

## Module 3 - Introduction to Data Warehousing and Multi-Dimensional Modelling.

### \* Operational Support System vs Decisional Support Sys-

O.S.S

D.S.S

- | O.S.S   | D.S.S  |
|---|--|
| ① In terms of <u>purpose</u> , O.S.S automates routine tasks and processes. | ① In terms of <u>purpose</u> , D.S.S provides data and tools for decision making and analysis.           |
| ② In terms of <u>data</u> , O.S.S has structural and transactional data.    | ② In terms of <u>data</u> , D.S.S has historical and predictive data.                                    |
| ③ The concurrent transaction volume in O.S.S is very high.                  | ③ concurrent transaction levels in D.S.S is medium or low.   |
| ④ Regarded as update transactions.  | ④ Regarded as query transactions. (Read-only)  |
| ⑤ O.S.S supports transactions which are happening in real-time              | ⑤ D.S.S supports transactions which have happened a long time ago to analyze and make decisions on them. |

## \* Data Warehousing

- Data Warehousing is a process of collecting, organizing and storing large amounts of data from various sources in a centralized repository.
- The data warehouse is designed to support business intelligence activities such as data analysis, reporting, and decision-making.

## \* Features of Data Warehouse

- ① Subject-oriented - A data warehouse is designed around specific subjects or topics related to an organization's business such as sales, customer data and financial data.
- ② Integrated - Data in a Data Warehouse is integrated from multiple sources, including internal systems and external data sources. Helps data stay accurate.
- ③ Time-variant - Data in a data warehouse is stored over time allowing analyst to track changes in data and identify trends over time.
- ④ Non-volatile - Data in a data warehouse is read-only, so data cannot be updated or changed directly.

## \* The Need for Data Warehousing

→ There are many reasons why organisations need data warehousing.

- ① Centralized data storage - organisations generally collect data from various sources, including transactional systems, websites, social media, articles etc., making it easier to access and analyze data.
- ② Improved data quality - inconsistent data is a major issue, but since data comes from various sources, the accuracy and quality of data is improved significantly.
- ③ Analytics and reports - Data warehouses are designed for analytics and reporting, with tools and techniques optimized for data analysis and visualisation.
- ④ Faster decision-making - Data warehouses provide timely access to the data they need. This is important for organisations working in competitive markets.
- ⑤ Historical analysis - Data warehouses store large amounts of data for a long time, allowing managers to analyze changes in trend patterns over many years.

## 4. Meta data

- Meta-data refers to the info about the data stored in the warehouse.
- It provides a description of data which is stored in the databases, such as origin, format, structure.
- Metadata plays a key role in data warehousing by providing valuable information. This helps in
  - ① Understanding the data - Helps users to analyze and report by knowing the origin, format and structure.
  - ② Managing data quality - Metadata also includes information regarding the accuracy, precision and completeness of the information.
  - ③ Data Integration - Metadata helps integrate data from multiple sources ensure data is consistent across the organization.
  - ④ Data governance - Metadata plays an imp. role in data governance by providing information regarding data ownership, data usage and data security.

① Data Lineage - Metadata provides info about the history of data helps user trace the lineage and understand concept.

#### \* Classification of Metadata

① Technical Metadata - Gives information regarding the technical details like format of data, storage, data model, data dictionary etc.

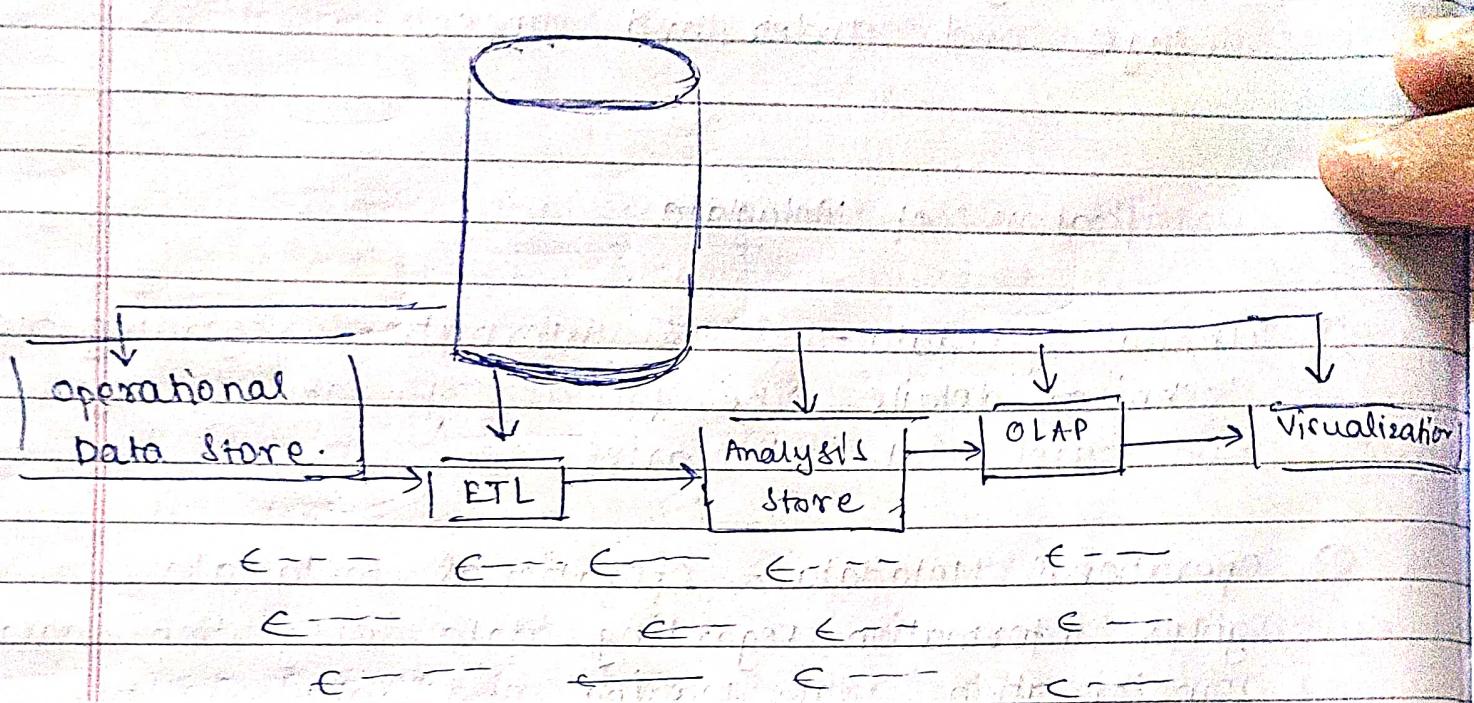
② Operational Metadata - Operational metadata gives information regarding data integration, data transformation, data loading and data quality.

③ Business Metadata - contains business concepts represented by data.

④ Usage Metadata - Contains info regarding query execution, report generation and data visualisation. used to analyse data.

⑤ Data Lineage Metadata - Describes the history of the data stored. transformation, date it came into the warehouse etc.

## \* Information Flow Mechanism and Architecture.



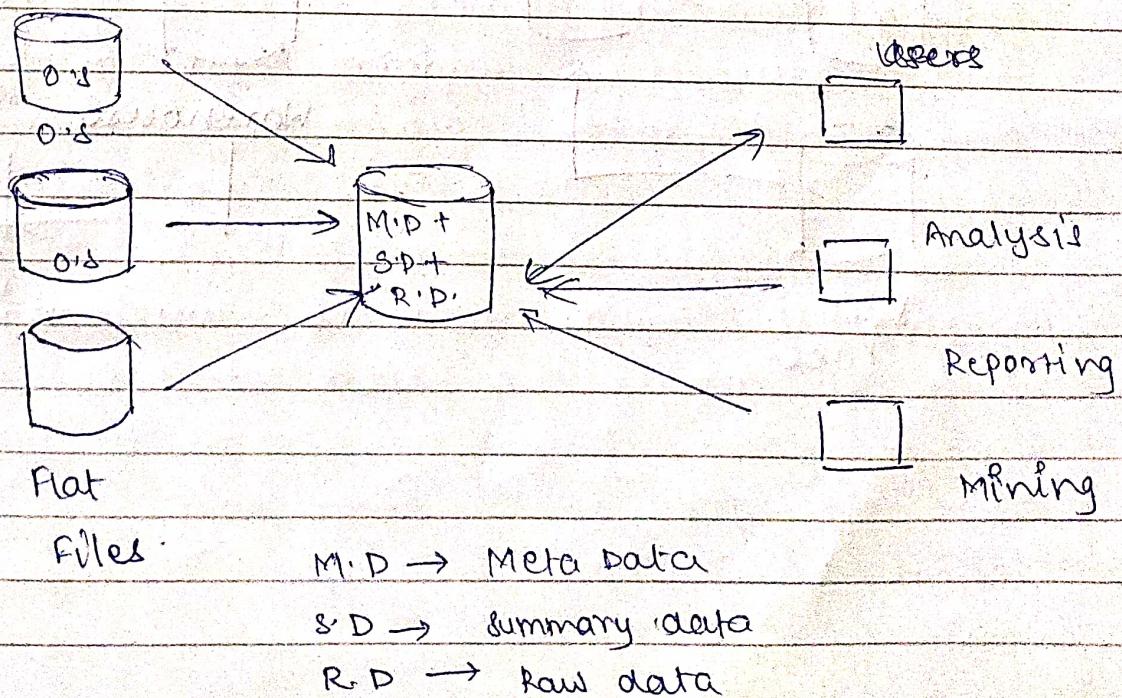
① Operational data store (ODS) - ODS is a database which stores the latest and real time data from various operating systems.

② ETL - contains 3 steps-

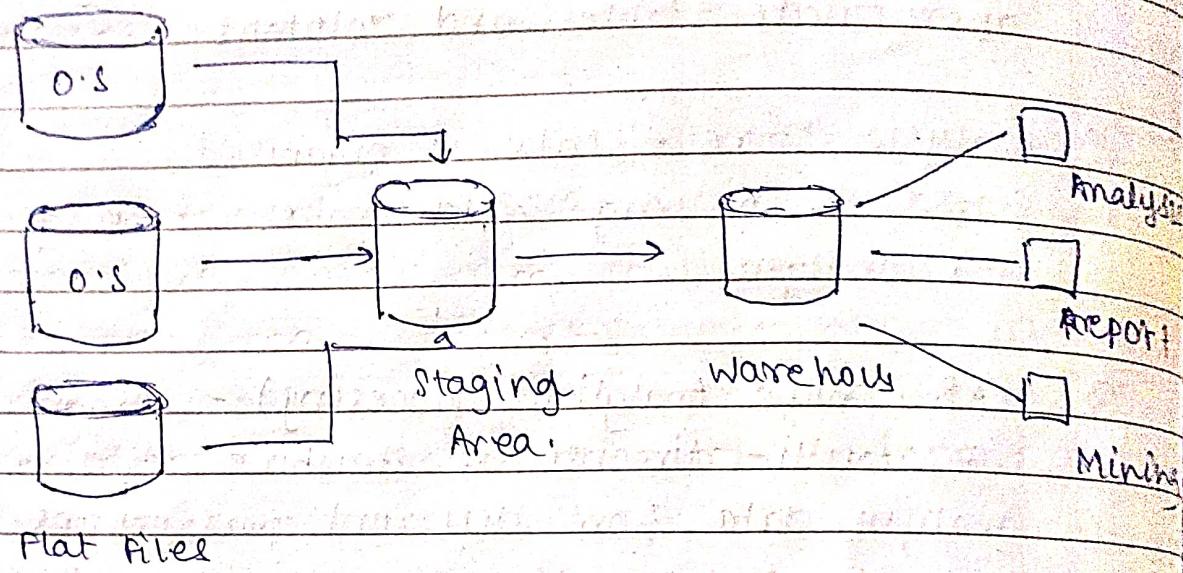
- a) Extract - Data is extracted from ODS and other source systems, such as CRM, ERP etc.
- b) Transform - The extracted data is then transformed into a format that is suitable for the data warehouse.

- c) Load - transformed data is then loaded onto the data warehouse. This involves putting data into appropriate tables and columns.
- ④ Analysis Store - Data is organised, structured and stored in a way which makes it easy to query and analyse.
- ⑤ OLAP (Online Analytical Processing) - organised data in a multi-dimensional structure. This allows users analyse data from different perspectives. OLAP provides tools to analyze trend patterns and anomalies.
- ⑥ Visualisation - creates reports, patterns, trends, insights, dashboard etc using raw data within the data warehouse.

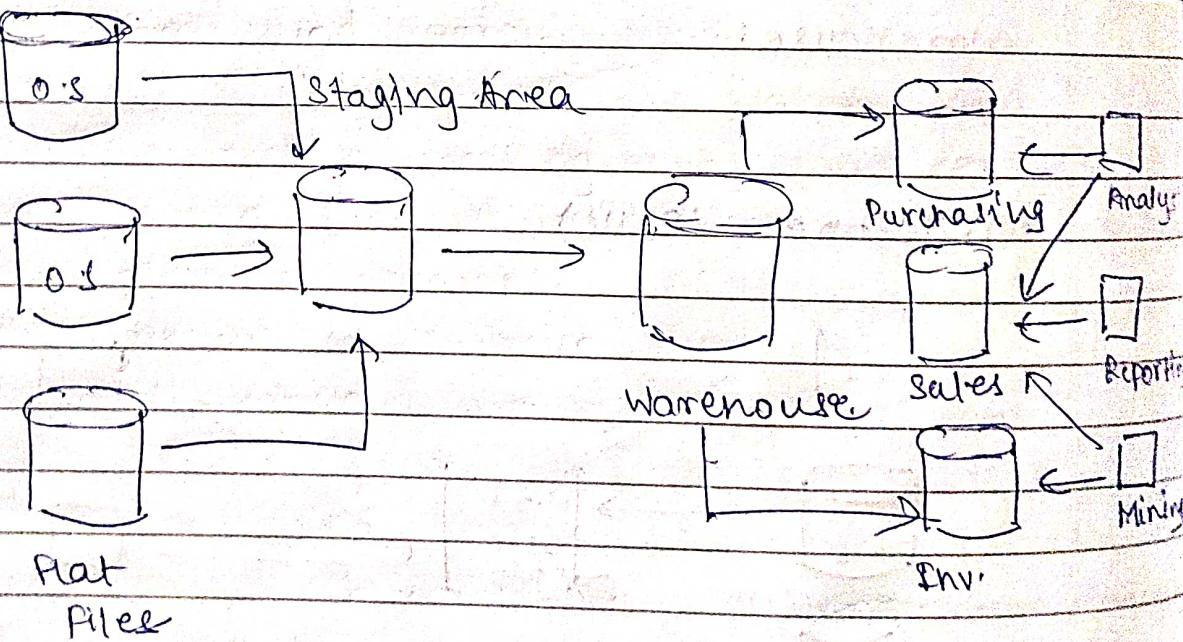
#### \* Architecture: (Basic)



## \* Architecture (with Staging Area)



## \* Architecture (Staging Area and Data Marts)



## Module 3 - (from 3.3 to 3.6)

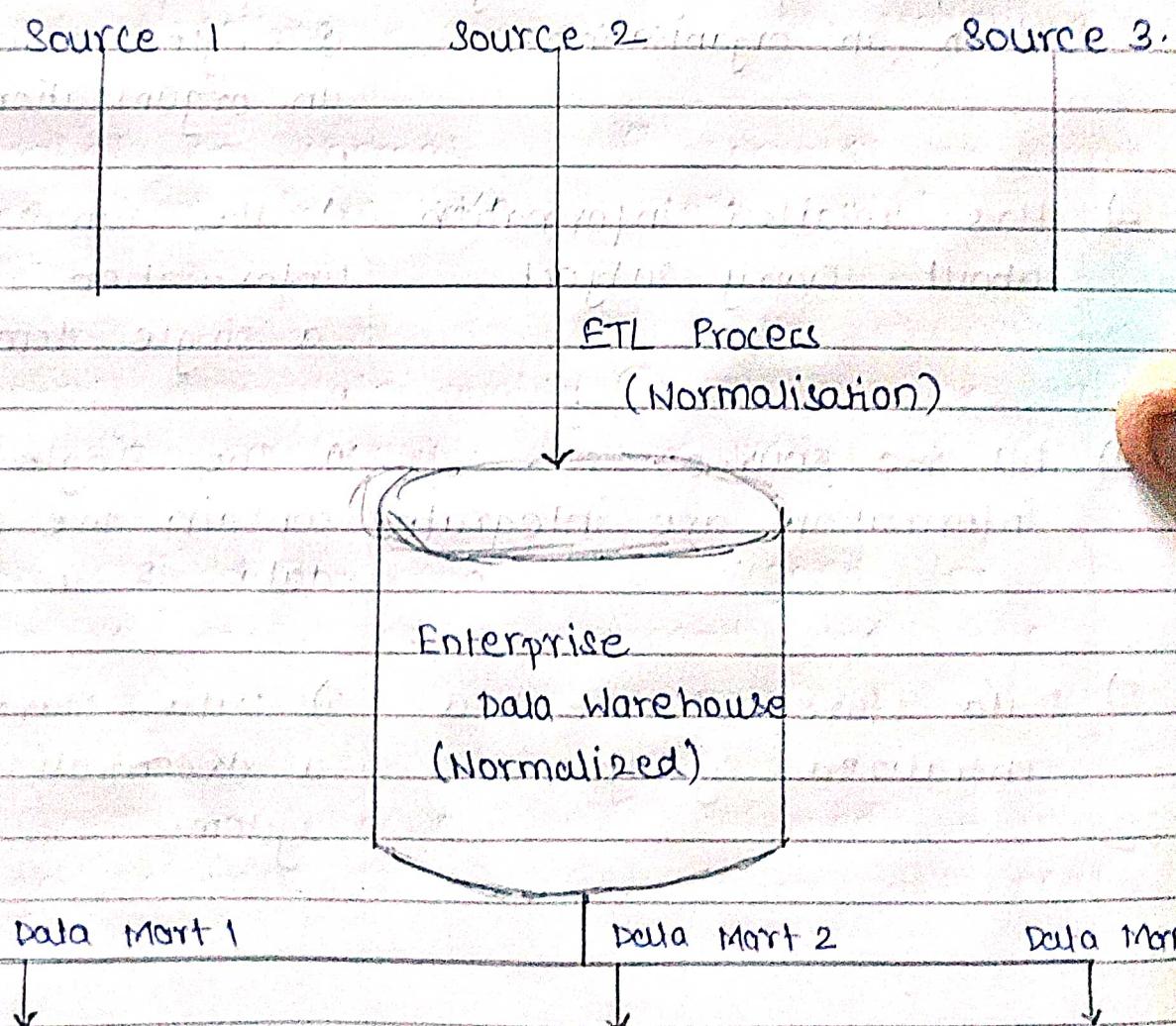
### \* Datawarehouse vs Data Marts

Data Warehouse	Data Mart
1) Data Warehouse have a vast repository of collected information for internal as well as external sources.	1) Data Mart is a sub category of data warehouse created to meet the requirements of a smaller sub-group.
2) Data Warehouse holds information related to multiple topics and subjects within an organisation.	2) Data Mart holds information related to a specific domain or a single subject within an organisation.
3) Has detailed information about every subject.	3) Has summarized information regarding a single domain.
4) All the sources of information are integrated.	4) The single source contain the specific data is integrated.
5) Data Warehouse has a centralized system.	5) Data Mart has a decentralized/distributed system.

## \* Data warehouse design strategies

- Data warehouse are used across multiple businesses to help with critical decision making and data analysis.
- It has two main approaches:

① Top- Down Approach: Also known as enterprise data warehouse (EDW), which involves creating a comprehensive, centralised data warehouse that serves an entire organisation.



② Bottom Advantages of Top-Down Approach.

→ Ensures data consistency and accuracy.

→ Developing new data marts is easier.

→ Disadvantages of Top-Down Approach.

→ Inflexible to change in departmental needs

→ cost is high.

③ Bottom-Up Approach :- Also known as the data-mart approach. This approach involves small, targeted data marts that serve specific business function and department.

Source 1                      Source 2                      Source 3

ETL Process

Process

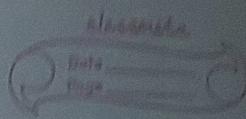
DM 1

DM 2

DM 3

Warehouse

(conformed facts)



### \* Advantages of Bottom Up Approach

- Documents are generated quickly.
- cheaper than Top-Down Approach.

### \* Disadvantages of Bottom Up Approach

- Lead to data inconsistency.
- difficult to manage and govern data.

## \* Data warehouse Modelling vs Operational Database Modelling

### DW Modelling

- ① Operates on schema like star, snowflake etc.
- ② Emphasis on data quality and data consistency
- ③ Designed for read intensive operations.
- ④ Designed to manage large volumes of data

### OD. Modelling

- ① Operates on normalized data model.
- ② Emphasis on data accuracy.
- ③ Designed for both read and write operations.
- ④ Designed to manage smaller but real-time data.

## \* Advantages of Bottom Up Approach

- Documents are generated quickly.
- Cheaper than Top-Down Approach.

## \* Disadvantages of Bottom-Up Approach

- Leads to data inconsistency.
- Difficult to manage and govern data.

## \* Data Warehouse Modelling vs Operational Database Modelling:

### DW Modelling

- ① Operates on schemas like star, snowflake etc.
- ② Emphasis on data quality and data consistency
- ③ Designed for read intensive operations.
- ④ Designed to manage large volumes of data.

### OD. Modelling

- ① Operates on normalized data model.
- ② Emphasis on data accuracy.
- ③ Designed for both read and write operations.
- ④ Designed to manage smaller but real-time data.

## \* Star Schema

→ In Star Schema, data is organized into a central fact table that contains the measures of interest, surrounded by dimension tables that describe the attributes of the measure.

Example:

		Product Dim.	
		Product ID	Product Name
		Product Cat.	Unit Price
Time Dimension			
Order ID			
Order Name			
Year	Sales		
Quarter	Product ID		
Month	Order ID		
	Customer ID	EMP Dimension	
	Employer ID	Emp ID	
	Total	Emp Name	
	Quantity	Title	
	Discount	Department	
Customer Dim		Region	
Customer ID			
Customer Name			
Address			
City			
Zip			

- In Star Schema, quantitative data is stored in fact tables and dimensions which are the context for the quantitative data is stored in dimension tables.
- Fact dimension table is joined to the fact table through a foreign key relationship. This allows query executions in fact table using attributes from dimension table.

#### \* Advantages of Star Schema.

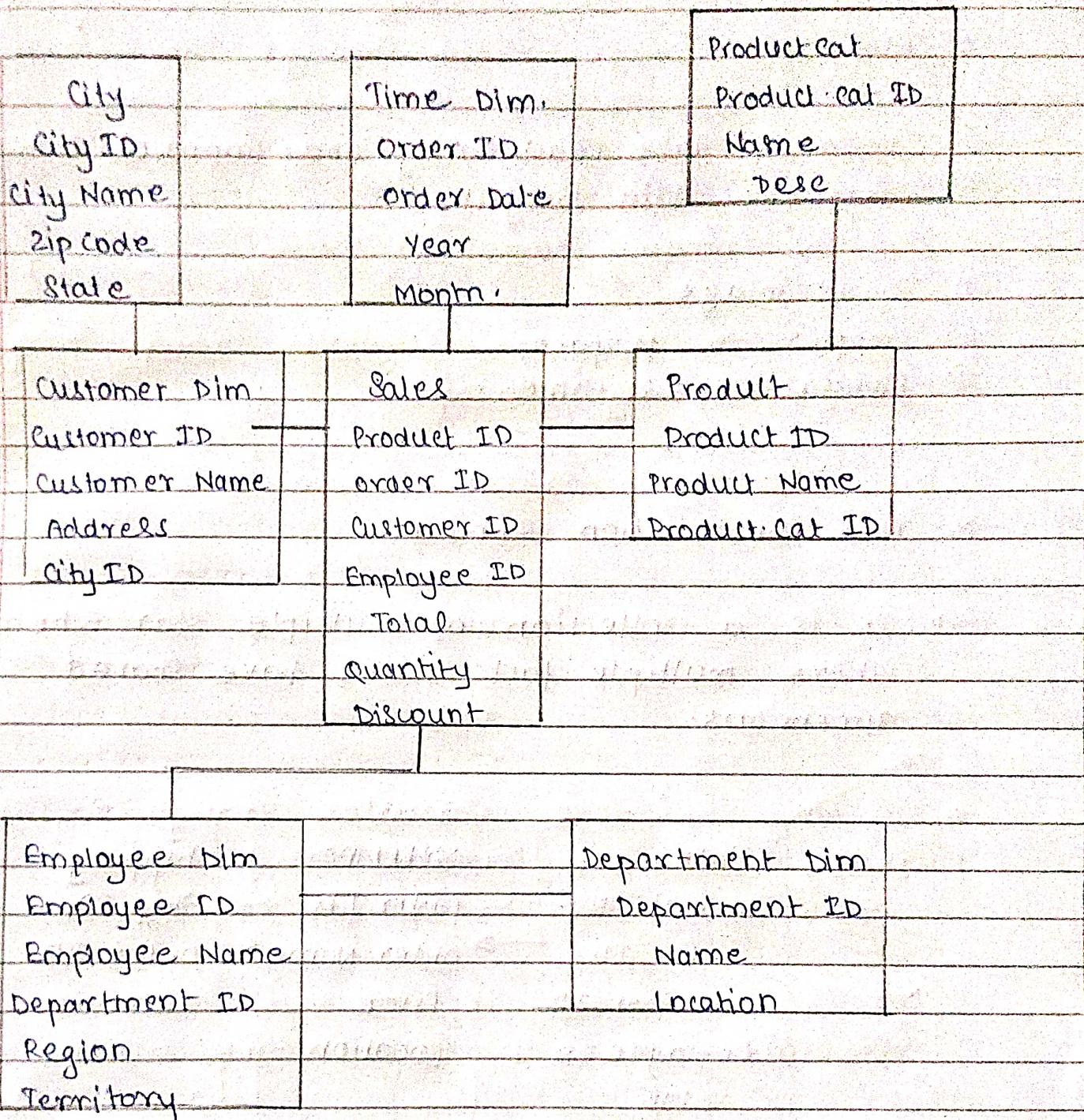
- ① Simpler queries.
- ② Simplified Business Reporting logic.
- ③ Feeding cubes in OLAP.

#### \* Disadvantages of Star Schema.

- ① Highly denormalized.
- ② Doesn't reinforce many-to-many relationships within business entities.

#### \* Snowflake Schema

- Snowflake Schema is the extended version of the Star Schema where the centralized fact table is connected to multiple dimensions. There are several levels of relationship and hierarchies.



→ The Fact table is located at the centre of the schema, surrounded by dimension tables. Each dimension table is broken down into smaller dimensions.

## \* Advantages

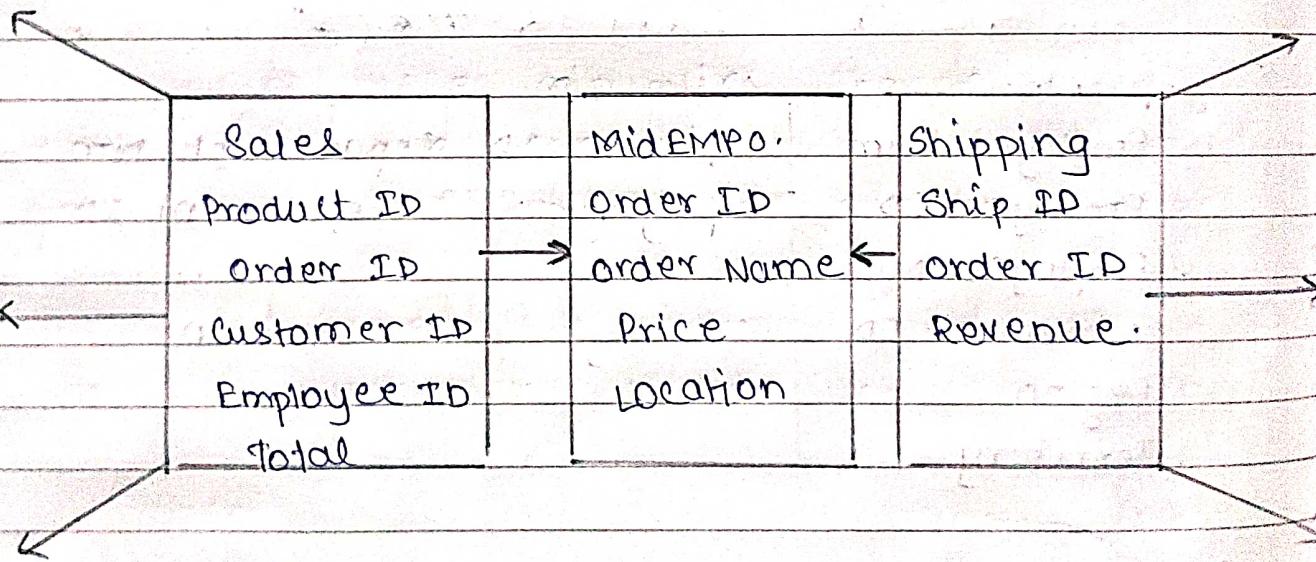
- ① Reduces data redundancy and improves consistency
- ② Improves data integrity.

## \* Disadvantage

- ① Complex to prepare.
- ② Slower response time.

## \* Fact constellation Schema

→ It is a collection of multiple star schema where multiple fact tables have shared dimensions.



### \* Factless - fact table.

→ Factless fact tables are ones with only foreign key column present while facts/measures are absent. Meaning they are only used to establish relations b/w 2 dimensions.

### \* slowly changing dimensions

→ As data constantly increases due to high transaction rates, the fact tables tend to change frequently but over time changes can occur in dimension tables due to change in their attributes.

i) Type 1 change: Overwriting existing data with new data having no record or log of the past and no changes except overwriting.

ii) Type 2 change: Tracking all the changes that are made by creating a new row for updated values.

iii) Type 3 change - The change in creating new fields in the dimension record to capture the last value for that attribute

#### \* Rapidly changing dimensions

→ When one or more attributes change frequently.  
e.g. BMI, height, weight.

For eg - Patient dimension

Pat. ID

Name

Gender

Weight

BMI

→ If any of these attributes change rapidly, many rows are created and dataset would get larger and larger for minute changes.

→ These attributes are removed and put in a separate dimension called Junk dimension.

Patient

Pat. ID

Name

Gender

↑

dimension

table

Patient JNK

Pat. SK

BMI

index.

- Pat SK acts as a surrogate key in junk dimension which is the primary key.
- \* Data Lake.
- A Data Lake is a large, centralised repository that allows you to store and manage vast amounts of raw data in its native format. Unlike traditional DW, Data Lake can store data in semi-structured, unstructured and structured format.
- Architecture

