

DW-Data-Modelling-Refresher-Notes

- DW-Data-Modelling-Refresher-Notes
 - 1. OLAP vs. OLTP
 - 2. Facts
 - 3. Dimensions
 - 4. Star Schema
 - Characteristics of a Star Schema
 - Advantages of Star Schema
 - Disadvantages of Star Schema
 - Conclusion
 - 5. Snowflake Schema
 - Characteristics of a Snowflake Schema
 - Advantages of Snowflake Schema
 - Disadvantages of Snowflake Schema
 - 6. Star Schema vs. Snowflake Schema – Which is Better?
 - 7. Fact Table and Its Types
 - 8. Granularity
 - 9. Data Virtualization
 - Think of it like a TV remote!
 - 10. Data Marts
 - Think of it like a small store inside a shopping mall!
 - 11. Data Lakes
 - Think of it like a big lake!
 - 12. Relational OLAP
 - Think of it like a Library with an Index!
 - 13. Multi dimensional OLAP
 - Think of it like a Vending Machine!
 - 14. Data Cubes
 - 1. Slice (Filtering one dimension)
 - 2. Dice (Filtering multiple dimensions)

- 3. Drill Down (Going deeper into details)
- 4. Roll Up (Summarizing the data)
- 5. Pivot (Rotating the view)
- Summary Table
- 15. Dimensional Data
- 16. Difference between table and cube data type
 - **1** Table Data Type (Relational Data Model)
 - ♦ Key Characteristics of Tables:
 - **2** Cube Data Type (Multidimensional Model)
 - ♦ Key Characteristics of Cubes:
 - 📦 Example of a Sales Cube (Multidimensional Data)
 - **3** Key Differences Between Table & Cube Data
- 17. Centipede Fact table
 - 🛠️ Characteristics of a Centipede Fact Table
 - 📌 Example: Centipede Fact Table in an E-commerce Business
 - Linked Dimension Tables
 - 📌 When to Use a Centipede Fact Table?
- 18. Factless Fact Table
 - **1** Why Use a Factless Fact Table?
 - 📌 Example: Event Tracking
 - **2** When to Use a Factless Fact Table?
- 19. Three-Tier Data warehouses
 - Three-Tier Data Warehouse Architecture
 - **1** Bottom Tier: Data Source Layer (Raw Data Storage)
 - **2** Middle Tier: Processing Layer (ETL & Data Storage)
 - **3** Top Tier: Presentation & Reporting Layer
 - 🔄 Summary Table: Three-Tier Data Warehouse
- 20. Junk dimension
 - **1** Why Do We Need a Junk Dimension?
 - 🛑 Without a Junk Dimension (Messy Fact Table)
 - **2** How Junk Dimension Solves This?
 - ✅ With a Junk Dimension (Clean Fact Table)
 - Junk Dimension Table (📦 Groups Unrelated Attributes Together)

- Fact Table (Now Cleaner)
- **3** When to Use a Junk Dimension?
- 21. ODS and its types
- 22. Major Parallel computing platforms
- 23. Quantifiable Data Indicator
 - Examples of Quantifiable Data Indicators
 - QDI vs. Qualitative Data Indicator
- 24. Cross Tab
 - Example of Cross Tab Analysis
- 25. Metadata in datawarehouse
- 26. Conformed Dimensions
- 26. Domain
- 27. Hierarchy in Data Warehousing
- 28. Role playing dimension
 - Example: Order Processing
 - Table Structure:
- 29. Degenerate dimension
 - Example: Order Transactions
 - Table Structure:
- 30. Data Lake vs Data Warehouse
- 31. In memory databases
- 32. Schematic Evolution
- 33. Dimensional Modelling
 - Dimensional modelling
- 34. Conceptual Data Model
- 35. Physical Data Model
- 36. Surrogate Key
- 37. ERWIN (Data Modeling Tool)
- 38. Softwares Used for Data Modeling
- 39. Data modelling types
 - **1** Conceptual Data Model
 - **2** Logical Data Model
 - **3** Physical Data Model

- Comparison Table
- 40. Examples of star and snowflake schema
 - Star Schema — (Simpler, Flattened)
 - What is it?
 - Example:
 - Snowflake Schema — (Normalized, More Tables)
 - What is it?
 - Example:
 - Summary Table:
- 41. Why Transformation is Needed?
- 42. ETL vs. ELT
- 43. ELT Procedures
- 44. SCDs (Slowly Changing Dimensions)
- 45. Types of SCDs
 - SCD Type 0 (Retain Original) – No Change
 - SCD Type 1 (Overwrite) – Keep Only Latest Data
 - SCD Type 2 (Versioning) – Maintain History with New Rows
 - SCD Type 3 (Add Column) – Keep Only One Previous Value
 - SCD Type 4 (Separate History Table) – Archive Old Data
 - SCD Type 5 (Hybrid – Type 1 + Type 4)
 - SCD Type 6 (Hybrid – Type 1 + Type 2 + Type 3)
- 46. What is ETL
 - Step 1: Extraction
 - Step 2: Transformation
 - Step 3: Loading
- 47. What is Change Data Capture

1. OLAP vs. OLTP

- ♦ **OLTP (Online Transaction Processing)** – Handles day-to-day transactions.
- ♦ **OLAP (Online Analytical Processing)** – Used for complex queries and reporting.

Feature	OLTP (e.g., Banking App)	OLAP (e.g., Data Warehouse)
Purpose	Fast transactions	Data analysis & reporting

Feature	OLTP (e.g., Banking App)	OLAP (e.g., Data Warehouse)
Data Size	Small (few MBs/GBs)	Huge (TBs/PBs)
Operations	Insert, Update, Delete	Read-heavy (Aggregation, Grouping)
Example	ATM withdrawal	Monthly sales report

Example:

- A shopping website (**OLTP**) records purchases in a database.
- A sales manager checks **yearly sales trends** using an OLAP system.



2. Facts

A **fact** is a measurable event stored in a data warehouse.

Example:

For an **e-commerce store**, facts could be:

- ✓ **Total Sales**
- ✓ **Number of Orders**
- ✓ **Revenue Generated**

Facts are usually stored in **Fact Tables**.



3. Dimensions

A **dimension** provides context to facts. It describes **who, what, where, and when**.

Example:

For sales data, dimensions could be:

- ✓ **Customer** (Name, Age, Location)
- ✓ **Product** (Category, Price)
- ✓ **Time** (Year, Quarter, Month)

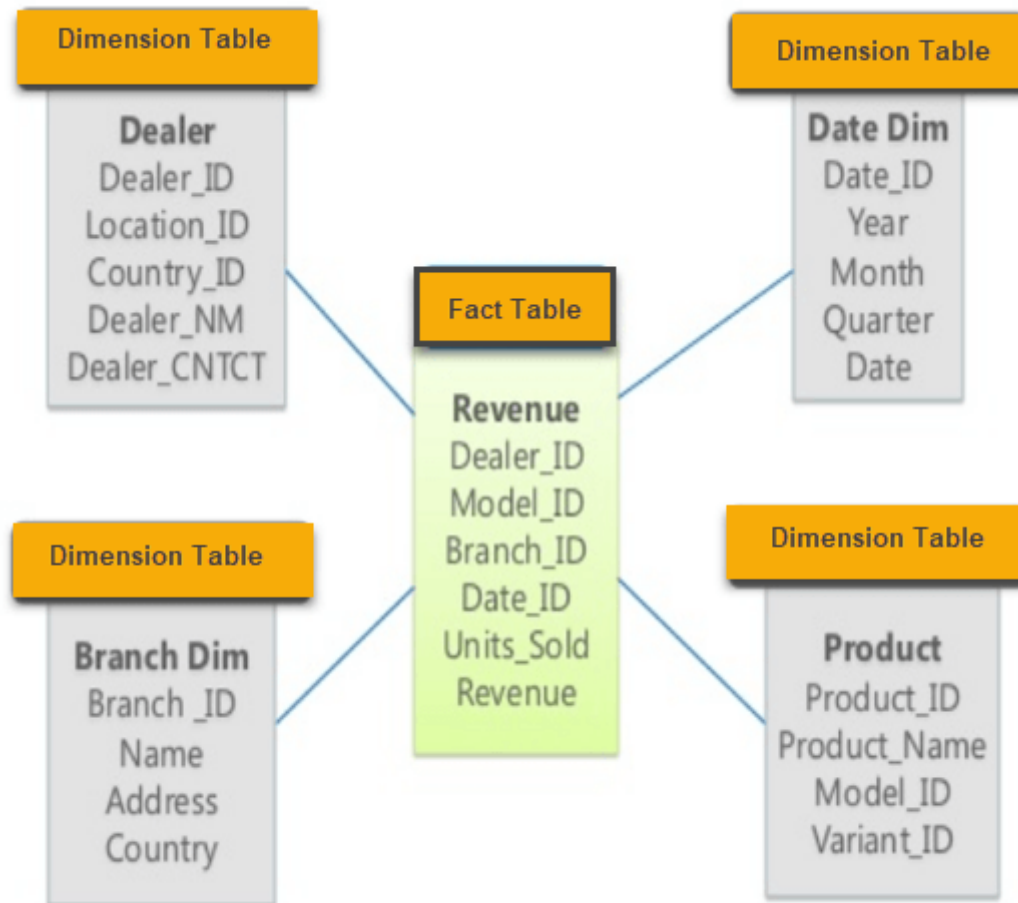
Dimensions are stored in **Dimension Tables** and linked to fact tables.



4. Star Schema

A **Star Schema** organizes data into:

- ✓ **A central Fact Table** (stores measurable data)
- ✓ **Multiple Dimension Tables** (descriptive details)



Example of Star Schema Diagram

Characteristics of a Star Schema

- ✓ **Single-level dimension tables** (no normalization).
- ✓ **Dimension tables connect only to the fact table** (not to each other).
- ✓ **Faster query performance** due to simple joins.
- ✓ **Easy to maintain and expand.**

Advantages of Star Schema

- ✓ **Faster Query Performance** – Simple structure enables efficient indexing and aggregation.
- ✓ **Easier for Reporting & Analysis** – Widely used in **BI tools** like Power BI, Tableau.

✔ **Supports OLAP Cubes** – Star schema helps in creating **multi-dimensional data cubes**.



Disadvantages of Star Schema

✗ **Data Redundancy** – Since dimensions are not normalized, some data gets repeated.

✗ **Not Ideal for Complex Relationships** – Cannot efficiently handle **many-to-many** relationships.

✗ **Limited Flexibility** – Adding new analysis fields might require schema redesign.

Conclusion

📌 **Star Schema** is the most popular **data warehouse model**, used in OLAP and reporting systems.

📌 It consists of **one fact table** and **multiple dimension tables** forming a star-like structure.

📌 Provides **fast querying**, **easy maintenance**, and **efficient OLAP cube design**.

🚀 **Ideal for:** Business Intelligence (BI), Data Warehousing, and Decision Support Systems.



5. Snowflake Schema

A **Snowflake Schema** is a more structured version of the Star Schema.

📌 **Example:**

Instead of storing a **Location** directly in the Customer table, we split it into:

✔ **Customer** → CustomerID, Name, LocationID

✔ **Location** → LocationID, City, Country

Characteristics of a Snowflake Schema

✔ **Dimension tables are normalized** into multiple sub-tables.

✔ **Minimizes redundancy** by reducing duplicate data storage.

✔ **More complex queries** due to multiple joins.

✔ **Better for hierarchical data** (e.g., Country → State → City).

Advantages of Snowflake Schema

- ✓ **Saves Disk Space** – Normalization reduces storage requirements.
- ✓ **No Data Redundancy** – Each value is stored only once.
- ✓ **Scalability** – Works well with large, complex datasets.

Disadvantages of Snowflake Schema

- ✗ **Slower Query Performance** – More joins lead to **longer execution times**.
- ✗ **Complex to Maintain** – Schema changes require updating multiple related tables.
- ✗ **Harder to Understand** – More complex relationships make querying difficult.



6. Star Schema vs. Snowflake Schema – Which is Better?

Feature	Star Schema	Snowflake Schema
Storage	Takes more space	Uses less space
Query Speed	Faster	Slower (More Joins)
Simplicity	Easy to understand	More complex

📌 Example:

- **Star Schema** is better for fast queries (e.g., sales reports).
- **Snowflake Schema** is better for **optimized storage**.



7. Fact Table and Its Types

A **Fact Table** stores numerical data (sales, revenue, etc.).

Types:

- ✓ **Transaction Fact Table** – Stores business events (e.g., a purchase)
- ✓ **Snapshot Fact Table** – Stores periodic data (e.g., monthly balance)
- ✓ **Accumulating Fact Table** – Tracks the status of a process (e.g., order lifecycle)



8. Granularity

- ♦ **Granularity** refers to the **level of detail** in a data warehouse, particularly in **fact tables**.

📌 **Example:**

- **High Granularity** → Data is very detailed (e.g., per second sales).
- **Low Granularity** → Data is summarized (e.g., monthly sales).

Choosing the right granularity affects **storage size and query speed**.



****9. Data Virtualization**

👉 **Access data from multiple sources without actually moving or copying it.**

Think of it like a TV remote!

Just like a TV remote lets you switch between different channels without moving to a different room, **data virtualization** allows users to access and query data from different databases without physically transferring it.

♦ **Example:**

A company has data stored in:

- MySQL (customer data)
- PostgreSQL (sales data)
- Google Sheets (marketing data)

Instead of copying all the data into one database, **data virtualization** lets employees **query all sources at once** without moving the data.

✅ **Pros:** No duplication, real-time access

❌ **Cons:** May be slower for big queries



10. Data Marts




👉 **A small, focused version of a data warehouse for a specific department.**

Think of it like a small store inside a shopping mall!

A **data warehouse** is like a giant mall with everything. A **data mart** is like a small store inside that mall, serving a specific group of people.







♦ **Example:**

A company has a large data warehouse, but different departments only need specific data:

-  **Sales Team:** Needs sales performance data
-  **Marketing Team:** Needs customer engagement data
-  **Finance Team:** Needs revenue and expense data

Instead of each team searching through the entire data warehouse, they get their own **data mart** with just the relevant information.

- ✓ **Pros:** Faster, easier to manage for teams
- ✗ **Cons:** Needs proper setup and maintenance

Type	How It's Built 	Where It Gets Data 	Best For 
1 Dependent Data Mart 	Built from a central data warehouse	From the main data warehouse	Large organizations that already have a data warehouse
2 Independent Data Mart 	Created separately, without a data warehouse	Directly from operational databases	Smaller businesses or when a data warehouse doesn't exist
3 Hybrid Data Mart 	Combination of Dependent & Independent	From both data warehouses & external sources	Companies that need a mix of internal & external data



11. Data Lakes




 A huge storage system that keeps raw data in its original format until needed.

Think of it like a big lake!

Just like a lake holds water from different sources (rivers, rain, streams) without filtering, a **data lake** stores all types of data—structured (tables), semi-structured (JSON, XML), and unstructured (videos, images).

♦ **Example:**

An e-commerce company stores:

-  Transaction data (structured)
-  Customer emails (semi-structured)
-  Product demo videos (unstructured)

A **data lake** keeps all this data **as-is** until it's needed for analysis.

✓ **Pros:** Flexible, stores all data types

✗ **Cons:** Can become messy (a "data swamp") if not managed well

12. Relational OLAP

👉 **Uses traditional relational databases (SQL) to store and analyze data.**

Think of it like a Library with an Index!

A library has books (data) stored in different sections (tables), and an **index** (SQL queries) helps find the information quickly.

♦ **How it works:**

- Stores data in relational databases (e.g., MySQL, PostgreSQL).
- Uses complex SQL queries to analyze large datasets.
- Works well with large and dynamic data



13. Multi dimensional OLAP

👉 **Uses pre-built cubes for fast analysis.**

Think of it like a Vending Machine!

Instead of searching for ingredients (raw data) and cooking, you get a ready-to-eat snack (pre-calculated data).

♦ **How it works:**

- Stores data in a **multidimensional cube** instead of tables.

- Pre-calculates summaries for quick access.
- Best for fast reporting and interactive analysis.

✅ **Pros:** Super fast queries since data is pre-aggregated.

❌ **Cons:** Takes up more storage space and is less flexible.

14. Data Cubes

A **data cube** is a multi-dimensional structure used for analyzing data efficiently, especially in **OLAP (Online Analytical Processing)**. It allows users to explore data from different perspectives, just like how we can look at a physical cube from different angles.

Now, let's break down the **key operations** used to analyze a data cube:

1. Slice (Filtering one dimension)

👉 Think of it as cutting a single slice from a cake.

- It selects a single value from **one** dimension and shows the remaining data.
- Example:
 - A data cube contains **sales data** with three dimensions: **Region, Product, and Time**.
 - If we take **only the sales data for "2024"** (fixing the "Time" dimension), we get a **slice** of the cube showing only "Region vs. Product" sales.

2. Dice (Filtering multiple dimensions)

👉 Like cutting a smaller block out of the cake.

- It selects specific values from **two or more** dimensions.
- Example:
 - From the same cube, we want **sales data for "2024" and only for "Mobile Phones" in "Asia"**.
 - This creates a **smaller sub-cube** focusing on those conditions.

3. Drill Down (Going deeper into details)

👉 Like zooming in on a map.

- It increases the level of detail by moving from **higher-level summaries** to **lower-level details**.

- Example:
 - We have **total sales per country** in a region.
 - If we "drill down" into India, we now see **sales per state** instead of just the country.

4. Roll Up (Summarizing the data)

👉 Like zooming out on a map.

- It decreases the level of detail by **aggregating** data into higher-level summaries.
- Example:
 - Instead of **sales per city**, we roll it up to show **total sales per country**.
 - This is the reverse of **Drill Down**.

5. Pivot (Rotating the view)

👉 Like turning a Rubik's cube to see a different side.

- It changes the way data is presented by **swapping rows and columns**.
- Example:
 - Suppose our table shows **Products in rows** and **Regions in columns**.
 - If we **pivot**, it can now show **Regions in rows** and **Products in columns**, giving a different perspective.

Summary Table

Operation	Action	Example
Slice	Filters one dimension	Show sales data only for 2024
Dice	Filters multiple dimensions	Show sales only for "Mobile Phones" in "Asia" during 2024
Drill Down	Increases detail	See sales per state instead of country
Roll Up	Summarizes	See sales per region instead of state
Pivot	Changes view	Swap "Regions" and "Products" in a table



15. Dimensional Data

Dimensional data is a way of organizing data to make it easier to analyze in a **data warehouse**. It focuses on how business users **query data** rather than how databases store it.

When data is structured for **reporting and analysis**, it is stored in two main types of tables:

1. **Fact Tables** (store business events & numbers)
2. **Dimension Tables** (store descriptive information)

16. Difference between table and cube data type

1 Table Data Type (Relational Data Model)

A **table** is a structured collection of rows and columns, like a spreadsheet. It follows the **relational database model** and is commonly used in databases like **PostgreSQL, MySQL, and SQL Server**.

♦ Key Characteristics of Tables:

- ✓ **Two-dimensional** (rows & columns).
- ✓ Stores **transactional data** in a normalized format.
- ✓ Uses **SQL queries** (JOIN, SELECT, WHERE, etc.).

Tables are **good for transactional systems (OLTP)** but **slow for complex analytics**, especially when dealing with large datasets.



2 Cube Data Type (Multidimensional Model)

A **cube** is a **multidimensional data structure** used in OLAP for fast data analysis. Instead of rows & columns, data is stored in multiple **dimensions**.

♦ Key Characteristics of Cubes:

- ✓ **Multi-dimensional** (more than just rows & columns).
- ✓ Optimized for **aggregated reporting and analytics**.
- ✓ Supports **OLAP operations** (Slice, Dice, Drill Down, Roll Up, Pivot).
- ✓ Pre-aggregated data for fast querying.

Example of a Sales Cube (Multidimensional Data)

A **Sales Cube** could have three dimensions:

- Time (Year, Month, Day)
- Product (iPhone, MacBook, iPad)
- Region (USA, India, Europe)


3 Key Differences Between Table & Cube Data

Feature	Table (Relational)	Cube (Multidimensional)
Structure	Rows & Columns (2D)	Multiple Dimensions (3D+)
Use Case	Storing transactional data	Fast analytical queries
Speed	Slower for complex queries	Faster due to pre-aggregated data
Operations	SQL (JOIN, WHERE, GROUP BY)	OLAP (Slice, Dice, Drill, Roll Up)
Optimization	For CRUD operations	For Business Intelligence (BI)

- Use **tables** when working with **transactional databases** (OLTP) for **storage and retrieval**.
- Use **cubes** when performing **fast analytical queries** on large datasets (OLAP).



17. Centipede Fact table

A **Centipede Fact Table** is a special type of fact table in a data warehouse that has **many foreign keys** pointing to multiple dimension tables. This structure makes it look like a centipede  because of its **long list of foreign key columns** acting like "legs."

Characteristics of a Centipede Fact Table

1. Many Dimension Tables

- A centipede fact table has a **large number of dimension tables** linked to it.
- Each dimension provides more details about the facts.

2. Wide Table (Many Foreign Keys)

- The fact table has **many foreign key columns**, making it **wide**.
- Each row in the table represents a business event (e.g., a sale, transaction, or shipment).

3. Optimized for Analysis, Not Performance

- While centipede fact tables provide **detailed** information, they can be **slower** because of many **joins** needed for querying.



Example: Centipede Fact Table in an E-commerce Business

Imagine an **E-commerce Sales Fact Table** with many dimension tables:

Order ID	Date Key	Customer Key	Product Key	Region Key	Payment Key	Shipment Key	Sales Amount	
1001	20240301	501	101	301	701	801	₹5000	:
1002	20240302	502	102	302	702	802	₹3000	:

Linked Dimension Tables

- **Date Dimension** (Date Key → Year, Month, Day)
- **Customer Dimension** (Customer Key → Name, Age, Location)
- **Product Dimension** (Product Key → Name, Category, Brand)
- **Region Dimension** (Region Key → Country, State, City)
- **Payment Dimension** (Payment Key → Payment Type, Bank Name)
- **Shipment Dimension** (Shipment Key → Courier, Delivery Time)

Because there are **so many foreign keys**, the fact table **looks like a centipede** with many "legs" (references to dimension tables).

When to Use a Centipede Fact Table?

- ✓ When you need **highly detailed data** for in-depth analytics.
- ✓ When your data warehouse supports **many dimensions** for better reporting.
- ✗ Avoid if performance is a concern—consider **denormalization** to reduce joins.



18. Factless Fact Table

A **Factless Fact Table** is a fact table that **does not have measurable numeric facts** (like sales amount, revenue, or quantity). Instead, it records **events or relationships** between

dimensions without any numerical metrics.

📌 Think of it as a fact table that captures "what happened" but without actual values.

1 Why Use a Factless Fact Table?

- ♦ To track **events or occurrences** (e.g., student attendance, product promotions).
- ♦ Useful for **analyzing trends and patterns** without needing numerical data.

📌 Example: Event Tracking

- Captures **something that happened** but has no numerical measures.
- Example: **Student Attendance** in a university.

Fact Table: Attendance Fact

Date Key	Student Key	Course Key	Professor Key
20240301	101	201	301
20240302	102	202	302

What this tells us?

- A student (101) attended a course (201) on a specific date (20240301).
- There's no measurable fact like "hours attended"—just the **event itself**.

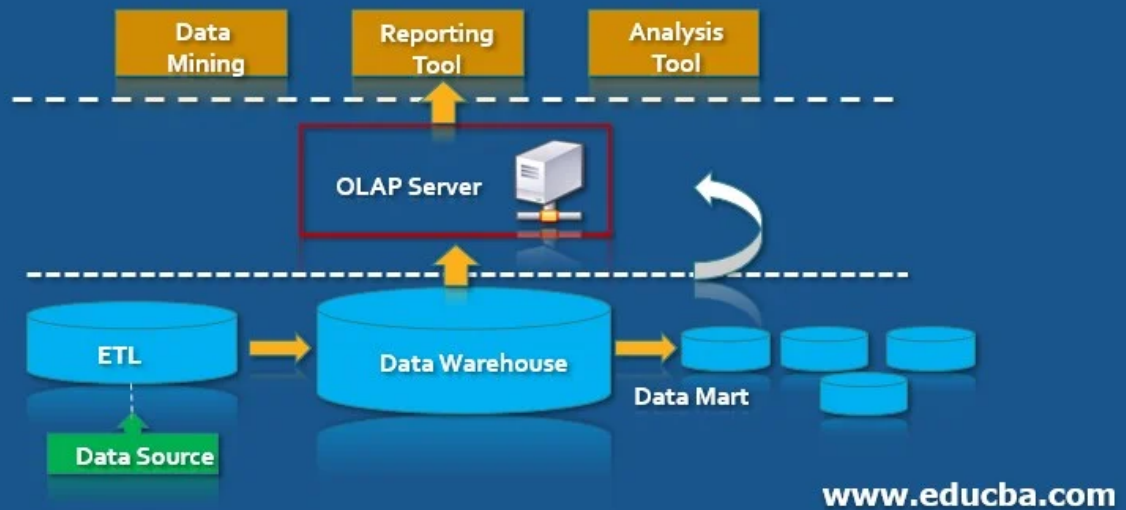
2 When to Use a Factless Fact Table?

- ✓ When tracking **events or actions** without measurable facts.
- ✓ When modeling **relationships** between dimensions (many-to-many).
- ✓ When analyzing **patterns, trends, or coverage data**.



19. Three-Tier Data warehouses

Three Tier Data Warehouse Architecture



Three-Tier Data Warehouse Architecture

A **Three-Tier Data Warehouse** is a structured way of organizing a data warehouse, dividing it into three layers:

- 1 **Bottom Tier (Data Source Layer)** – Where data is collected
- 2 **Middle Tier (Processing Layer)** – Where data is transformed and stored
- 3 **Top Tier (Presentation Layer)** – Where users access the data

📌 **Think of it like a restaurant:**

- **Kitchen (Bottom Tier)** → Raw ingredients (data) are stored.
- **Chef (Middle Tier)** → Prepares and organizes the food (transforms data).
- **Waiter & Menu (Top Tier)** → Serves the final dishes (reports & dashboards).

1 Bottom Tier: Data Source Layer (Raw Data Storage)

- ♦ This is where the data warehouse gets **raw data** from multiple sources.

📌 **Example:** A retail company collects data from **sales, inventory, and customer transactions** in various databases.



2 Middle Tier: Processing Layer (ETL & Data Storage)

- ◆ This layer does:
 - **ETL (Extract, Transform, Load)** → Cleans, formats, and loads data into the warehouse.
 - **Data Storage** → Stores processed data in **Star Schema** or **Snowflake Schema**.
 - **OLAP Processing** → Enables multidimensional data analysis (cubes, aggregations).

📌 Example:

- Sales data from different stores is cleaned and converted into **fact and dimension tables**.
- The data warehouse organizes the data for fast querying.



3 Top Tier: Presentation & Reporting Layer

- ◆ This is where users **query and analyze** the data.
- ◆ Includes:
 - **BI Tools & Dashboards** (Power BI, Tableau, Looker, etc.).
 - **SQL Query Interfaces** (for data analysts).
 - **Reports & Visualizations** (for business users).

📌 Example:

- A sales manager runs a **dashboard report** to see monthly revenue trends.
- A data analyst runs an **SQL query** to find the top-selling products.

🔄 Summary Table: Three-Tier Data Warehouse

Tier	Purpose	Example Technologies
Bottom Tier	Collect & store raw data	PostgreSQL, MySQL, CSV files
Middle Tier	Process & organize data (ETL, OLAP)	Apache Spark, Snowflake, Amazon Redshift
Top Tier	User access (reports, dashboards)	Power BI, Tableau, Looker, SQL Queries



20. Junk dimension

A **Junk Dimension** is a dimension in a data warehouse that combines **low-cardinality** attributes (flags, indicators, or miscellaneous data) into a single table.

📌 **Think of it as a "dumping ground" for small, unrelated attributes** that don't fit well in other dimensions.

1 Why Do We Need a Junk Dimension?

- To **reduce clutter** in the fact table (avoid too many foreign keys).
- To **improve performance** by grouping unrelated attributes together.
- To **reduce storage space** by avoiding multiple small dimension tables.

🔴 Without a Junk Dimension (Messy Fact Table)

Order ID	Product Key	Customer Key	Discount Applied?	Promo Code Used?	Return Status	Gift Wrapped?
1001	101	501	Yes	BLACKFRIDAY	No	Yes
1002	102	502	No	-	Yes	No

❌ Problems:

- Too many small attributes in the fact table.
- Increases the number of **columns** and makes queries slower.
- Harder to manage flags and indicators.



2 How Junk Dimension Solves This?

Instead of keeping all flags and indicators in the fact table, we **combine them into one Junk Dimension** table.

✅ With a Junk Dimension (Clean Fact Table)

Junk Dimension Table (📦 Groups Unrelated Attributes Together)

Junk Dimension Key	Discount Applied?	Promo Code Used?	Return Status	Gift Wrapped?
1	Yes	BLACKFRIDAY	No	Yes
2	No	-	Yes	No

Fact Table (Now Cleaner)

Order ID	Product Key	Customer Key	Junk Dimension Key
1001	101	501	1
1002	102	502	2

✓ Benefits:

- Fewer columns in the fact table.
- Faster queries due to indexing on **Junk Dimension Key**.
- Easier to maintain and extend in the future.

3 When to Use a Junk Dimension?

- ✓ When you have **small attributes** (flags, indicators, status fields) that don't belong to any specific dimension.
- ✓ When you want to **optimize storage** and **improve performance**.
- ✓ When attributes **don't change frequently** (if they do, consider a Slowly Changing Dimension instead).



21. ODS and its types

An **Operational Data Store (ODS)** is a type of database that **stores current and integrated data** from multiple sources for **real-time or near-real-time reporting**.

📌 **Think of it as a “mini-data warehouse”** that provides fresh data for operational decision-making.

Classification of ODS

Classifications	Comment
Class 1	<ul style="list-style-type: none">• Refresh cycle: Real-time• Degree of transformation: Low due to compressed timeframes
Class 2	<ul style="list-style-type: none">• Refresh cycle: ½ - 1 hour store and forward mechanism• Degree of transformation: Medium
Class 3	<ul style="list-style-type: none">• Refresh cycle: Daily. Traditional batch process• Degree of transformation: High
Class 4	<ul style="list-style-type: none">• Refresh cycle: Ad hoc, often involving preprocessed, value-added information from a data warehouse• Degree of transformation: Highest

- **Class I** – Updates of data from operational systems to ODS are synchronous.
- **Class II** – Updates between operational environment & ODS occurs between 2-3 hour frame.
- **Class III** – synchronization of updates occurs overnight.
- **Class IV** – Updates into the ODS from the DW are unscheduled.



22. Major Parallel computing platforms

1 Cluster Computing – A group of connected computers that work together as a single system. They are located in the same place and connected through a high-speed network. Used for AI, simulations, and databases.

Example: Google's data centers, Hadoop clusters.

2 Grid Computing – A collection of computers spread across different locations, working together on a shared problem. Each computer can have different hardware and operating systems. Used in scientific research and distributed computing.

Example: CERN's computing grid, Folding@home.

3 Massively Parallel Processing (MPP) – A system where many processors work independently on different parts of a large problem. Each processor has its own memory and

communicates with others through a network. Used in big data analytics and data warehousing.

Example: Snowflake, Amazon Redshift, Teradata.

4 High-Performance Computing (HPC) – Using supercomputers or large clusters to solve highly complex problems at extreme speeds. Requires specialized processors and high-speed networking. Used in scientific simulations, weather forecasting, and AI training.

Example: Summit Supercomputer, Fugaku, Cray systems.



23. Quantifiable Data Indicator

A **Quantifiable Data Indicator (QDI)** is a measurable value that helps track performance, trends, or progress over time. These indicators are based on **numerical data** and can be analyzed to make informed decisions.

Examples of Quantifiable Data Indicators

Business:

- Monthly sales revenue (\$50,000 in March).
- Customer retention rate (85% customers returned).
- Website traffic (100,000 visitors per month).

Healthcare:

- Patient recovery rate (92% success).
- Hospital bed occupancy rate (75% occupied).

Education:

- Student pass percentage (90% passed).
- Average test scores (85 out of 100).

Manufacturing:

- Defect rate in production (2% defective items).
- Number of units produced per hour (500 units/hour).

QDI vs. Qualitative Data Indicator

- **QDI (Quantifiable)** → Numerical, measurable (e.g., "500 units sold").
- **Qualitative Indicator** → Descriptive, subjective (e.g., "Customer satisfaction is high").



24. Cross Tab

A **Cross Tab (Cross Tabulation)** is a method of organizing and analyzing data by comparing two or more variables in a table format. It helps identify patterns, trends, and relationships between data points.

Example of Cross Tab Analysis

Scenario: A company surveys 1,000 customers about their preferred beverage (Tea or Coffee) based on age groups.

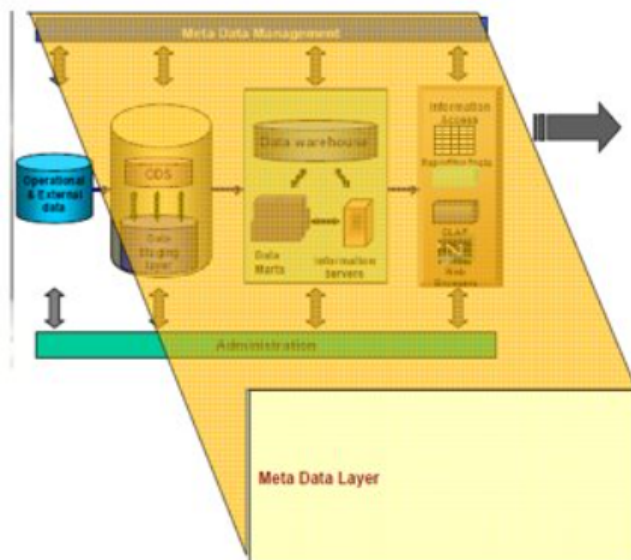
Age Group	Prefer Tea	Prefer Coffee	Total
18-25	200	300	500
26-40	150	250	400
41+	50	50	100

Insights from this Cross Tab:

- ✓ Young people (18-25) prefer Coffee more than Tea.
- ✓ Older people (41+) have an equal preference for both beverages.
- ✓ Total survey participants = 1,000.

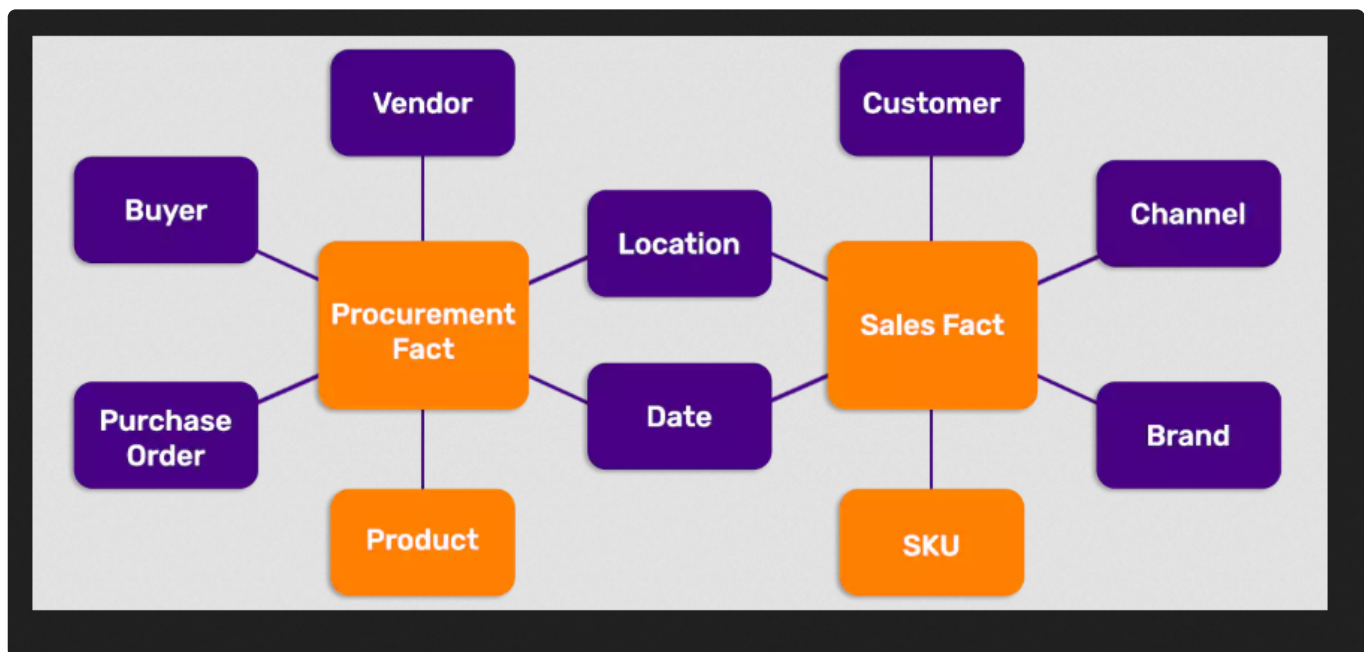


25. Metadata in datawarehouse



- Metadata is data about data.
- Stored in a repository.
- Contains all corporate metadata resources: database catalogs and data dictionaries.

26. Conformed Dimensions



A **Conformed Dimension** is a **dimension table** that is **shared across multiple fact tables** in a data warehouse. It allows different fact tables to be analyzed using the same dimension data, ensuring **consistency** across reports and queries.

📌 Key Characteristics:

- ✓ Shared across multiple fact tables
- ✓ Consistent across different business processes
- ✓ Enables accurate cross-domain reporting

Example:

Imagine a retail company with **Sales** and **Returns** fact tables. Both need **Date** and **Customer** dimensions. Instead of creating separate tables, a single **conformed dimension** is used.

```
Sales_Fact    ← Date_Dim
Returns_Fact ← Date_Dim
```

This ensures that the **Date dimension** is consistent for both sales and returns analysis.

Benefits of Conformed Dimensions:

- ✓ Eliminates redundancy by reusing dimensions.
- ✓ Ensures consistency across different fact tables.
- ✓ Supports enterprise-wide reporting.




26. Domain

A **Domain** in a data warehouse defines the **valid set of values** for a specific data field. It ensures **data consistency and quality** by restricting values to a predefined set.

Example:

- **Customer_Gender** domain: {Male, Female, Other}
- **Order_Status** domain: {Pending, Shipped, Delivered, Canceled}

 A **domain** prevents invalid data entry, ensuring uniformity across the data warehouse.



27. Hierarchy in Data Warehousing

A **Hierarchy** defines a **parent-child relationship** within a dimension, allowing **drill-down analysis** from a broad to a detailed level.

Example of a Geographic Hierarchy:

```
Continent → Country → State → City → Store
```

A report can be analyzed at different levels:

- ✓ Global sales (Continent Level)
- ✓ Country-wise sales (Country Level)
- ✓ City-wise sales (City Level)



28. Role playing dimension

A **role-playing dimension** is when **one dimension table is used multiple times** in the same fact table, but for **different roles**.

Example: Order Processing

A sales order might have:

- ✓ **Order Date** (When the order was placed)
- ✓ **Shipment Date** (When the order was shipped)
- ✓ **Delivery Date** (When the order reached the customer)

All these dates come from the **same "Date Dimension" table**, but they represent **different roles** in the fact table.

Table Structure:

Fact Table: Order_Fact

Order_ID	Order_Date_ID	Shipment_Date_ID	Delivery_Date_ID	Total_Amount
101	20240201	20240203	20240205	500
102	20240202	20240204	20240206	700

Dimension Table: Date_Dim

Date_ID	Full_Date	Month	Year
20240201	01-Feb-2024	Feb	2024
20240203	03-Feb-2024	Feb	2024

- ◆ Why is this useful?

- Instead of creating separate tables like `Order_Date_Dim` and `Shipment_Date_Dim`, we **reuse the same date table** with different column names in the fact table.
- Helps maintain **consistency** and **reduces redundancy**.

29. Degenerate dimension

A **degenerate dimension (DD)** is a dimension that exists in the fact table but does **NOT** have a separate dimension table.


Example: Order Transactions

Every order has a **unique order number**, but we **don't need a separate table for it**. Instead, it is stored **directly in the fact table**.

Table Structure:

Fact Table: `Order_Fact`

Order_Number	Date_ID	Customer_ID	Product_ID	Total_Amount
ORD1001	20240201	501	301	500
ORD1002	20240202	502	302	700

 **No separate `Order_Dim` table!**

♦ Why is this useful?

- Some values, like **Order Number** or **Invoice Number**, don't have additional attributes, so creating a separate table is unnecessary.
- Improves **query performance** because we avoid unnecessary joins.

30. Data Lake vs Data Warehouse

- **Data Lake** is ideal for storing **large volumes of diverse data** for future exploration.
- **Data Warehouse** is optimized for **structured data analysis** and **business reporting**.

31. In memory databases

- An in-memory database (IMDB) is a type of database that primarily relies on main memory (RAM) for data storage, rather than traditional disk-based storage.

- This makes it extremely fast for data access and manipulation, which is ideal for applications requiring real-time performance.

32. Schematic Evolution

- When you store structured data (like in a Parquet file), the schema defines the structure: field names, types, and order. Over time, you might need to:
 - Add new fields
 - Remove old fields
 - Change data types
 - Rename fields
- Schema evolution allows these changes to happen safely, so older data can still be read and new data can be written without compatibility issues.

33. Dimensional Modelling

Dimensional modelling

1. Fact Tables:

- Contain **measurable, quantitative data** (e.g., sales amount, quantity sold).
- Often very large.
- Include **foreign keys** to dimension tables.

2. Dimension Tables:

- Contain **descriptive attributes** (e.g., customer name, product category).
- Help provide context to facts.
- Typically denormalized for performance.

3. Star Schema:

- Central fact table surrounded by dimension tables.
- Simple and fast for querying.

4. Snowflake Schema:

- Dimensions are normalized into multiple related tables.
- More complex but can save space.

34. Conceptual Data Model

A **conceptual model** is a high-level design that shows how different entities relate to each other. It does **not** include technical details like data types or keys.

Example:

Imagine you're designing a university database. A conceptual model would include:

- **Entities:** Student, Course, Professor
- **Relationships:** A student **enrolls in** a course, a professor **teaches** a course
- ♦ Think of it as a rough sketch or blueprint of the system before adding technical details.



35. Physical Data Model

A **physical model** is the **detailed** implementation of the database, including table structures, columns, and data types.

Example:

For a **Student** entity, the physical model may look like this in SQL:

```
CREATE TABLE Student (  
    StudentID INT PRIMARY KEY,  
    Name VARCHAR(100),  
    Age INT,  
    Email VARCHAR(255) UNIQUE  
);
```

- ♦ The physical model ensures the database is optimized for **storage and performance**.



36. Surrogate Key

A **surrogate key** is a unique, system-generated identifier (usually a number). It is **not** derived from real-world data.

Example:

- Instead of using **email** as the primary key, we use a **StudentID** (1, 2, 3...).
- If a student changes their email, it won't affect the database structure.

- ♦ It's useful for large databases where natural keys (like email or phone numbers) can change.



37. ERWIN (Data Modeling Tool)

ERWIN is a tool for designing databases visually using **Entity-Relationship Diagrams (ERD)**.



38. Softwares Used for Data Modeling

There are several tools used for **designing databases visually**:

- **ERWIN** → Most popular tool for ER diagrams
- **Microsoft Visio** → Used for conceptual designs
- **Lucidchart** → Web-based modeling tool
- **MySQL Workbench** → Used for designing and managing MySQL databases

39. Data modelling types

Data modeling is classified into three levels:

1. **Conceptual Data Model** – Focuses on *what* the system contains.
2. **Logical Data Model** – Defines *how* data should be structured.
3. **Physical Data Model** – Describes *how* data is stored in a specific database.---

1 Conceptual Data Model

- High-level representation of business entities, attributes, and relationships.
- Created by business analysts and data architects.
- Independent of hardware, software, or database technologies.
- Helps stakeholders understand data at a business level.

Example:

- Entities: **Customer** and **Product**
- Attributes:

- Customer: `Customer_ID` , `Name` , `Email`
- Product: `Product_ID` , `Name` , `Price`
- Relationship: **A customer purchases a product.**

Characteristics:

- ✓ Business-oriented, not technical
- ✓ No details about tables, columns, or keys
- ✓ Provides a common vocabulary for stakeholders

2 Logical Data Model

- Defines the structure of data elements and relationships in detail.
- Still independent of the database management system (DBMS).
- Adds data attributes, data types, and normalization rules.

Example:

- Tables: `Customer` , `Product` , `Sales`
- Data Types:
 - `Customer_ID` → Integer (Primary Key)
 - `Name` → Varchar(50)
 - `Price` → Decimal(10,2)
- Relationship:
 - `Sales` table connects `Customer` and `Product` .

Characteristics:

- ✓ More detailed than the conceptual model
- ✓ Includes data types and normalization rules
- ✓ Helps in refining business rules before database implementation

3 Physical Data Model

- Database-specific representation of the logical data model.
- Defines tables, columns, constraints, indexes, triggers, etc.
- Developed for a specific DBMS (MySQL, PostgreSQL, SQL Server, etc.).

Example:

- **Table: Customer**

```
CREATE TABLE Customer (
    Customer_ID INT PRIMARY KEY,
    Name VARCHAR(50) NOT NULL,
    Email VARCHAR(100) UNIQUE
);
```

- **Table: Sales**

```
CREATE TABLE Sales (
    Sale_ID INT PRIMARY KEY,
    Customer_ID INT,
    Product_ID INT,
    Sale_Date DATE,
    FOREIGN KEY (Customer_ID) REFERENCES Customer(Customer_ID),
    FOREIGN KEY (Product_ID) REFERENCES Product(Product_ID)
);
```

Characteristics:

- ✓ Specific to a DBMS (includes SQL scripts)
- ✓ Defines indexes, constraints, and storage details
- ✓ Helps developers implement the actual database

Comparison Table

Feature	Conceptual Model	Logical Model	Physical Model
Purpose	High-level business view	Detailed structure and relationships	Implementation in a specific DBMS
Who Creates It?	Business Analysts, Data Architects	Data Architects, Analysts	DBAs, Developers
Focus	Entities, Attributes, Relationships	Data types, Normalization, Keys	Tables, Indexes, Constraints
Independence	Independent of DBMS	Independent of DBMS	Specific to DBMS

40. Examples of star and snowflake schema

Star Schema — (Simpler, Flattened)

What is it?

- A **central fact table** connected directly to **dimension tables**.
- Dimension tables are **not normalized** (they contain all descriptive info).

Example:

- Imagine you're building a **Sales Data Warehouse**.

Fact Table: `sales_fact`

<code>sale_id</code>	<code>product_id</code>	<code>customer_id</code>	<code>store_id</code>	<code>date_id</code>	<code>amount</code>
----------------------	-------------------------	--------------------------	-----------------------	----------------------	---------------------

Dimension Tables:

- `product_dim(product_id, product_name, category)`
- `customer_dim(customer_id, customer_name, gender, city)`
- `store_dim(store_id, store_name, location)`
- `date_dim(date_id, full_date, month, year)`

🎯 These dimension tables are **flat** — all data in one place.

Snowflake Schema — (Normalized, More Tables)

What is it?

- A **variation of star schema** where dimension tables are **normalized** into multiple related tables.
- Saves space, but is a bit more complex.

Example:

Same `sales_fact` table as before.

Now look at the dimensions:

- `product_dim(product_id, product_name, category_id)`
- `category_dim(category_id, category_name)`

- `customer_dim(customer_id, name, gender, city_id)`
- `city_dim(city_id, city_name, state_id)`
- `state_dim(state_id, state_name)`

🎯 Now dimensions are **split into smaller tables**, connected through foreign keys.

Summary Table:

Feature	Star Schema	Snowflake Schema
Complexity	Simple (flat)	Complex (normalized)
Joins	Fewer joins	More joins
Performance	Faster for queries	Slightly slower
Storage Space	More	Less (due to normalization)
Use Case	Best for BI tools	Good for space efficiency

“Star schema is simpler and used when we want fast reads.

Snowflake schema saves space by normalizing the dimensions, but involves more joins.”

41. Why Transformation is Needed?

- ♦ Data from different sources (databases, APIs, files) comes in different formats and structures.
- ♦ We **transform** data to make it **consistent, clean, and usable** for analysis.

📌 Example:

- A sales database stores `Date` as `YYYY-MM-DD` , but another system uses `DD/MM/YYYY` .
- To analyze data properly, we need to **standardize the date format** during transformation.



42. ETL vs. ELT

- ♦ **ETL (Extract → Transform → Load)**
- Data is **transformed first**, then loaded into the warehouse.
- Best for **structured** data.

♦ **ELT (Extract → Load → Transform)**

- Data is **loaded first**, then transformed inside the warehouse.
- Best for **Big Data & Cloud** environments.

Feature	ETL	ELT
Processing	Before loading	After loading
Storage	Smaller data	Requires large storage
Speed	Slower	Faster
Best for	Traditional Data Warehouses	Big Data, Cloud

Example:

- A **bank** may use **ETL** to clean and load financial data.
- A **cloud-based system** (like Google BigQuery) may use **ELT** to store raw data first and process it later.



43. *ELT Procedures*

♦ **ELT Procedures** define the steps for:

- ✓ **Extracting** data from multiple sources
- ✓ **Loading** data into the warehouse
- ✓ **Transforming** it inside the warehouse

Example Steps:

- 1 Extract** sales data from an e-commerce website and a CRM system.
- 2 Load** raw data into a cloud data warehouse (like Snowflake).
- 3 Transform** the data to standardize customer names, dates, and sales formats.



44. *SCDs (Slowly Changing Dimensions)*

- ♦ **SCDs** handle changes in dimension data (e.g., customer address, job title).

Example:

- A customer **moves to a new city** or **gets promoted at work**.
- How do we store the old and new values in a data warehouse?



45. Types of SCDs

Slowly Changing Dimensions (SCDs) are techniques used in data warehousing to manage changes in dimension data over time. They are categorized into different types (SCD 0 to SCD 6), based on how changes are tracked and stored. Let's break them down with simple explanations and examples.

SCD Type 0 (Retain Original) – No Change

👉 **No changes are allowed in the dimension table.** Data remains the same even if the source system updates it.

♦ **Example:** A product table where historical prices must never change, even if the company updates them.

Product_ID	Product_Name	Price
101	Laptop	50000
102	Phone	20000

♦ If the price of "Laptop" changes to 55000, it **won't** be updated in the dimension table.

SCD Type 1 (Overwrite) – Keep Only Latest Data

👉 **Old data is replaced with new data, and history is lost.**

♦ **Example:** A customer table where only the latest address is stored.

Customer_ID	Name	Address
201	John	Pune

♦ If John moves to Mumbai, the table updates:

Customer_ID	Name	Address
201	John	Mumbai

Old address is lost.

SCD Type 2 (Versioning) – Maintain History with New Rows

👉 A new row is added for each change, maintaining historical records.

👉 Usually, start and end dates or version numbers are used.

- ♦ **Example:** Employee salary history

Emp_ID	Name	Salary	Start_Date	End_Date	Is_Current
301	Alice	50000	2023-01-01	2024-01-01	No
301	Alice	55000	2024-01-02	NULL	Yes

- ♦ If Alice's salary increases to 60000, a new row is added with updated dates.

SCD Type 3 (Add Column) – Keep Only One Previous Value

👉 Stores the previous value in a separate column, but doesn't keep full history.

- ♦ **Example:** Product price changes

Product_ID	Name	Current_Price	Previous_Price
401	TV	40000	35000

- ♦ If the price updates to 45000, the table updates:

Product_ID	Name	Current_Price	Previous_Price
401	TV	45000	40000

Only the last change is kept.

SCD Type 4 (Separate History Table) – Archive Old Data

👉 Current data stays in the main table, and historical data is moved to a separate history table.

♦ Example:

Main Table (Current Data):

Customer_ID	Name	Address
501	Sam	Mumbai

History Table (Old Data):

Customer_ID	Name	Old_Address	Changed_Date
501	Sam	Pune	2024-02-01

If Sam moves again, the old address moves to the history table.

SCD Type 5 (Hybrid – Type 1 + Type 4)

👉 Uses both Type 1 (overwrite current data) and Type 4 (history table).

👉 Adds a surrogate key to track changes.

♦ Example:

Current Table (Type 1 Overwrite):

Customer_SK	Customer_ID	Name	Address
1001	601	Raj	Delhi

History Table (Type 4 Archive):

Customer_SK	Customer_ID	Name	Old_Address	Changed_Date
1000	601	Raj	Jaipur	2024-02-15

SCD Type 6 (Hybrid – Type 1 + Type 2 + Type 3)

👉 Combines Type 1 (overwrite), Type 2 (new row for changes), and Type 3 (store previous value).

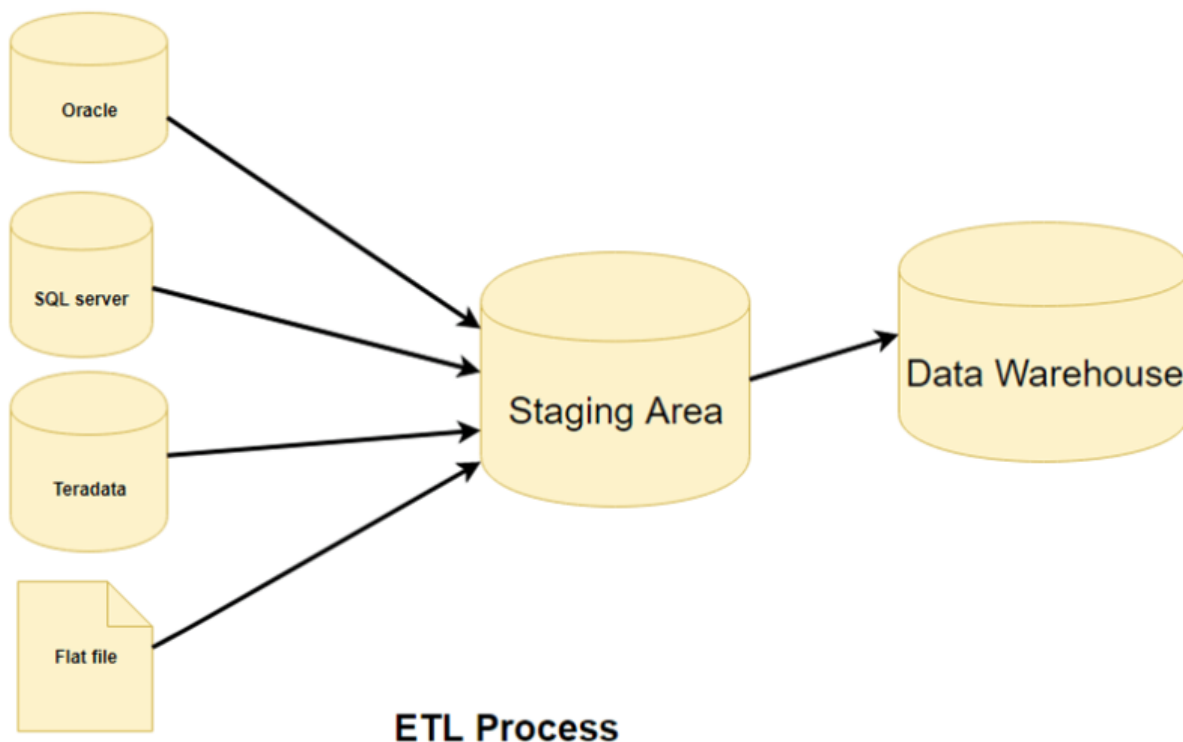
♦ Example:

Emp_ID	Name	Salary	Previous_Salary	Start_Date	End_Date	Is_Current
701	Mike	70000	NULL	2023-01-01	2024-02-01	No
701	Mike	75000	70000	2024-02-02	NULL	Yes

♦ If Mike gets a raise, a new row is added (Type 2), the current row is updated (Type 1), and the old salary is stored in a separate column (Type 3).



46. What is ETL



ETL (Extract, Transform, Load) is a process used in data warehousing to collect data from multiple sources, clean and transform it, and then load it into a data warehouse for analysis. Here's a breakdown of each step:

Step 1: Extraction

- Extracts data from multiple sources like databases, legacy systems, flat files, ERP systems, etc.
- Uses a **staging area** to avoid impacting source system performance and allow validation before loading.
- Data mapping is done to define the relationship between source and target.

Extraction Methods:

1. **Full Extraction** – Extracts all data each time.
2. **Partial Extraction (without update notification)** – Extracts only changed data, but changes are identified manually.
3. **Partial Extraction (with update notification)** – Extracts only changed data using timestamps or database triggers.

Validations during Extraction:

- Ensure data integrity and remove duplicates.
- Validate data types and formats.
- Filter out irrelevant or incorrect data.



Step 2: Transformation

- Converts raw extracted data into a usable format by applying business rules, cleansing, and data integration.
- Includes processes like:
 - **Data Cleansing** – Removing duplicates, handling missing values.
 - **Data Mapping** – Aligning data structures from different sources.
 - **Aggregation** – Summarizing data (e.g., calculating total sales).
 - **Joining & Splitting** – Merging data from multiple sources or splitting complex fields.

Example: Combining `First Name` and `Last Name` into a single `Full Name` field.



Step 3: Loading

- Moves transformed data into the **data warehouse** for analytics and reporting.
- Needs to be optimized for performance, ensuring integrity and consistency.

Types of Loading:

1. **Initial Load** – First-time full data population.
2. **Incremental Load** – Only updates changed or new data.
3. **Full Refresh** – Deletes existing data and reloads from scratch.

Load Verification:

- Ensure primary keys and relationships are correct.
- Validate BI reports and ensure aggregated values are accurate.
- Check historical data updates in slowly changing dimensions (SCD).

ETL plays a crucial role in Business Intelligence (BI) and decision-making by ensuring high-quality, structured data is available for analysis.

47. What is Change Data Capture

Change Data Capture (CDC) is a technique used in databases to **track and capture changes** (INSERTs, UPDATEs, DELETEs) made to data in real-time or near-real-time.

“Change Data Capture means **capturing only the changes** in the source database, instead of reloading everything.

It helps keep downstream systems, like data warehouses or real-time dashboards, up to date efficiently.”