Delhi

Delhi

5000

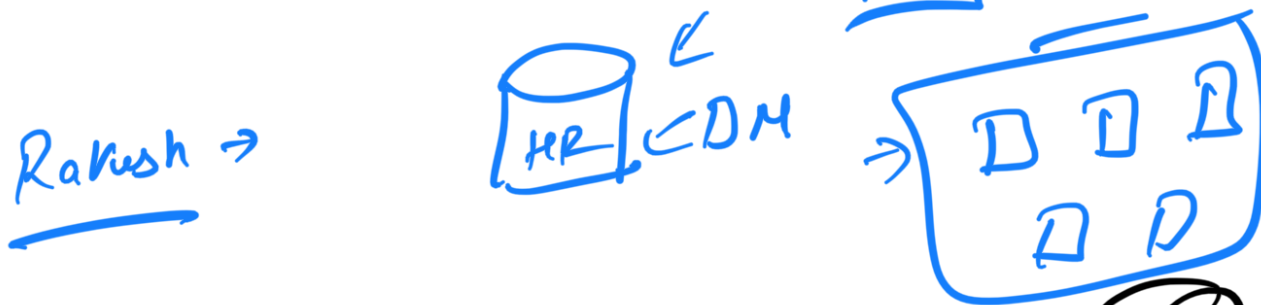
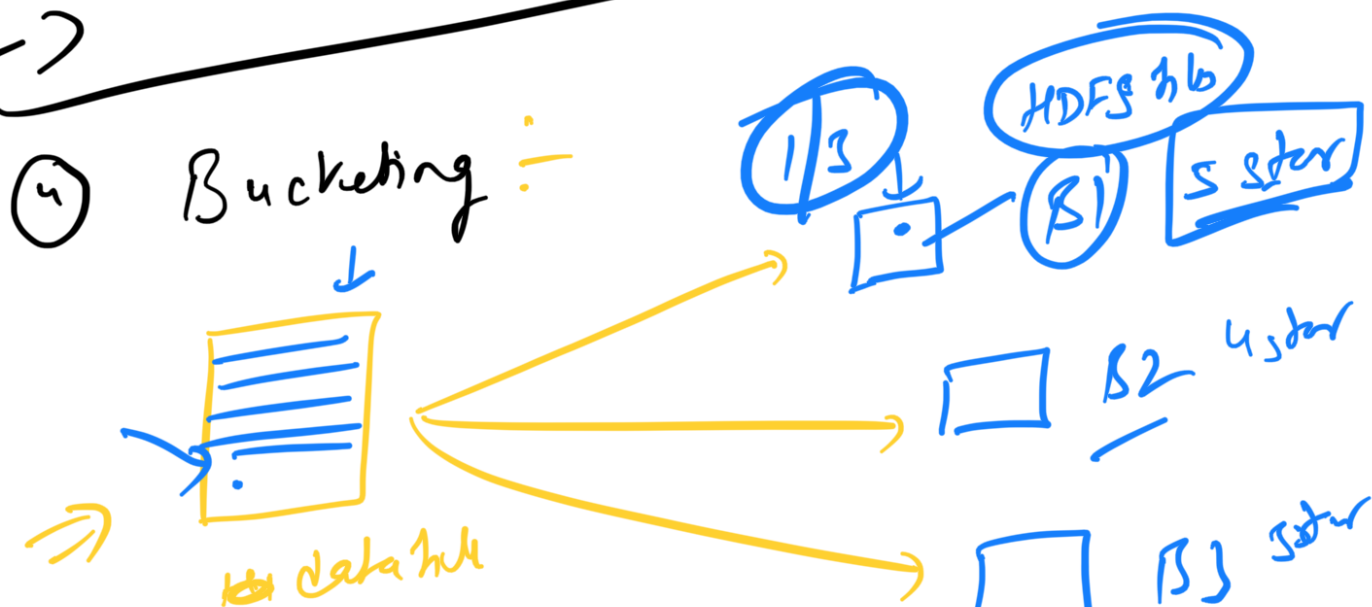
Null
5000

CRUD

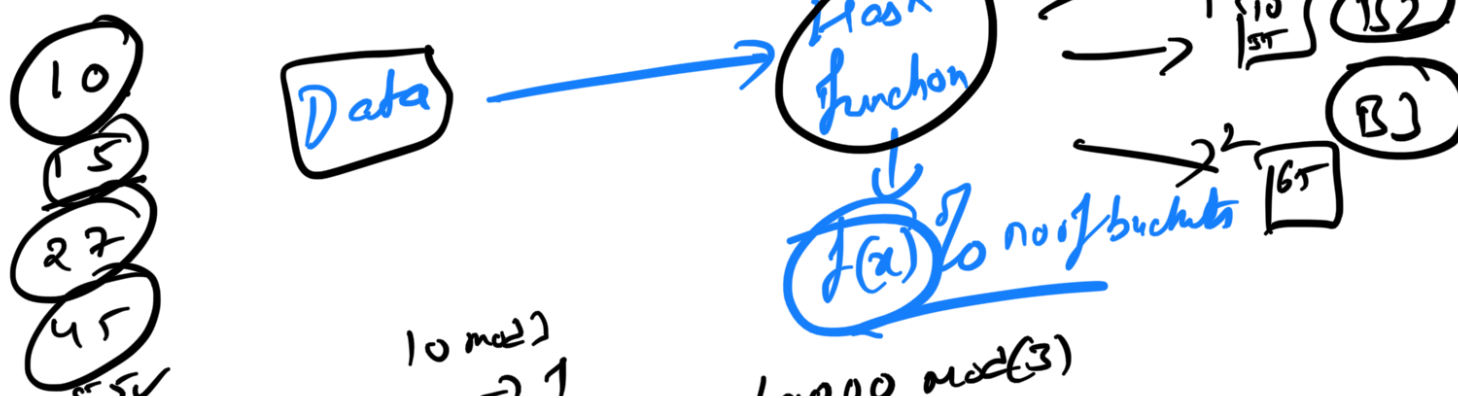
⇒ update

IS
25

User profile Avatar
 Registration_date = Accr ... UA | AF4wUA
 id = 1 | 2



Hash function



55
65
75

75

1000

⑤

Cost Based optimizer

① Query Vectorization

Set hive.vectorized.execution.enabled = true ; **false**

→ one by one → Batch

Hive

↳ TikTok

↳ Walmart → \$1 billion

CDH
→

Apache Hive

AWS = Redshift

Azure = Synapse

GCP = Big Query

10/10

⇒ 8



0008-0

121008
122005

↓

→ delete

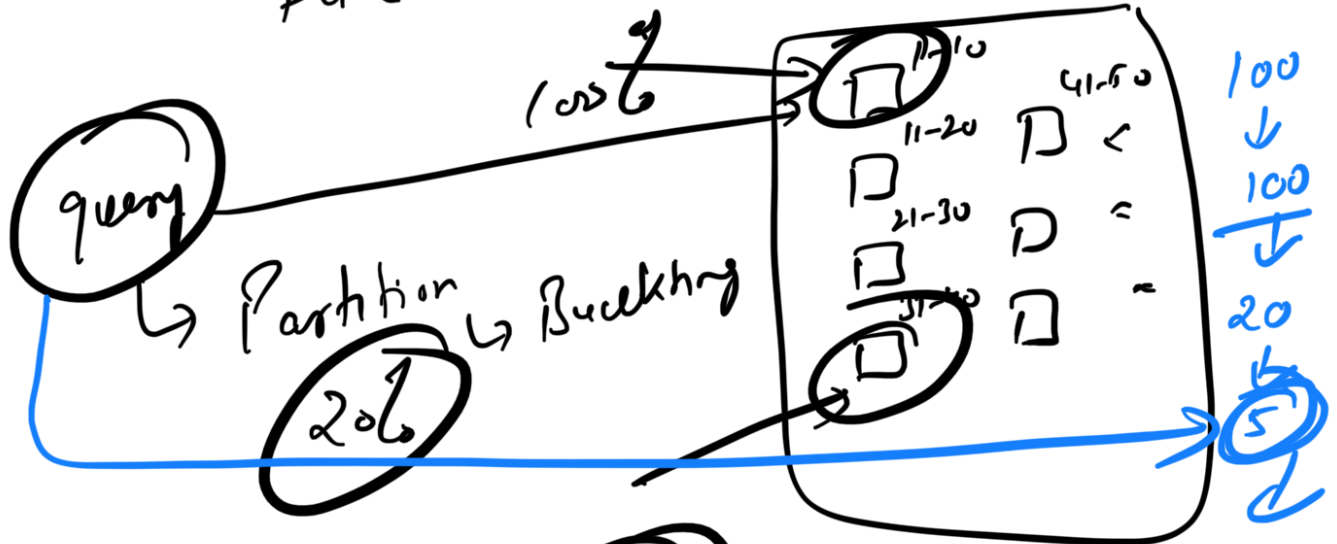
56000
44000

Part 0-0000
0000
0001
1 mapper

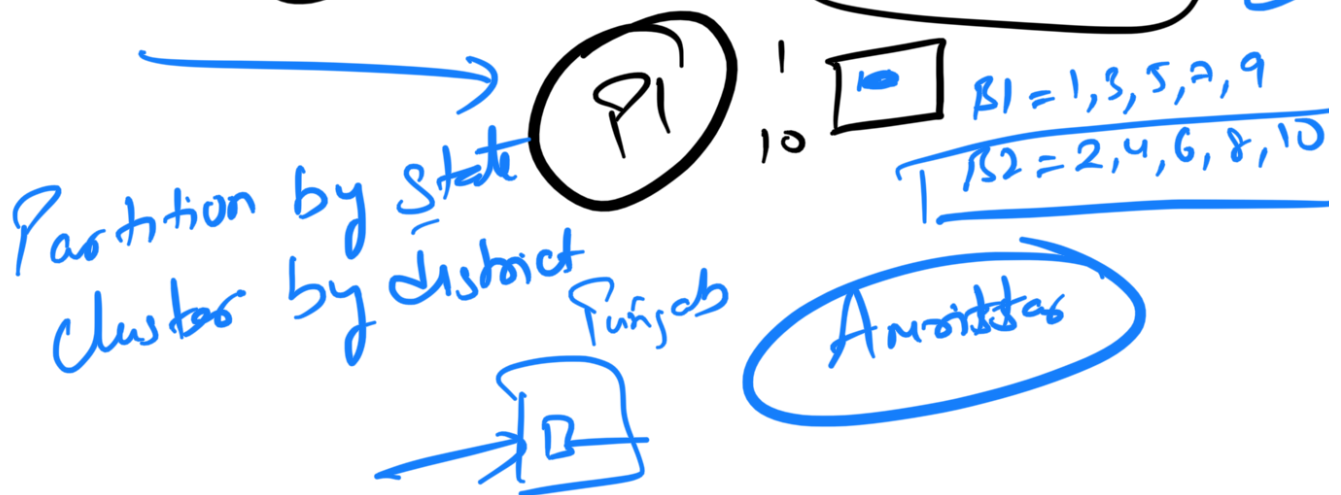
Part 0001
0002
0003
...

Here

100



⇒



AWS Redshift

Hardware

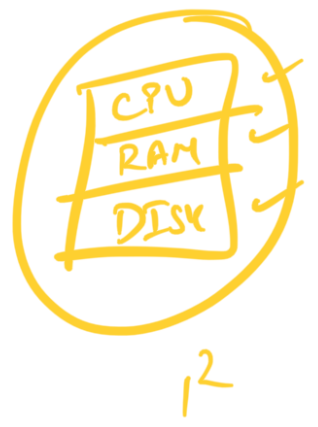
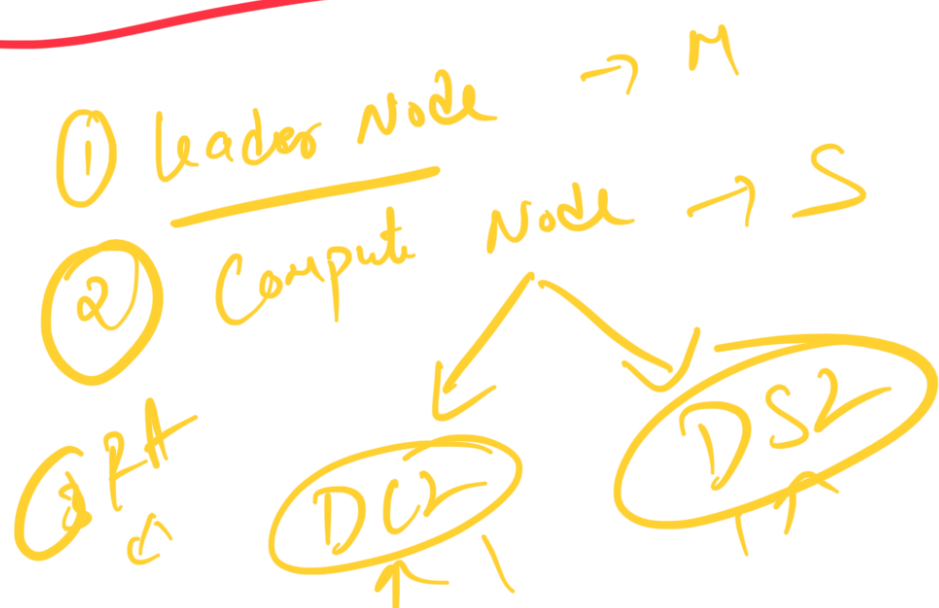
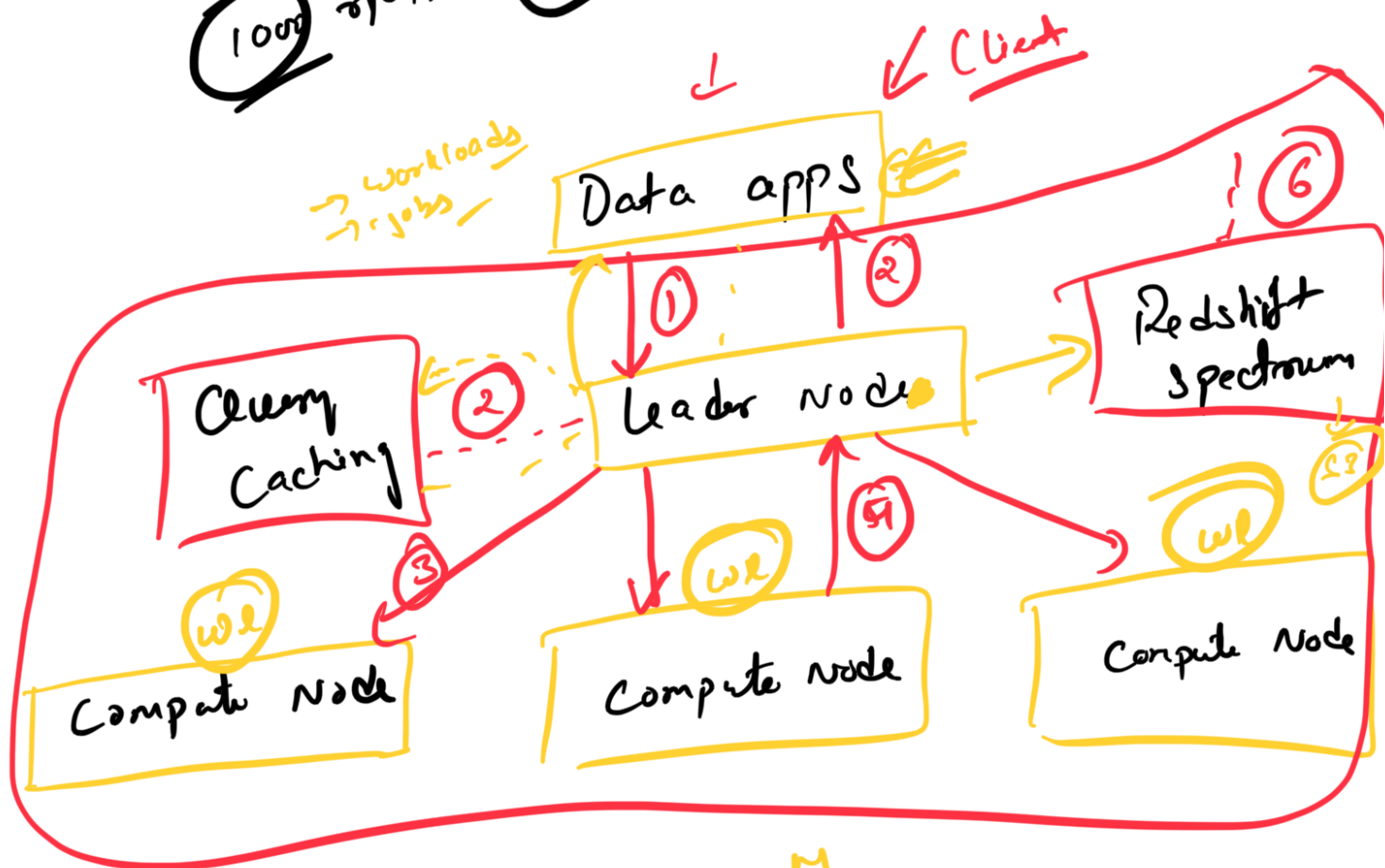
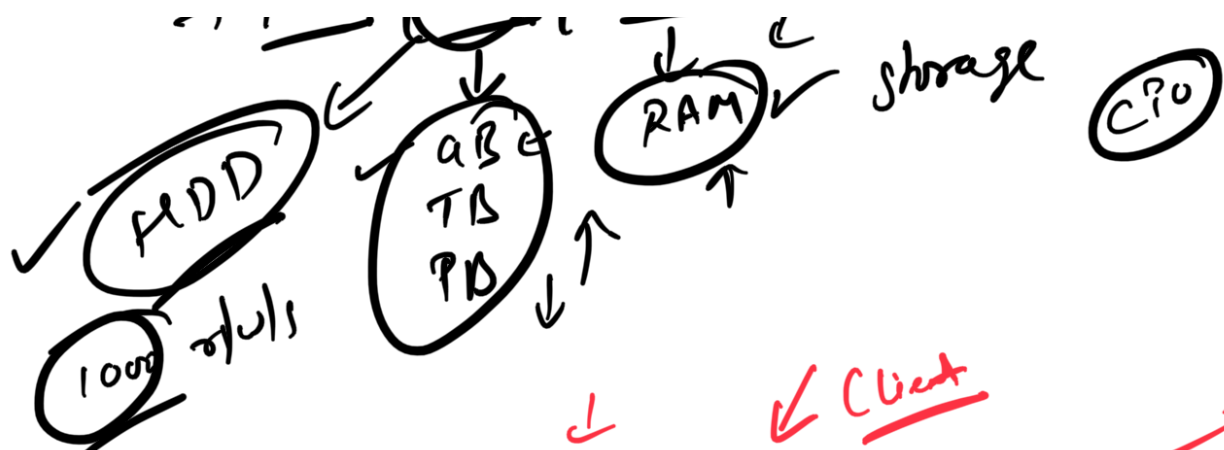
RA

DS

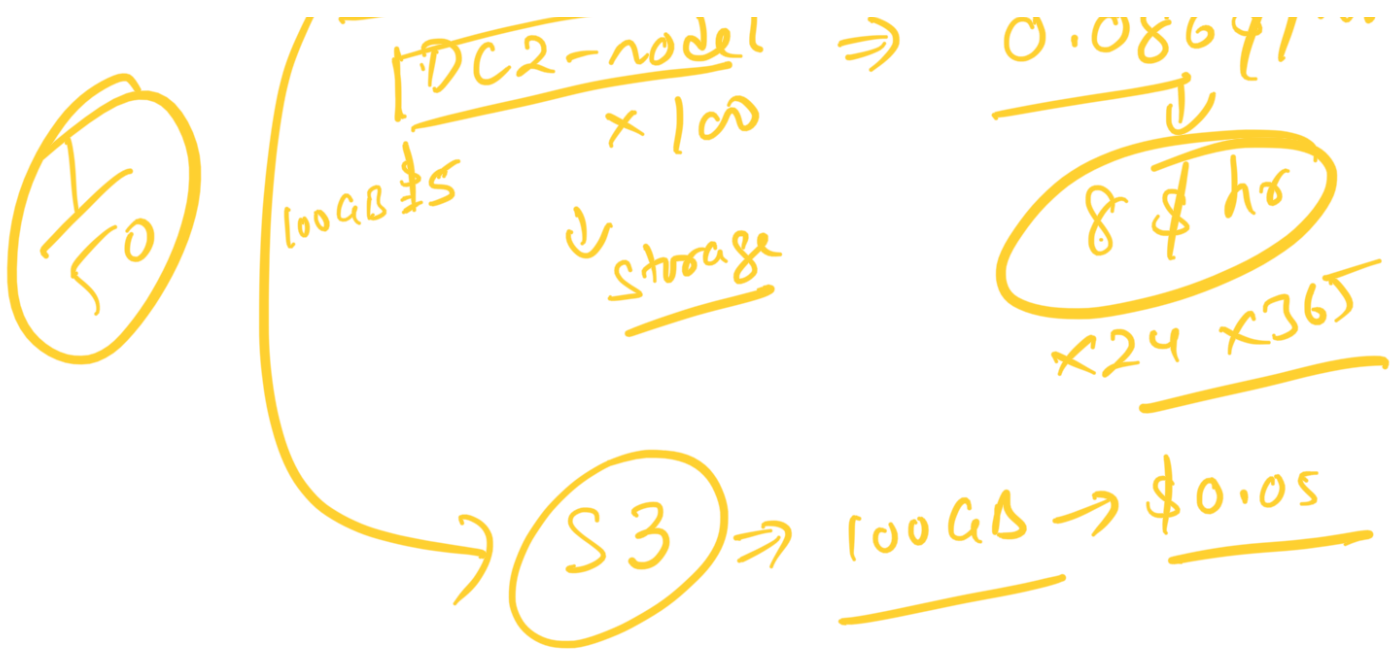
DC

SSD

6000

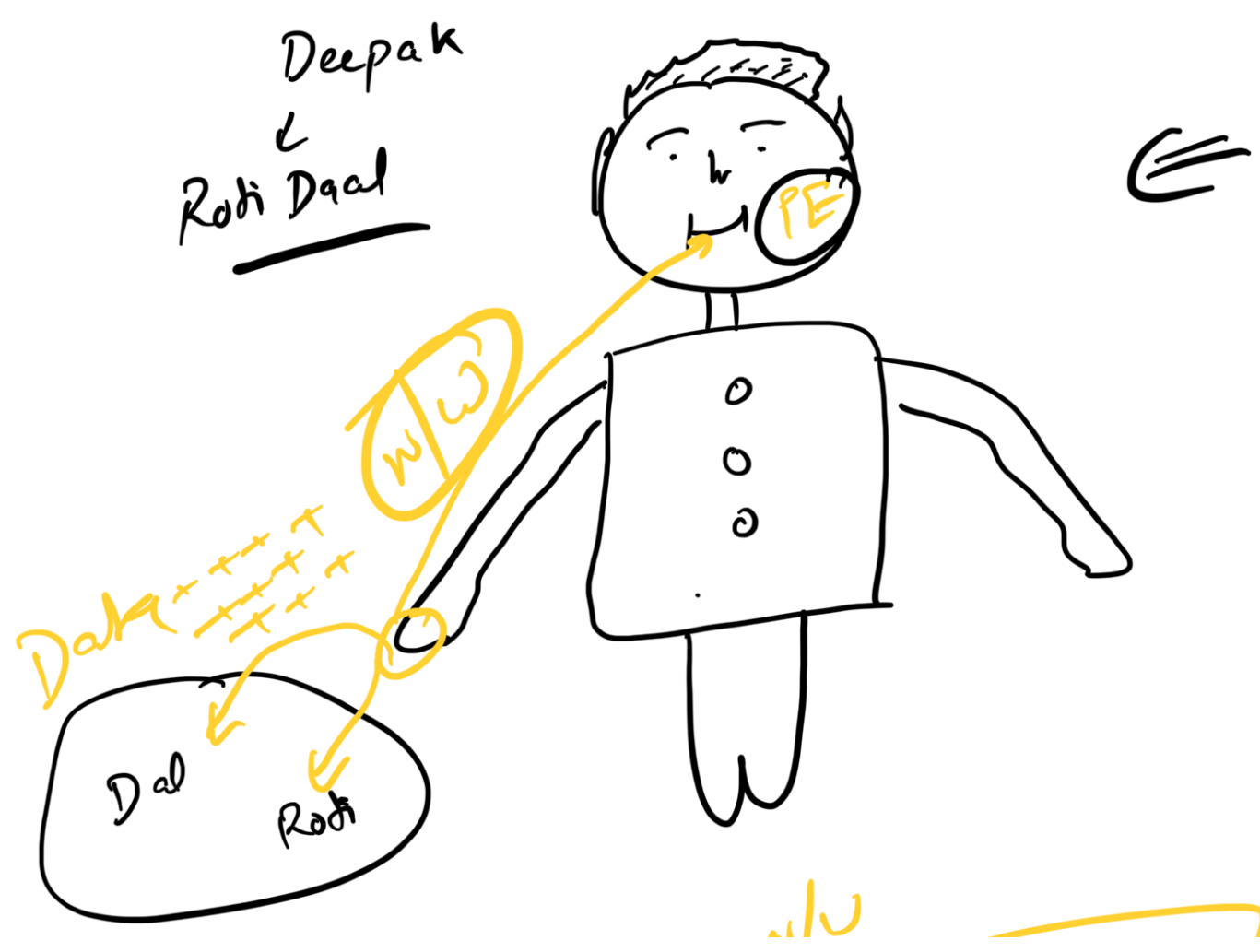


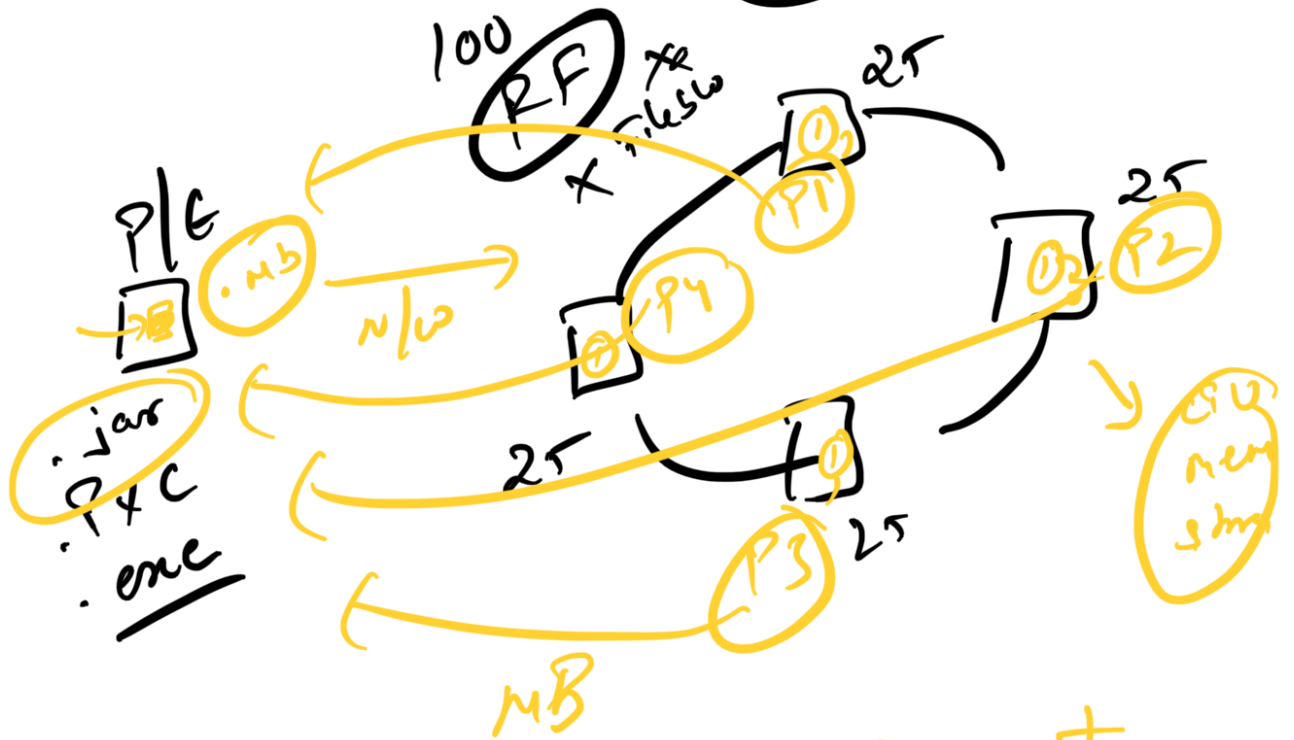
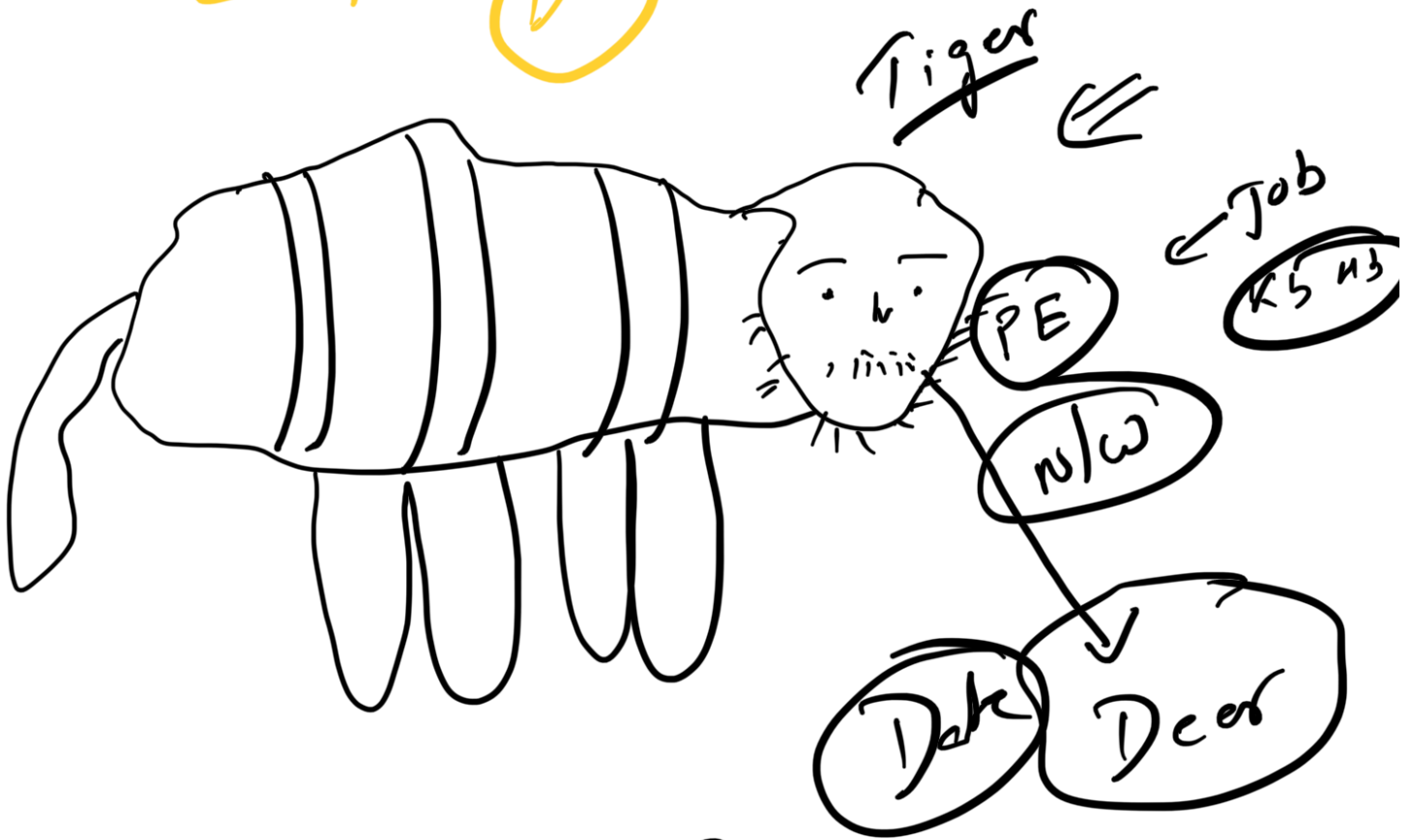
Redshift Spectrum



Redshift \hookrightarrow SQL \hookrightarrow ML \leftarrow

How Hadoop Process Data ?





Hadoop \rightarrow Ecosystem
components

