

Agenda

Problem Statement

Existing Problem

So what is Data Lake ?

How are they different from Data warehouses?

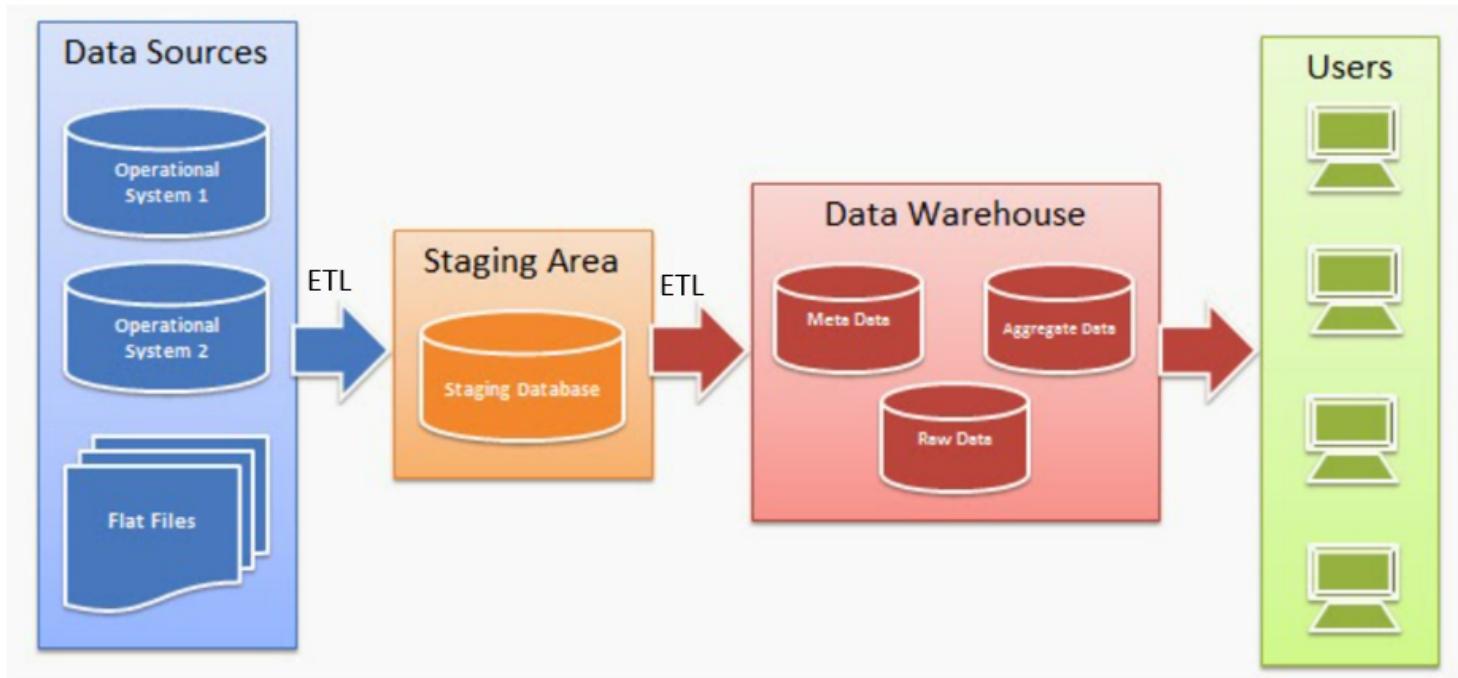
Why Netflix choose data lake

What are the components involved in Data Lake Architecture

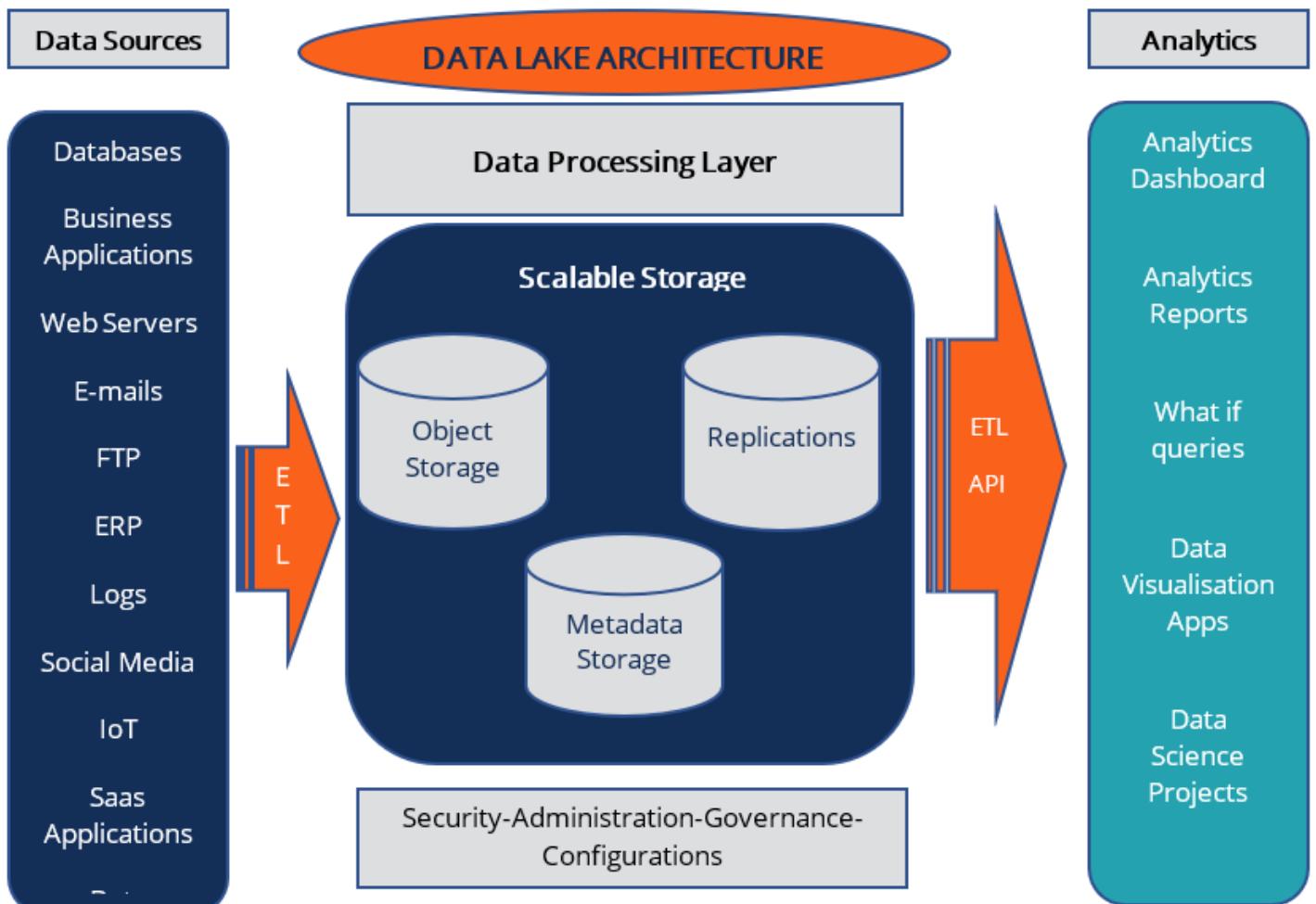
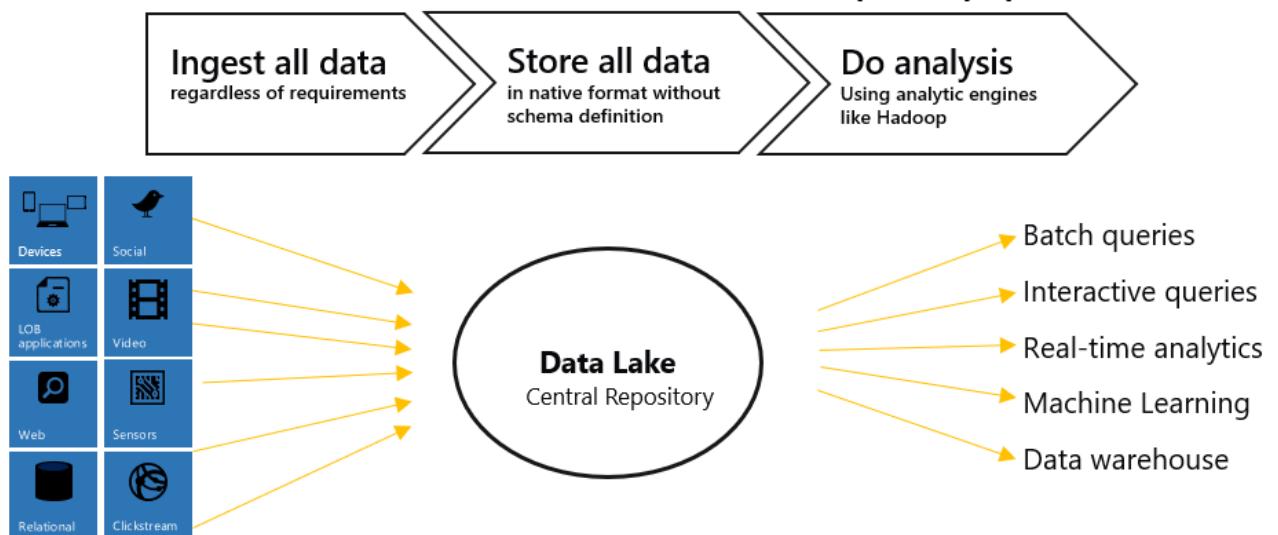
How to create a Successful Data Lake

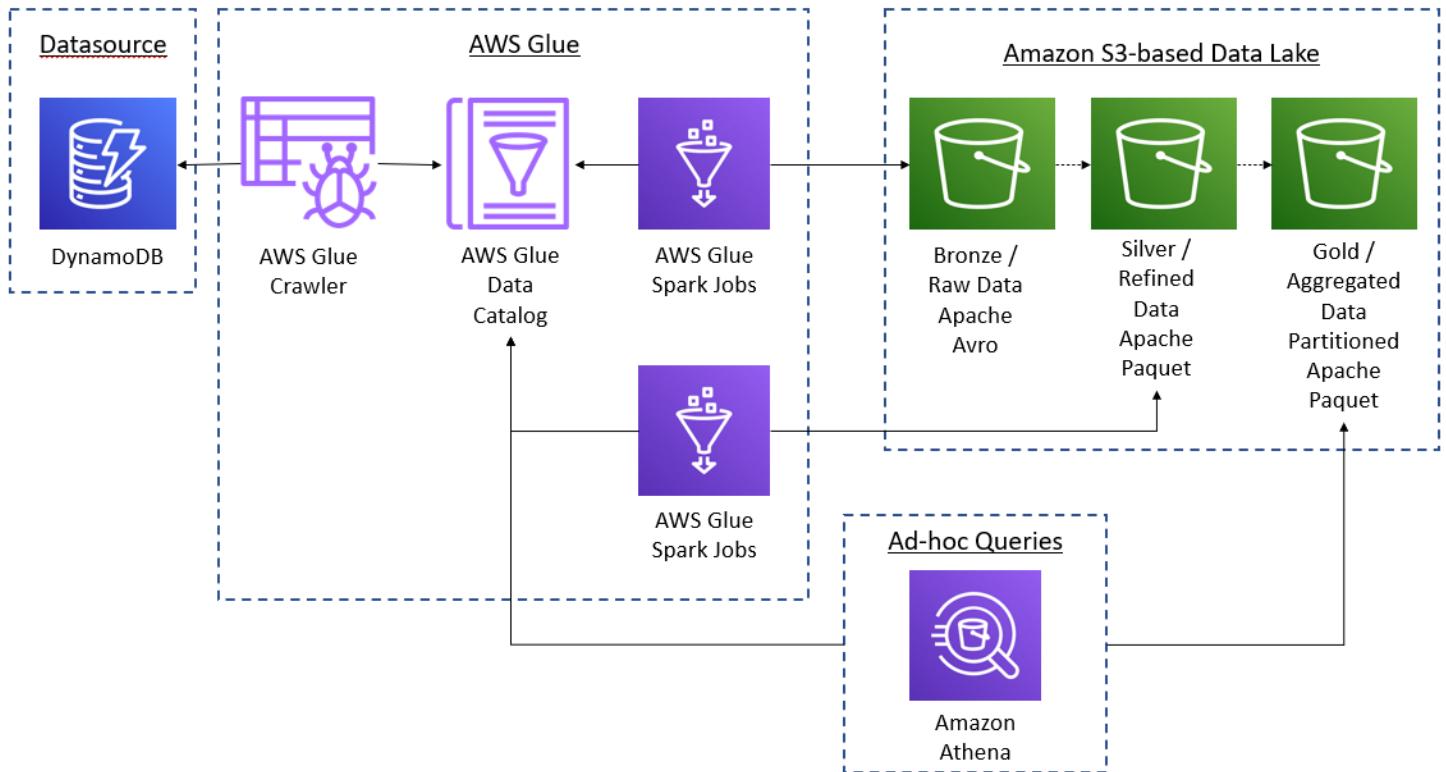
How to form a Data lake with AWS Components

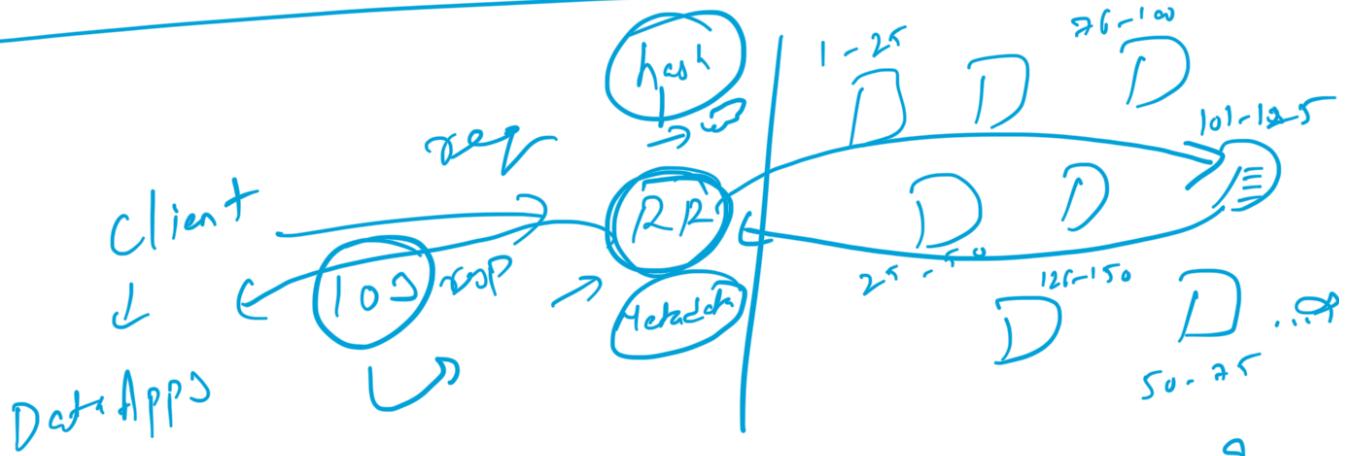




The Data Lake Uses A Bottom-Up Approach





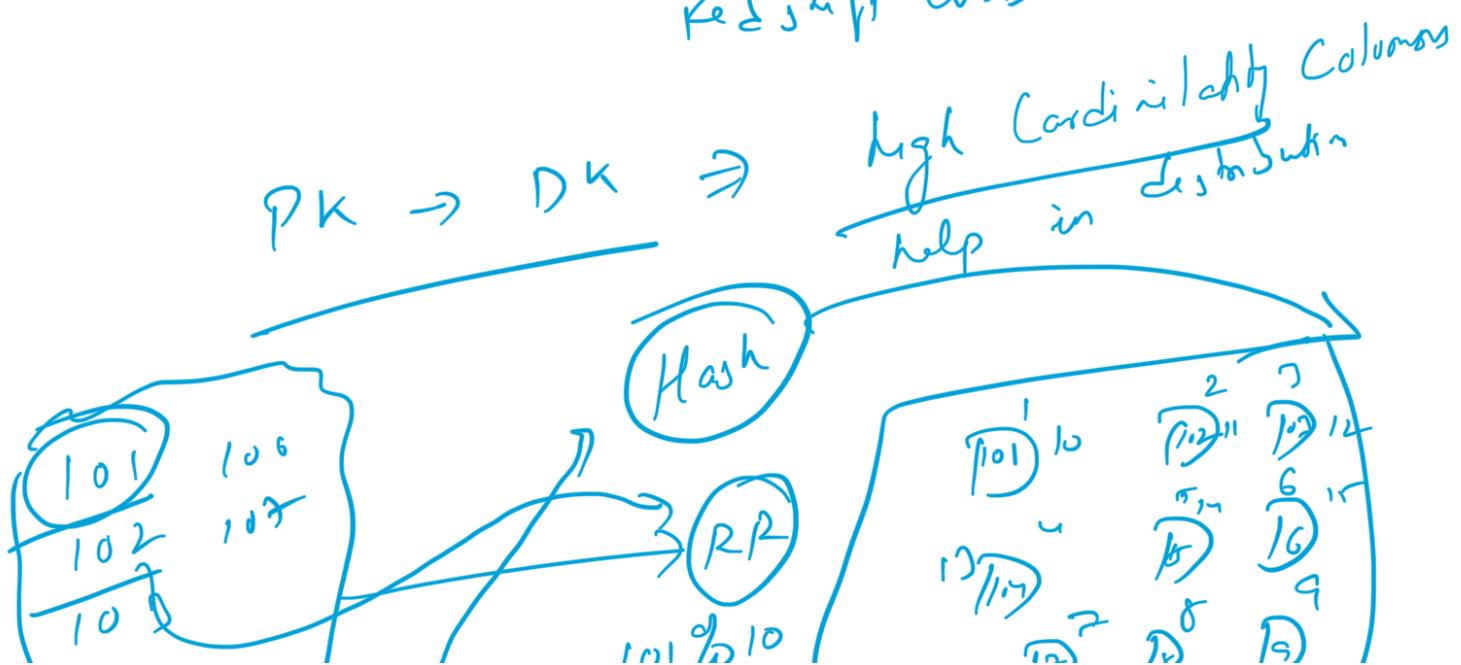


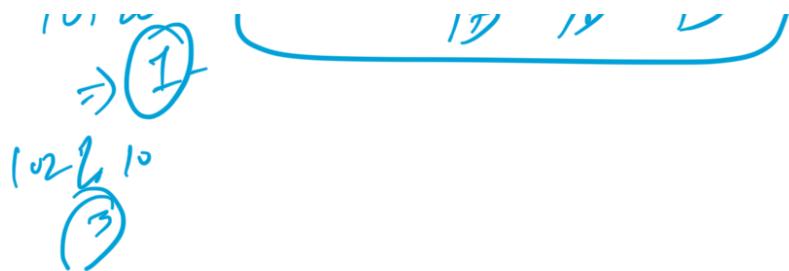
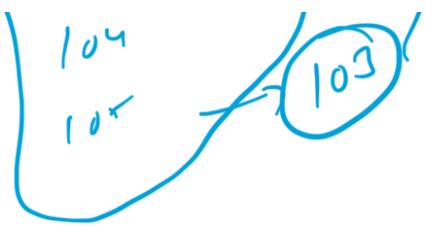
→ How Data is actually Distributed?

- Distribution key / Partition key
- Root key

(a) Distribution key :

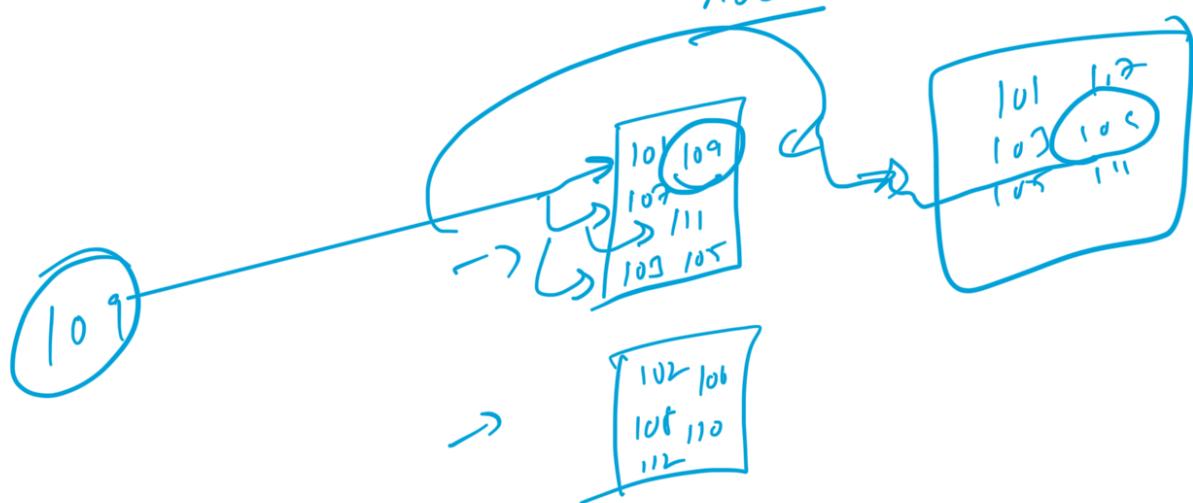
↳ determines how data is distributed across different nodes in your Redshift cluster.





② Sort Key

↳ how the data within each node is sorted on disk.



⇒ Distribution style

↳ how the data in a table is distributed across nodes in a cluster.

① Even

→ Dist Style Even
→ Data is distributed evenly across all in the cluster

② Key

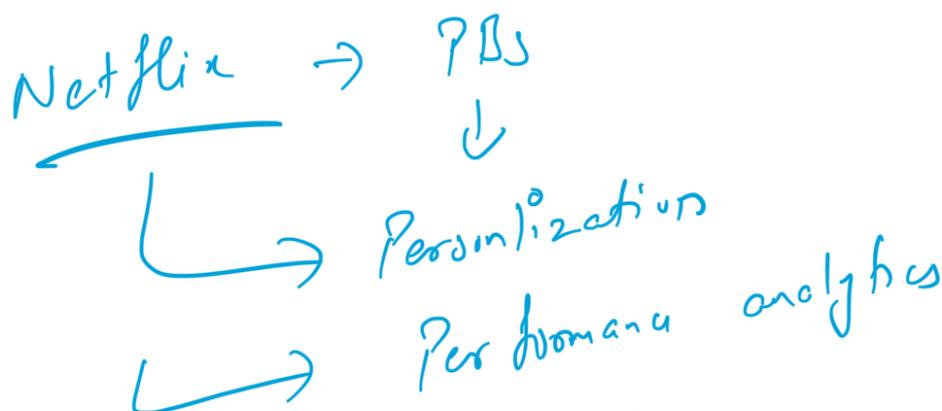
→ Dist Style Key
→ DistKey (Per node)
→ Data is distributed based on value of

③ All

→ DistStyle All
→ A full copy of entire table is shared among users.

- Each node will get an approximate equal portion of data.
- Which one to choose
- ↳ Key : Use when you have a clear join col that can act as dist key. Particulay for large fact table.
 - ↳ Even : Use when you want even data distribution across nodes and there is no clear DK.
 - ↳ All : Use for small tbs that are frequently joined with large tables, to avoid data movement during joins.

Data Lake



↳ *Implementation*

Issues

- ↳ To store all data
- ↳ No early access to data
- ↳ Scaling
- ↳ Price

Solutions :-

- ↳ Data streaming ✓
- ↳ No Sql DB
- ↳ HDFS

Data Lakes

- ↳ are central data repositories used to store any data at any scale.

Different from Data warehouse



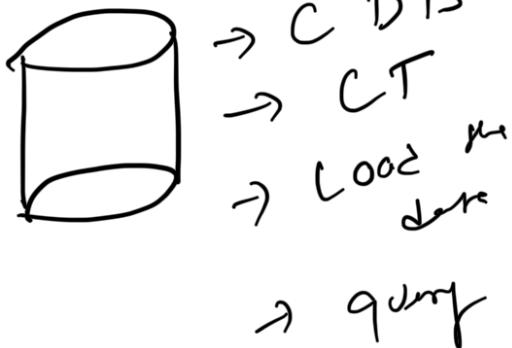
↓

↓ Data Lakes

Data Warehousing

① Store & structure Data

② Top down Approach



③ ETL

④ BA | DA

⑤ TBS | PBS

⑥ Expansion

① Store structure,
Semi structure or
Unstructured.

② Bottom up approach

③ ELT

④ DS
PBS | GR, HIS
⑤ cheap

Yule BroadCon Netlink Arrested UDo

Paytm

SS →
Azure blob
...NE)

DL

X
DL = ELT

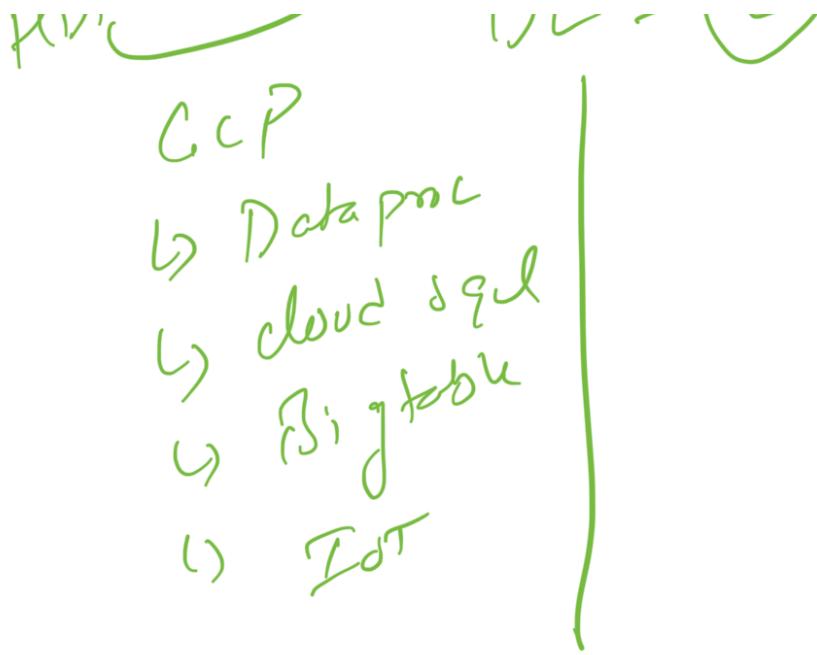
AWD

↳ Glue

↳ S3

→ DynamoDB

↳ Athena



① Data Sources :- providers for real time business
↳ auto to DL.

↳ IoT | sensor | DD | Kinesis
nodel | FTP | CSV, ...

② Data storage :-

ⓐ Object storage
↳ objects (large volumes of unstructured data)

ⓑ Metadata storage
↳ provides properties about the object.

* ⓒ Replication
↳ replicates objects across multiple locations of DL for DR and backup

③

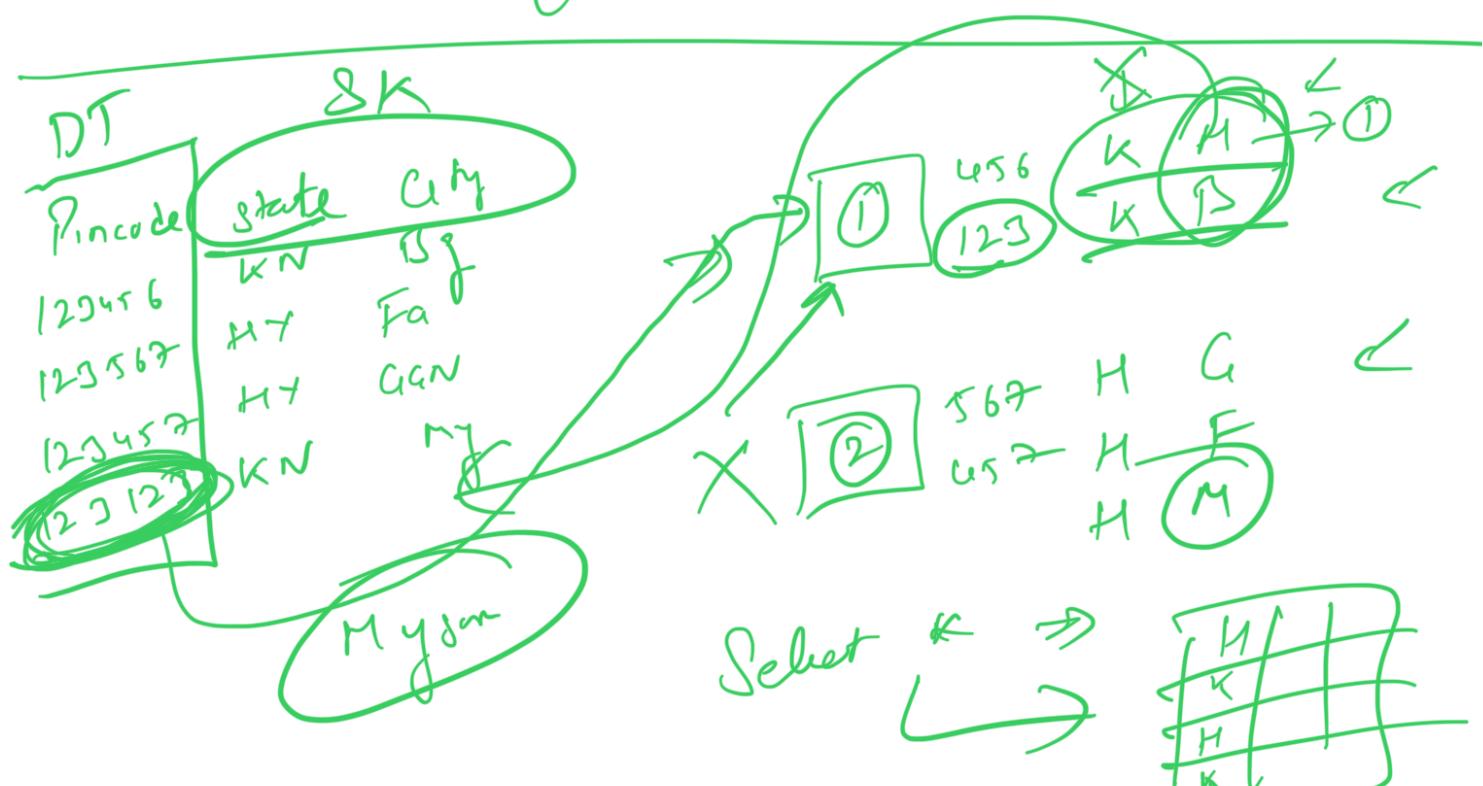
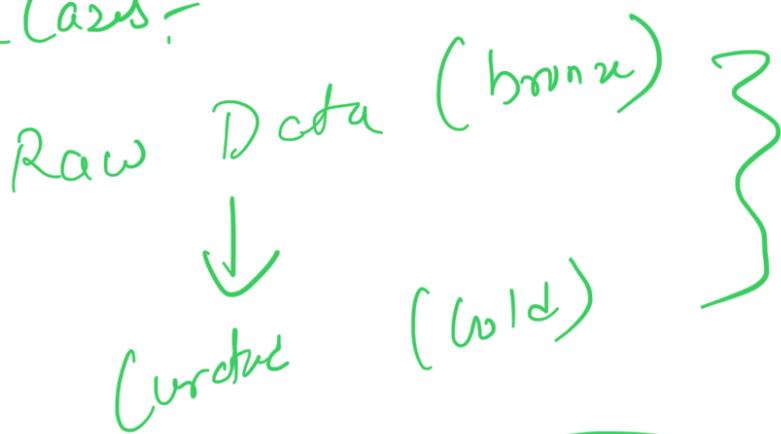
Data Processing Layer

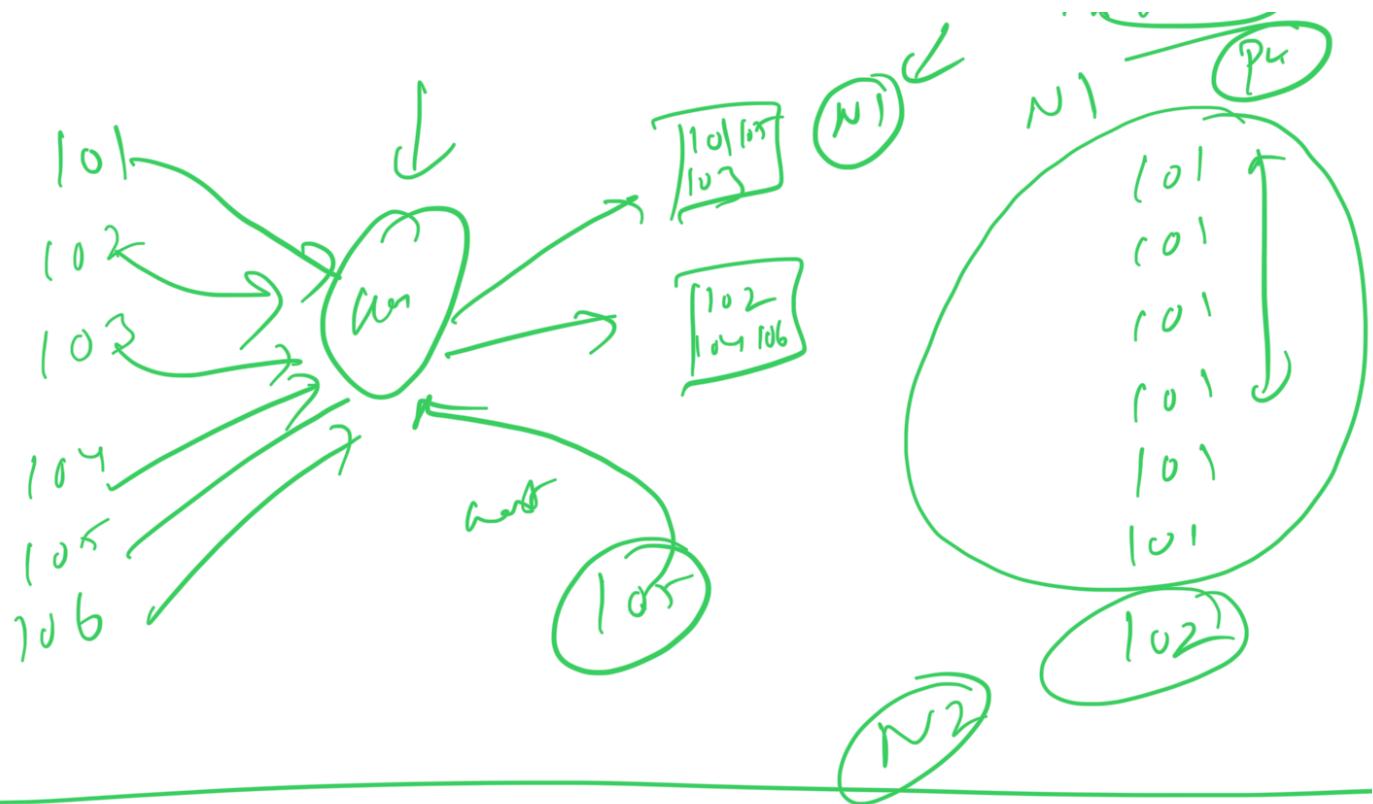
↳ Retrieved data is processed to produce info that is applicable to business.

④

Data Analytics Layer: Dashboards + Reports

Use Cases:





DCS

Sales

↳ txn-id $\circledcirc P_4$

↳ C_id

↳ txn-date

↳ P_id

↳ amount

$\rightarrow \text{DK}$

txnid	C_id	txn-date	i_id
1	1	02-dep	iR
2	1	05-dep	Cover
3	2	02-dept	i1-o2
4	3	06-dept	Cover

↓
1 Billion Rows

Sort Key	
1	02 Dept
2	04 Dept
3	03 Dept
4	06 Dept

$\circledcirc PR$

↓
Housing
1.9.2

$\boxed{13}$

$\Rightarrow N_1$

3	03
1	02

$\boxed{24}$

$\Rightarrow N_2$

Date \rightarrow Desc

1 \xrightarrow{w} 1