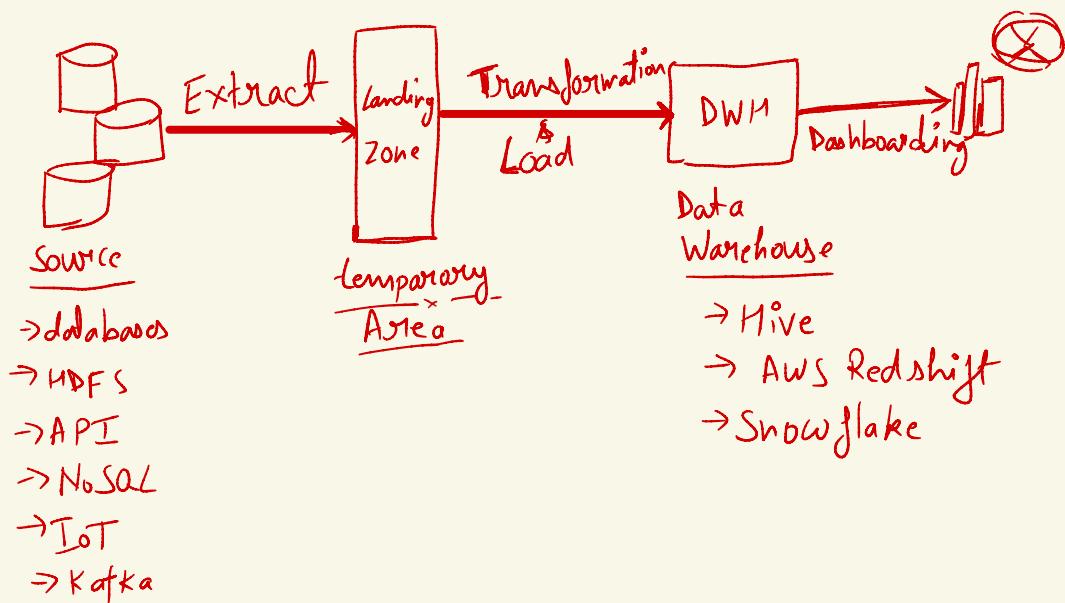

2017

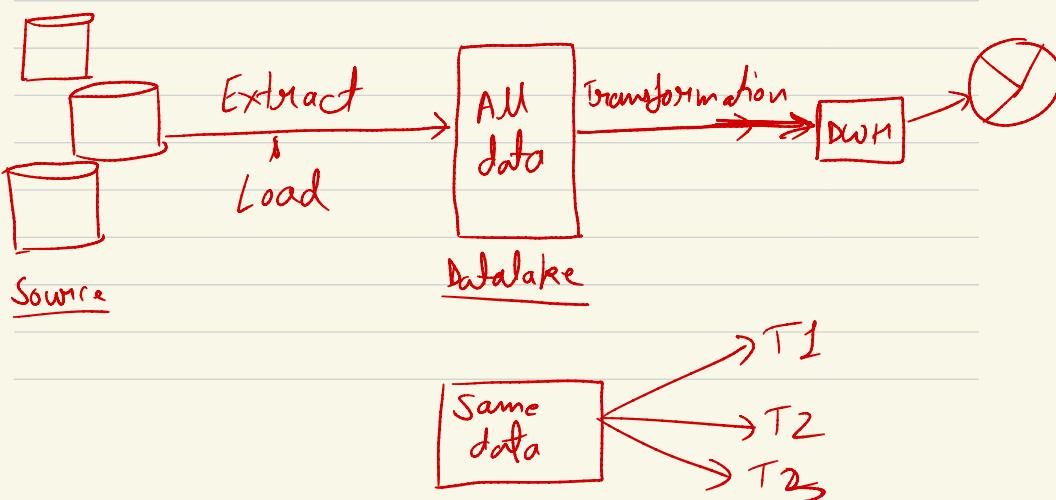


ETL vs ELT

ETL → Extract transform Load



ELT → Extract Load Transform



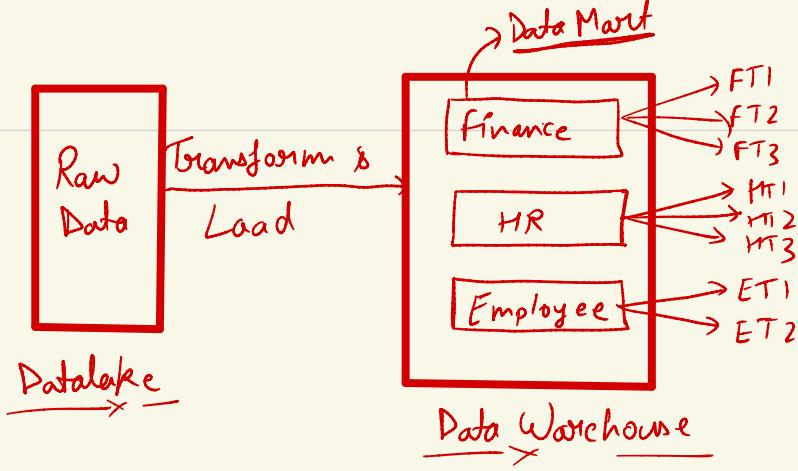
* Datalake: Any system/tool/file system where we keep all the data coming from different sources and different formats. Ex- structured, unstructured, semi-structured. Examples for datalakes:-

- HDFS
- AWS S3
- Azure Blob
- GCP Storage

* Data warehouse: Any system/tools/framework which holds transformed data will be known as Data warehouse. It represents data in only structured form.

Example: → Hive
→ Aws Redshift
→ Snowflake

* Data Marts: Category level separation inside a Data warehouse is known as Data Marts.



Code Vs NO-Code

↓

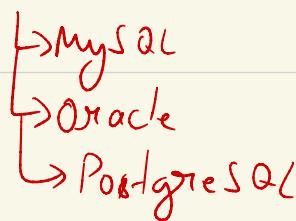
Spark (ETL) Informatica

OLTP vs OLAP

x) OLTP (Online Transaction Processing) : Systems

designed to support transaction level operations like Insert - Update - Delete etc. They support ACID properties and main goal is to capture data in real time. In OLTP data resides in its normalized form.

Example: RDBMS



x) OLAP (Online Analytical Processing)

It is designed to support Analytical Queries

It keeps data in denormalized form and doesn't care about ACID properties. Main goal is to help users to query historical data & insights out of it.

Ex: Redshift

Hive

NoSQL

Snowflake

Important Terminologies in Data Warehousing -

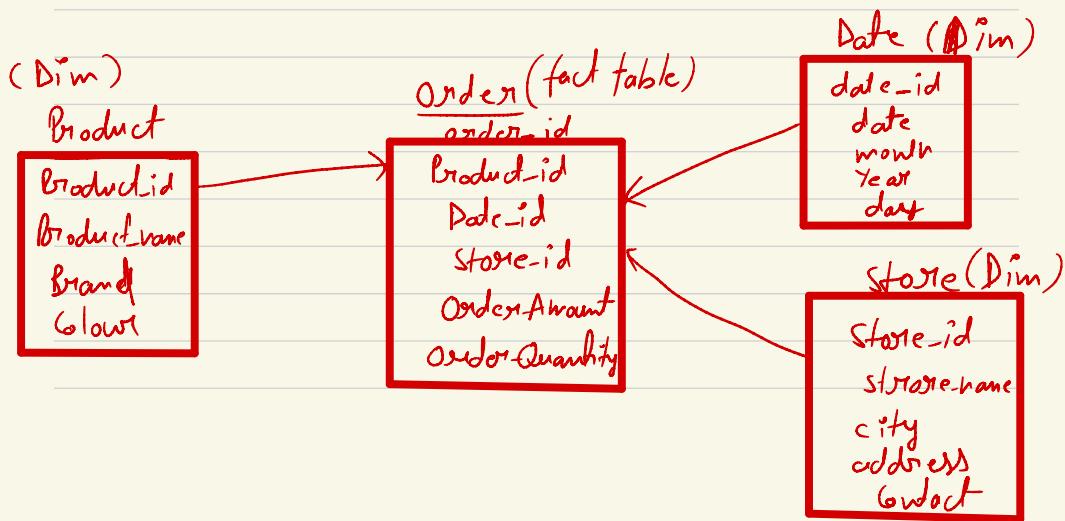
- ① Dimension Table
- ② Fact Table
- ③ Relationship
- ④ Star Schema
- ⑤ Snowflake Schema

* Dimension Table:

- It contains dimensions of a fact.
- They are joined with fact tables using foreign keys
- Dimension tables are de-normalized
- It offers descriptive characteristic of the facts with the help of attributes

*) Fact Tables

- Contain those numerical attributes which helps to derive meaningful values.
- It holds foreign key of dimension table.



Employee detail

emp-id	emp-name	dept-name
1	abc	Software
2	xyz	Software
3	klm	HR
4	MNO	HR
5	efg	Software

Engineering Department

De-hominalized

For heavy analytical
queries it is good fit.

another approach

Department

emp-id	name	dept-id
1	X-Y-Z	100
2	abc	100
3	klm	200
4	MNO	200
5	efg	100

dept-id, dept-name

100, Software

200, HR

Normalized

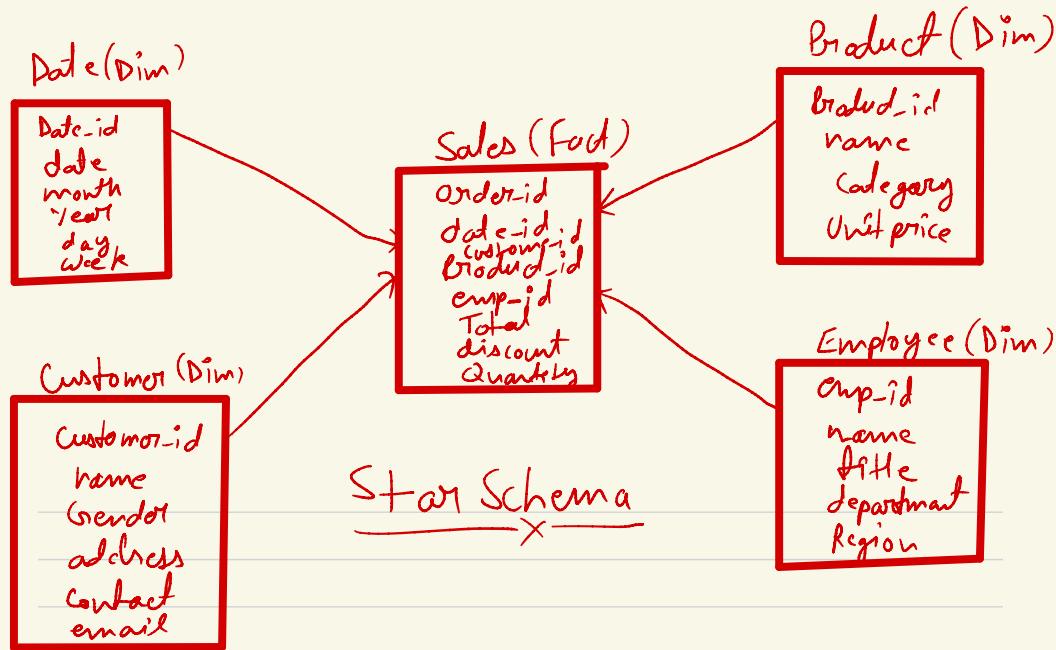
We need to perform

joins which is

bad for analytical
queries

Star Schema: In this kind of schema design

One or more fact tables will there and they can be connect with any number of dimension tables.

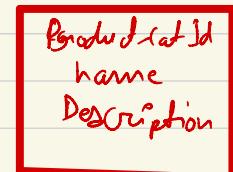


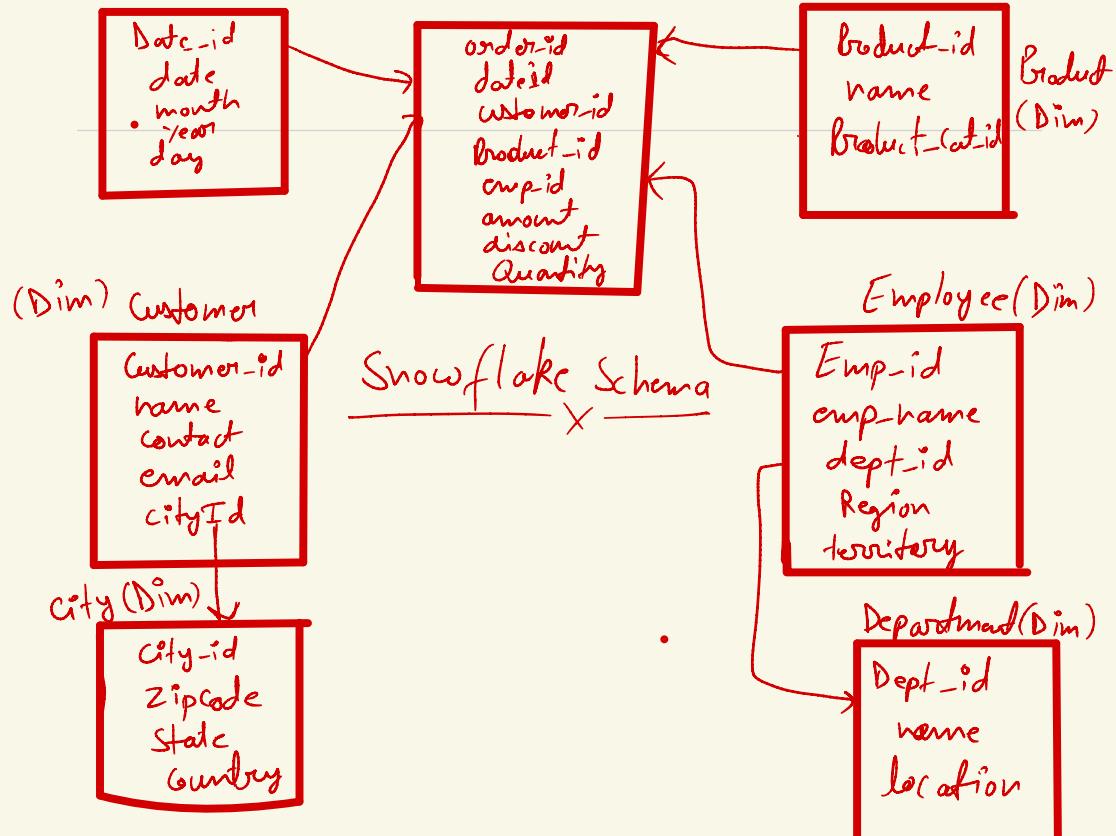
Showflake Schema: In this schema design

dimension tables will be in normalized form

(Dim)
Date

Sales (fact)





Star Schema Vs Snowflake Schema difference

- 1.) Star Schema dim tables are not normalized, Snowflake schemas are normalized.
- 2.) Snowflake schemas will use less space to store dimension data but it is complicated
- 3.) Star Schema will only join the fact table with dim tables and that is why queries are faster in star schema than snowflake.

- 4) Snowflake schemas have no redundant data that's easy to maintain.
- 5) Snowflake schemas are good for data warehouse, star schemas are better for data marts with simple relationship.

Data Warehousing Case Study -

- ① Four Step Dimensional Design Process
 - 1) Select the Business Process
 - 2) Define the grain ("How do you describe a single row in the fact table?")
 - ③ Identify Dimensions - Who, What, Where, When, Why and how?
 - ④ Identify facts

Solution

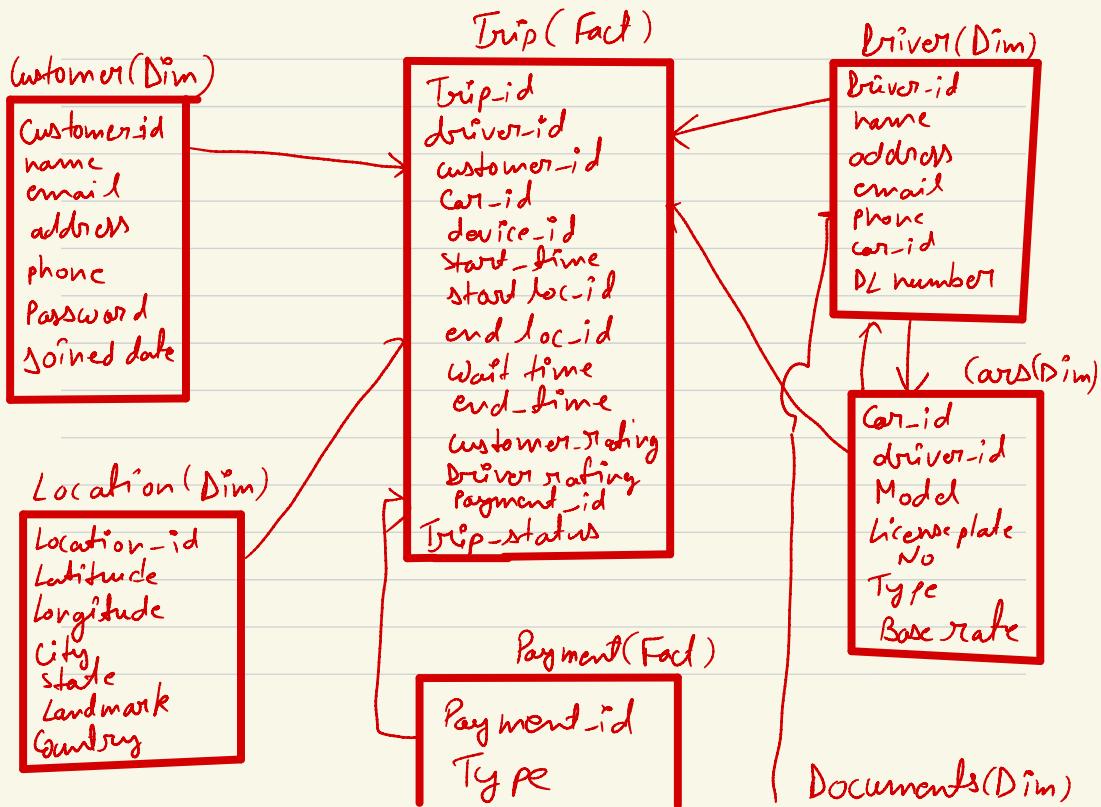
- ① Business Process :- Taxi company would like to design a data model to capture all critical data elements.
 - Track rides done by driver & their performance
 - How many rides are happening to a common or famous destination each day

- How many trips cancelled per day
- How many rides and the average price during the peak hour per day

② Grain - Individual trip on each transaction level

③ Dimensions - Date, Customers, Drivers, Cars, Documents, Devices, locations

④ Facts - Trips, Payments



Base Rate
Surge rate
Trip Amount
Total Amount
Transaction_id

Doc-id
name
driver_id
Category
code
Country
expiry-date