

Topic 3: Dimensional Modeling

At the end of this topic, students should be able to discuss:

- Star schema and snowflake schema
- Fact and dimension tables
- Measures and metrics
- Hierarchies and aggregates
- Case study: Dimensional modeling for a retail business

1. Overview of Dimensional Modeling

Dimensional Modeling

Dimensional modeling is a design technique used to structure data for reporting and analytical purposes in a way that is user-friendly and optimized for fast querying. It primarily involves organizing data into **facts** (measurable data) and **dimensions** (descriptive data) that users can query easily for business insights. Key components of dimensional modeling include **star schemas**, **snowflake schemas**, and the understanding of **fact tables**, **dimension tables**, **measures**, and **metrics**.

1. Star Schema and Snowflake Schema

Star Schema:

- **Overview:** The star schema is a simple and widely used dimensional model where a central fact table is connected to multiple dimension tables. The dimension tables are denormalized, meaning they store redundant data to simplify querying and improve performance.
- **Structure:**
 - The **fact table** contains the primary data (e.g., sales transactions), and each row in the fact table corresponds to a unique event or transaction.
 - The **dimension tables** are descriptive tables (e.g., time, product, store) that provide context to the data in the fact table. Each dimension table is linked to the fact table by a foreign key.
- **Advantages:**
 - Simple and intuitive design.
 - Fast query performance due to denormalization of dimensions.
 - Ideal for end-users, as it is easy to understand and query.
- **Disadvantages:**
 - Data redundancy in dimension tables.
 - Maintenance challenges when dimension data changes.

Snowflake Schema:

- **Overview:** The snowflake schema is a more complex version of the star schema. It normalizes the dimension tables, meaning related data is stored in separate tables

(e.g., a product dimension could be split into product_category and product_subcategory tables).

- **Structure:**
 - The **fact table** remains the same as in the star schema.
 - The **dimension tables** are normalized into multiple related tables, forming a snowflake-like structure.
- **Advantages:**
 - Reduces data redundancy by normalizing dimensions.
 - More efficient in terms of storage.
- **Disadvantages:**
 - More complex queries due to the need to join multiple tables.
 - Slower query performance compared to the star schema.

2. Fact and Dimension Tables

Fact Tables:

- **Overview:** A fact table stores quantitative data that represents business transactions or events (e.g., sales, purchases, inventory movements). It is typically large and contains data that can be aggregated for analysis.
- **Key Characteristics:**
 - Contains **measures** (e.g., sales revenue, quantity sold) as columns.
 - Includes **foreign keys** linking to dimension tables.
 - Typically has a primary key consisting of a combination of foreign keys.
- **Example:** In a retail business, a **sales fact table** might include columns like:
 - Transaction_ID
 - Product_ID
 - Customer_ID
 - Store_ID
 - Date_ID
 - Sales_Amount
 - Quantity_Sold

Dimension Tables:

- **Overview:** Dimension tables store descriptive or categorical data that provide context for the measures in the fact table. These tables typically contain attributes that describe the entities involved in the business event.
- **Key Characteristics:**
 - Contains **descriptive attributes** (e.g., product name, customer address).
 - Includes a **primary key** that is referenced by foreign keys in the fact table.

- **Example:** In a retail business, a **product dimension table** might include columns like:
 - Product_ID (primary key)
 - Product_Name
 - Product_Category
 - Product_Subcategory
 - Brand
 - Supplier

3. Measures and Metrics

Measures are quantitative values stored in the fact tables that are typically analyzed or aggregated in reporting. They represent the “what” of the data (e.g., sales amounts, units sold).

Metrics are calculations or derived measurements that are often based on the raw measures. Metrics can be used to derive insights and KPIs (Key Performance Indicators) for business analysis (e.g., average sales per customer, profit margin).

- **Example Measures:**
 - Sales_Amount: The total dollar value of sales transactions.
 - Quantity_Sold: The number of items sold.
- **Example Metrics:**
 - **Profit Margin:** $(\text{Sales_Amount} - \text{Cost_Amount}) / \text{Sales_Amount}$.
 - **Average Sales per Customer:** $\text{Sales_Amount} / \text{Number_of_Customers}$.

4. Hierarchies and Aggregates

Hierarchies refer to the structured levels of granularity in dimension tables, allowing data to be grouped or aggregated at different levels for analysis. They provide a way to drill down or roll-up data in reports.

- **Example of Hierarchy in the Product Dimension:**
 - **Product Category → Product Subcategory → Product.**

This hierarchy allows for analysis at different levels, such as total sales by product category, subcategory, or individual product.

Aggregates are pre-computed values that summarize the data in the fact table, typically based on different levels of granularity (e.g., daily, monthly, quarterly). Aggregating data can improve query performance, especially for large datasets.

- **Example of Aggregates:**
 - **Monthly Sales Total:** Total sales amount for each month.
 - **Quarterly Sales Average:** Average sales amount for each quarter.

- These aggregates can be stored in separate tables or computed dynamically based on user queries.

5. Case Study: Dimensional Modeling for a Retail Business

Let's take a **retail business** as an example to demonstrate dimensional modeling.

Scenario: A retail company wants to analyze its sales data to understand trends, product performance, and customer behavior. The data comes from various transactional systems, and the company wants to create a data warehouse for reporting.

Fact Table: The central fact table in this case could be a **Sales Fact Table**, which would record the sales transactions.

Column Name	Data Type	Description
Sales_ID	Integer	Unique identifier for each sales transaction.
Product_ID	Integer	Foreign key to the Product dimension.
Customer_ID	Integer	Foreign key to the Customer dimension.
Store_ID	Integer	Foreign key to the Store dimension.
Date_ID	Integer	Foreign key to the Date dimension.
Sales_Amount	Decimal	The total dollar value of the sale.
Quantity_Sold	Integer	Number of units sold.

Dimension Tables: Here are some of the key dimension tables in this scenario:

- **Product Dimension Table:** Stores details about products.

Column Name	Data Type	Description
Product_ID	Integer	Unique identifier for the product.
Product_Name	String	Name of the product.
Category	String	Product category (e.g., Electronics).
Subcategory	String	Product subcategory (e.g., Laptops).

- **Customer Dimension Table:** Stores details about customers.

Column Name	Data Type	Description
Customer_ID	Integer	Unique identifier for the customer.
Customer_Name	String	Name of the customer.
Gender	String	Gender of the customer.
Location	String	Customer's location (e.g., city).

- **Store Dimension Table:** Stores details about stores.

Column Name	Data Type	Description
Store_ID	Integer	Unique identifier for the store.
Store_Name	String	Name of the store.
Region	String	Region where the store is located.

- **Date Dimension Table:** Stores information about dates for reporting.

Column Name	Data Type	Description
Date_ID	Integer	Unique identifier for the date.
Date	Date	The date of the sale.
Month	String	The month of the sale.
Year	Integer	The year of the sale.

Aggregates and Hierarchies: The company might want to create reports at various levels of granularity, such as:

- **Monthly Sales:** Total sales for each month.
- **Product Sales by Category:** Aggregating sales by product category.
- **Customer Segment Analysis:** Aggregating sales by customer demographics (e.g., gender, location).

This can be achieved using **hierarchies** in the dimension tables (e.g., Product Category → Product) and **aggregated fact tables** (e.g., monthly sales totals).

Note:

Dimensional modeling, with its focus on organizing data into facts and dimensions, allows for efficient and intuitive reporting, especially for large data sets in environments like retail.

2. Star Schema vs Snowflake Schema

Both the **Star Schema** and **Snowflake Schema** are popular data modeling techniques used in **dimensional modeling** for data warehouses. They both have the same fundamental goal: to organize data for efficient querying and reporting, especially for analytical purposes. However, they differ significantly in their structure and complexity.

1. Star Schema

Overview:

The **star schema** is a simple and widely used database schema that organizes data into fact tables and dimension tables. The fact table is at the center of the schema, and the dimension tables surround it like the points of a star.

Structure:

- **Fact Table:** The central table in the schema contains the quantitative data or measures (e.g., sales amount, quantity sold).
- **Dimension Tables:** Surrounding the fact table are dimension tables that contain descriptive attributes (e.g., customer details, product names, time periods). These dimension tables are directly linked to the fact table via foreign keys.

Key Features:

- **Denormalization:** Dimension tables are typically denormalized, meaning they store all attributes in a single table. This reduces the need for joins and improves query performance.
- **Simplicity:** The structure is easy to understand and use, making it ideal for end-users and business analysts.
- **Performance:** Since dimension tables are denormalized, queries tend to be faster as they involve fewer joins.

Example of Star Schema in Retail:

- **Fact Table:** Sales Fact Table (contains Sales_Amount, Quantity_Sold, Date_ID, Product_ID, Customer_ID).
- **Dimension Tables:**
 - **Product Dimension** (contains Product_ID, Product_Name, Product_Category).
 - **Customer Dimension** (contains Customer_ID, Customer_Name, Customer_Gender).
 - **Date Dimension** (contains Date_ID, Date, Month, Year).

Advantages:

- Simple to design and easy to query.
- Fast query performance due to fewer joins.
- Intuitive for business users and analysts.

Disadvantages:

- Data redundancy in dimension tables (because of denormalization).
- More storage is required for the redundant data.
- Changes in dimension data can be challenging to manage (e.g., changes in product names or categories).

2. Snowflake Schema

Overview: The **snowflake schema** is a more normalized version of the star schema. In this schema, the dimension tables are broken down into smaller, related tables, creating a “snowflake-like” structure. It is a more complex design compared to the star schema.

Structure:

- **Fact Table:** Just like in the star schema, the fact table contains quantitative data and links to the dimension tables through foreign keys.
- **Dimension Tables:** The dimension tables are normalized, meaning they are split into multiple related tables. For example, instead of having a single Product Dimension table, it might be split into Product, Product_Category, and Product_Subcategory tables.

Key Features:

- **Normalization:** Dimension tables are normalized, meaning they are split into multiple tables to eliminate redundancy. This ensures that related data is stored in separate tables.
- **Complexity:** The snowflake schema is more complex than the star schema because it requires more joins between tables. This can make it harder for end-users to understand.
- **Storage Efficiency:** Since dimension data is normalized, the snowflake schema uses less storage compared to the star schema.

Example of Snowflake Schema in Retail:

- **Fact Table:** Sales Fact Table (contains Sales_Amount, Quantity_Sold, Date_ID, Product_ID, Customer_ID).
- **Dimension Tables:**
 - **Product** (contains Product_ID, Product_Name).
 - **Product Category** (contains Category_ID, Category_Name).
 - **Customer** (contains Customer_ID, Customer_Name, Customer_Gender).
 - **Date** (contains Date_ID, Date, Month, Year).

Advantages:

- Reduces redundancy and saves storage space.
- Ensures better data consistency (no duplicated data in dimension tables).
- More flexible in handling changes in dimension data.

Disadvantages:

- More complex design, which may be harder to maintain.
- Requires more joins, which can result in slower query performance.
- More difficult for business users to understand and work with.

Comparison: Star Schema vs Snowflake Schema

Feature	Star Schema	Snowflake Schema
---------	-------------	------------------

Feature	Star Schema	Snowflake Schema
Complexity	Simple and intuitive	More complex and normalized
Normalization	Denormalized (redundant data in dimensions)	Normalized (reduces redundancy)
Query Performance	Faster (fewer joins)	Slower (more joins due to normalization)
Storage Efficiency	Less efficient (more storage required)	More efficient (saves storage)
Data Redundancy	Higher redundancy in dimension tables	Less redundancy (normalized data)
Maintenance	Easier to maintain, but harder to handle dimension changes	More difficult to maintain, but easier to handle dimension changes
Usability	Easier for business users to understand	More difficult for business users to navigate

When to Use Star Schema vs Snowflake Schema

- **Star Schema:**
 - Use when simplicity and fast query performance are critical.
 - Ideal for **smaller to medium-sized data warehouses** with stable and relatively simple data structures.
 - Best suited for **end-users** and **business analysts** who need easy access to data and fast queries.
- **Snowflake Schema:**
 - Use when storage efficiency is a priority, and the data is complex or prone to frequent changes.
 - Ideal for **large data warehouses** where dimensional data is subject to frequent changes (e.g., product catalogs, customer details).
 - Suitable for environments where **data integrity** and **consistency** are paramount and where normalization can reduce redundancy.

Note:

The choice between a **star schema** and **snowflake schema** depends on the specific requirements of the organization or business, such as performance, storage, and complexity needs. While the **star schema** offers simplicity and faster performance, the **snowflake schema** offers better storage efficiency and data integrity but requires more complex queries and maintenance.

Example: Retail Business Sales Data

Scenario:

A retail business wants to analyze its sales data. The business needs to track sales, customers, products, and time. The data collected will be stored in a **Data Warehouse** for reporting and analysis.

1. Star Schema Example

In the **Star Schema**, we have the **fact table** at the center and **dimension tables** surrounding it. The **fact table** contains the sales data (measures), while the **dimension tables** contain descriptive data about the products, customers, and time.

Fact Table: Sales Fact Table

Sales_ID	Product_ID	Customer_ID	Date_ID	Sales_Amount	Quantity_Sold
1	101	5001	2023-01	100	2
2	102	5002	2023-01	150	3
3	103	5001	2023-02	200	5

Dimension Tables:

• Product Dimension Table

• Product_ID Product_Name Product_Category
• 101 Laptop Electronics 102 Phone Electronics
• 103 Tablet Electronics

• Customer Dimension Table

• Customer_ID Customer_Name Customer_Gender Customer_Location
• 5001 John Doe Male New York
• 5002 Jane Smith Female Los Angeles

• Date Dimension Table

• Date_ID Date Month Year
• 2023-01 2023-01-01 January 2023
• 2023-02 2023-02-01 February 2023

In the **Star Schema**:

- The **fact table** (Sales) contains quantitative data, such as Sales_Amount and Quantity_Sold.
- The **dimension tables** (Product, Customer, Date) contain descriptive attributes like product names, customer information, and date details.
- The dimension tables are **denormalized**, meaning all the related data is in a single table.

2. Snowflake Schema Example

In the **Snowflake Schema**, the **dimension tables** are **normalized** into multiple related tables, breaking down the data further to reduce redundancy.

Fact Table: Sales Fact Table

The **Sales Fact Table** remains the same as in the **Star Schema**:

Sales_ID	Product_ID	Customer_ID	Date_ID	Sales_Amount	Quantity_Sold
1	101	5001	2023-01	100	2
2	102	5002	2023-01	150	3
3	103	5001	2023-02	200	5

Dimension Tables:

- **Product Dimension Table** (Normalized)

• Product_ID Product_Name Category_ID
• 101 Laptop 1
• 102 Phone 1
• 103 Tablet 1

-

- **Product Category Dimension Table** (Newly normalized)

• Category_ID Category_Name
• 1 Electronics

- **Customer Dimension Table**

• Customer_ID Customer_Name Customer_Gender Location_ID
• 5001 John Doe Male 101
• 5002 Jane Smith Female 102

-

- **Location Dimension Table** (Newly normalized)

• Location_ID Location_Name
• 101 New York 102 Los Angeles

-

- **Date Dimension Table**

• Date_ID Date Month_ID Year
• 2023-01 2023-01-01 1 2023
• 2023-02 2023-02-01 2 2023

-

- **Month Dimension Table** (Newly normalized)

• Month_ID Month_Name
• 1 January 2 February

Differences Between Star and Snowflake Schema

Feature	Star Schema	Snowflake Schema
Fact Table	Same in both schemas	Same in both schemas
Dimension Tables	Denormalized, one table per dimension	Normalized, multiple tables per dimension
Redundancy	More redundancy in dimension tables	Less redundancy due to normalization
Query Complexity	Simpler queries, fewer joins	More complex queries, more joins
Performance	Faster queries due to fewer joins	Slower queries due to more joins
Storage Efficiency	Less efficient (more storage required)	More efficient (less storage required)

Note:

- **Star Schema** is simple, with fewer tables and fast queries, but can lead to data redundancy and higher storage requirements.
- **Snowflake Schema** is more complex, with normalized dimension tables, leading to better storage efficiency but slower query performance due to the need for more joins.

Each schema has its place depending on the **size** and **complexity** of the data warehouse, as well as **query performance** and **storage** requirements.

3. Fact and Dimension Tables in Dimensional Modeling

In a data warehouse, **fact** and **dimension** tables are core components of the schema. Together, they enable efficient organization and retrieval of data for analysis and reporting.

1. Fact Tables

A **fact table** is the central table in a star or snowflake schema. It stores **quantitative data** or **measures** that are used for analysis, such as sales revenue, quantities, or performance metrics.

Characteristics:

- **Measures/Key Performance Indicators (KPIs):** Contains numerical data that can be aggregated (e.g., sum, average, count).
- **Foreign Keys:** Links to dimension tables to provide context to the measures.
- **Grain:** Represents the level of detail captured in the table (e.g., daily sales, per-transaction data).

Example: Sales Fact Table

Fact Table Attributes		Description			
Sales_ID		Unique identifier for each sale			
Product_ID		Links to the Product Dimension			
Customer_ID		Links to the Customer Dimension			
Date_ID		Links to the Date Dimension			
Sales_Amount		Total sales value			
Quantity_Sold		Number of items sold			

Sales_ID	Product_ID	Customer_ID	Date_ID	Sales_Amount	Quantity_Sold
1	101	5001	2023-01	100	2
2	102	5002	2023-01	150	3
3	103	5001	2023-02	200	5

2. Dimension Tables

Dimension tables provide descriptive or categorical information to provide context to the measures in the fact table. They answer the “who, what, where, when, and why” questions about the data.

Characteristics:

- **Descriptive Attributes:** Contain textual or categorical data (e.g., product names, customer demographics).
- **Surrogate Keys:** Unique identifiers for each row, often used as primary keys.
- **Denormalized Structure:** In a star schema, dimension tables are usually denormalized for simpler and faster querying. In a snowflake schema, they are normalized.

Example: Product Dimension Table

Dimension Table Attributes		Description	
Product_ID		Unique identifier for each product	
Product_Name		Name of the product	
Product_Category		Category of the product	
Product_Brand		Brand of the product	

Product_ID	Product_Name	Product_Category	Product_Brand
101	Laptop	Electronics	Dell
102	Phone	Electronics	Apple

Product_ID	Product_Name	Product_Category	Product_Brand
103	Tablet	Electronics	Samsung

Example: Customer Dimension Table

Customer_ID	Customer_Name	Customer_Gender	Customer_Location
5001	John Doe	Male	New York
5002	Jane Smith	Female	Los Angeles

Relationship Between Fact and Dimension Tables

- Fact tables are linked to dimension tables through **foreign keys**.
- Dimension tables provide the **context** for interpreting the numerical data in the fact table.
- This relationship forms the basis of **dimensional modeling** and enables analytical queries.

Example Query:

“What was the total sales amount for all laptops sold in January 2023?”

1. **Fact Table:** Query the Sales_Amount for the Product_ID that corresponds to “Laptop” and Date_ID in “January 2023.”
2. **Dimension Tables:**
 - Use the **Product Dimension Table** to find the Product_ID for “Laptop.”
 - Use the **Date Dimension Table** to identify the Date_ID for “January 2023.”

Summary

Aspect	Fact Table	Dimension Table
Purpose	Stores quantitative data (measures).	Provides descriptive context to the facts.
Data Type	Numeric and foreign keys.	Textual and categorical attributes.
Grain	Represents the level of detail.	Denormalized or normalized attributes.
Examples	Sales amount, quantity, revenue.	Product name, category, customer details.
Joins	Joins with dimension tables via foreign keys.	Links to fact tables through surrogate keys.

This structure allows for fast and intuitive analytical queries in a data warehouse environment.

4. Measures and Metrics in Data Warehousing

In a data warehouse, **measures** and **metrics** are essential for analyzing performance, identifying trends, and making data-driven decisions. These concepts are closely related but have distinct roles in dimensional modeling and reporting.

1. Measures

Measures are **numerical values** that represent quantitative data stored in the **fact tables** of a data warehouse. They are the raw data points that can be aggregated, calculated, or analyzed to derive insights.

Characteristics:

- Stored in **fact tables**.
- Represent **quantitative data** (e.g., sales revenue, profit, number of transactions).
- Can be aggregated using functions like **SUM**, **AVERAGE**, **COUNT**, **MAX**, and **MIN**.
- The level of granularity depends on the grain of the fact table (e.g., daily sales vs. hourly sales).

Example Measures:

- **Sales Amount:** The revenue from a sale.
- **Quantity Sold:** The number of items sold.
- **Profit:** Revenue minus costs.
- **Order Count:** The number of orders.

Example:

Sales_ID	Product_ID	Date_ID	Sales_Amount	Quantity_Sold
1	101	2023-01	100	2
2	102	2023-01	150	3

2. Metrics

Metrics are derived values or **calculated indicators** based on one or more measures. They are used to evaluate performance, measure progress, or track trends over time. Metrics often involve business logic and are more meaningful for decision-making than raw measures.

Characteristics:

- Calculated using one or more **measures**.
- Provide **context** and are aligned with business goals or Key Performance Indicators (KPIs).
- Typically stored in reports or dashboards rather than fact tables.

Example Metrics:

Example Metrics:

- Average Order Value (AOV): $AOV = \frac{\text{Total Sales Amount}}{\text{Total Orders}}$
- Profit Margin: $\text{Profit Margin} = \frac{\text{Profit}}{\text{Sales Amount}} \times 100\%$
- Conversion Rate: $\text{Conversion Rate} = \frac{\text{Number of Purchases}}{\text{Number of Visitors}} \times 100\%$
- Year-over-Year Growth: Measures growth percentage compared to the previous year.

Example:

Metric	Formula	Example Value
Average Order Value	Total Sales Amount / Total Orders	\$125
Profit Margin	(Profit / Sales Amount) \times 100%	20%

Key Differences Between Measures and Metrics

Aspect	Measures	Metrics
Definition	Raw numerical data.	Calculated or derived indicators.
Source	Stored in fact tables.	Derived from measures using formulas.
Aggregation	Directly aggregable (SUM, AVG, etc.).	Requires business logic or calculations.
Examples	Sales amount, quantity sold.	Average order value, profit margin.
Purpose	Basis for analysis.	Evaluation of performance or trends.

Use Case in Retail Business

Measures:

- **Sales Amount:** Total revenue from transactions.
- **Quantity Sold:** Number of products sold.

Metrics:

- Total Sales for January: SUM(Sales Amount for January).
- Average Sales per Customer: Total Sales Amount/Number of Customers.
- Profit Margin: (Sales Amount – Cost of Goods Sold)/Sales Amount \times 100%.

Note:

- **Measures** are raw, aggregable data points that represent **what happened** (e.g., how much revenue was generated).
- **Metrics** are derived, calculated indicators that help evaluate **why or how it happened** (e.g., is the revenue meeting goals?).

Both measures and metrics are crucial for **data analysis** and **decision-making** in data warehousing and business intelligence systems.

5. Hierarchies and Aggregates in Data Warehousing

In a data warehouse, **hierarchies** and **aggregates** are used to structure data and enhance query performance, making it easier for users to analyze information at different levels of detail.

1. Hierarchies

A **hierarchy** defines a logical organization of data, allowing users to navigate from a high-level summary to more detailed data or vice versa. Hierarchies are often present in **dimension tables** and are key to organizing attributes for **drill-down** and **roll-up** operations in OLAP (Online Analytical Processing).

Characteristics:

- Represent **levels of granularity** in data.
- Allow users to perform operations like **drilling down** (exploring detailed data) and **rolling up** (summarizing data).
- Commonly used in time, geography, product, or organizational data.

Example Hierarchies:

1. **Time Hierarchy:** Year → Quarter → Month → Week → Day
2. **Geography Hierarchy:** Country → State → City → Store
3. **Product Hierarchy:** Category → Subcategory → Product → SKU

Example: Time Dimension Table

Date_ID	Day	Week	Month	Quarter	Year
2023-01-01	Monday	Week 1	January	Q1	2023
2023-02-01	Wednesday	Week 5	February	Q1	2023

In this example, users can:

- Drill down from **Year** to **Month** to **Day**.
- Roll up from **Day** to **Month** to **Year**.

2. Aggregates

Aggregates are **precomputed summaries** of data stored in the fact table. They improve the performance of queries that involve aggregations (e.g., SUM, AVERAGE, COUNT) by reducing the volume of data processed at runtime.

Characteristics:

- Precomputed and stored in aggregate tables.
- Reduce query response time for summary-level reports.
- Can be based on specific hierarchies (e.g., monthly sales, sales by product category).

Types of Aggregates:

- **Summed Aggregates:** Total sales, total quantity.
- **Average Aggregates:** Average sales per month, average order size.
- **Count Aggregates:** Number of transactions, number of customers.

Example Aggregate Table: Sales Summary by Month and Product Category

Year	Month	Product_Category	Total_Sales	Total_Quantity
2023	January	Electronics	200,000	1,500
2023	February	Electronics	180,000	1,300

In this table:

- Instead of querying the detailed fact table, users can directly access precomputed totals for faster reporting.

Relationship between Hierarchies and Aggregates

1. **Hierarchies Enable Aggregates:** Hierarchies define how data is organized and what levels of aggregation are possible. For example:
 - In a **time hierarchy**, aggregates might include monthly or yearly sales totals.
 - In a **geographic hierarchy**, aggregates might include sales by city or country.
2. **Aggregates Use Hierarchies:** Aggregates are often precomputed at specific levels of a hierarchy to enhance query performance. For instance:
 - Sales summaries can be precomputed by **Year** and **Quarter** in the time hierarchy.
 - Inventory levels can be summarized by **Store** and **City** in the geography hierarchy.

Use Case Example

Retail Business

- **Hierarchy:**
 - Time Dimension: Year → Quarter → Month → Day
 - Geography Dimension: Country → State → City → Store

- **Aggregate Tables:**

- Sales Summary by Month and Product Category:

Year Month Category Total_Sales

2023 January Electronics 150,000

- Sales Summary by City and Product:

City	Product	Total_Sales
New York	Laptop	50,000
Los Angeles	Phone	30,000

Query Optimization:

- Without aggregates, querying total monthly sales requires scanning the detailed **fact table**.
- With aggregates, the precomputed data provides immediate results.

Benefits of Hierarchies and Aggregates

Aspect	Hierarchies	Aggregates
Purpose	Organize data for drill-down/roll-up.	Enhance query performance by precomputing summaries.
Data Granularity	Defines levels of granularity.	Summarizes data at specific levels.
Performance Impact	Simplifies navigation in reports.	Reduces runtime query load.
Use Cases	Reporting and OLAP analysis.	High-frequency summary queries.

Note:

- **Hierarchies** allow users to explore data across levels of detail and are foundational for designing effective dimensional models.
- **Aggregates** leverage these hierarchies to precompute and store summarized data, optimizing the performance of analytical queries.

Both play critical roles in ensuring data warehouses support **scalable, fast, and meaningful analysis**

6. Information Package in Data Warehousing

An **information package** is a structured way to define the key aspects of data required for business analysis. It is a planning tool used during the design phase of a data warehouse to

capture the essential business requirements. The concept revolves around identifying facts, dimensions, and metrics relevant to a business process and organizing them into a cohesive structure.

Key Components of an Information Package

1. **Facts**
 - **Definition:** Quantitative data that represents a business measure or metric.
 - **Examples:** Sales revenue, profit, quantity sold, or order count.
 2. **Dimensions**
 - **Definition:** Qualitative attributes or descriptive data that provide context to facts.
 - **Examples:** Time, product, location, customer, or department.
 3. **Attributes**
 - **Definition:** Detailed properties or characteristics of dimensions.
 - **Examples:** For the "Product" dimension, attributes might include product name, category, and brand.
 4. **Hierarchies**
 - **Definition:** Organized levels within dimensions that allow data to be analyzed at various granularities.
 - **Examples:** For the "Time" dimension, hierarchies might include year → quarter → month → day.
 5. **Granularity**
 - **Definition:** The level of detail stored in the data warehouse.
 - **Examples:** Data stored at the daily transaction level or aggregated to monthly sales.
 6. **Measure Types**
 - **Definition:** Metrics can be additive (e.g., sales), semi-additive (e.g., inventory), or non-additive (e.g., ratios).
-

Purpose of an Information Package

- Helps identify the key questions the data warehouse needs to answer.
 - Serves as a blueprint for designing fact tables and dimension tables.
 - Aligns data warehouse design with business objectives.
-

Steps to Develop an Information Package

- 1. **Identify Business Processes**
 - Understand the business operations that need analysis, such as sales, inventory, or customer behavior.
- 2. **Define the Facts**
 - Determine what numerical measures are critical for the business.
- 3. **Identify Dimensions**
 - Choose the dimensions needed to analyze the facts effectively.
- 4. **Determine Granularity**
 - Decide the level of detail at which facts and dimensions will be stored.
- 5. **Document Attributes and Hierarchies**
 - Record all attributes and hierarchies for each dimension.

Example Information Package

Business Process: Sales Analysis

- **Facts:** Sales revenue, quantity sold, discount amount.
- **Dimensions:**
 - **Time:** Year, quarter, month, day.
 - **Product:** Product name, category, brand.
 - **Location:** Region, country, state, city.
 - **Customer:** Customer name, segment, loyalty level.

Table Representation of an Information Package

Below is a structured **table format** for representing an information package:

Component	Details
Business Process	Sales Analysis
Facts	Sales Revenue, Quantity Sold, Discount Amount
Dimensions	
Time Dimension	Year, Quarter, Month, Day (Hierarchy: Year → Quarter → Month → Day)
Product Dimension	Product Name, Category, Brand (Hierarchy: Category → Brand → Product Name)

Component	Details
Location Dimension	Region, Country, State, City (Hierarchy: Region → Country → State → City)
Customer Dimension	Customer Name, Segment, Loyalty Level (No hierarchy in this example)
Granularity	Transaction level (e.g., each sale recorded per day, per customer, per product)

Explanation

- **Business Process:** Focuses on the specific process to be analyzed (e.g., "Sales Analysis").
- **Facts:** Contains numerical data or metrics central to the analysis.
- **Dimensions:** Each dimension adds contextual data for the facts, with attributes and hierarchies for drill-down analysis.
- **Granularity:** Specifies the lowest level of detail at which data is stored (e.g., daily sales transactions).

7. Case Study: Dimensional Modeling for a Retail Business

Dimensional modeling is an essential technique for structuring data in a retail business data warehouse to support analytics and reporting. This case study focuses on designing a star schema for a fictional retail business, **ShopSmart**, which wants to analyze sales performance.

Business Requirements

1. **Analyze Sales:**
 - Total revenue by product category, time period, and store location.
 - Number of items sold by product and region.
2. **Track Customer Insights:**

- Customer purchasing patterns.
 - Revenue contribution by customer demographics.
3. **Inventory Monitoring:**
- Stock levels by product category.
 - Sales trends to inform restocking decisions.

Dimensional Model Design

The dimensional model for **ShopSmart** includes the following key components:

1. **Fact Table:** Sales facts with measures.
2. **Dimension Tables:** Descriptive data to provide context for analysis.

Star Schema

1. **Fact Table:** **Sales_Fact**
The fact table contains measures such as sales amount, quantity sold, and links to dimensions.
2. **Dimension Tables:**
 - **Time Dimension:** Organizes sales by day, month, quarter, and year.
 - **Product Dimension:** Contains details about products (e.g., category, brand).
 - **Customer Dimension:** Provides customer demographics (e.g., age, gender, location).
 - **Store Dimension:** Contains information about store locations.

Schema Illustration

Here's a diagrammatic representation of the star schema:



Fact Table: Sales_Fact

Attribute	Description
Sales_ID	Unique identifier for each transaction.
Product_ID	Foreign key to Product_Dim.
Customer_ID	Foreign key to Customer_Dim.
Store_ID	Foreign key to Store_Dim.
Date_ID	Foreign key to Time_Dim.
Sales_Amount	Total revenue for the transaction.
Quantity_Sold	Total quantity of items sold.

Sample Data:

Sales_ID	Product_ID	Customer_ID	Store_ID	Date_ID	Sales_Amount	Quantity_Sold
1	101	201	301	20230101	500	2
2	102	202	302	20230101	200	1

Dimension Tables

1. Time_Dim

Attribute	Description
Date_ID	Unique identifier for the date.
Day	Day of the transaction.
Month	Month of the transaction.
Quarter	Quarter of the year.
Year	Year of the transaction.

Sample Data:

Date_ID	Day	Month	Quarter	Year
20230101	Monday	January	Q1	2023

2. Product_Dim

Attribute	Description
Product_ID	Unique identifier.

Attribute	Description
Product_Name	Name of the product.
Category	Product category (e.g., Electronics).
Brand	Product brand.

Sample Data:

Product_ID	Product_Name	Category	Brand
101	Laptop	Electronics	Dell

3. Customer_Dim

Attribute	Description
Customer_ID	Unique identifier.
Name	Name of the customer.
Gender	Gender of the customer.
Age	Age of the customer.
Location	Location of the customer.

Sample Data:

Customer_ID	Name	Gender	Age	Location
201	John Doe	Male	30	New York

4. Store_Dim

Attribute	Description
Store_ID	Unique identifier.
Store_Name	Name of the store.
Location	Location of the store.

Sample Data:

Store_ID	Store_Name	Location
301	ShopSmart NY	New York

Sample Query and Analysis

Query:

“What were the total sales and quantity sold by product category in January 2023?”

1. Join **Sales_Fact** with **Product_Dim** and **Time_Dim**.
2. Filter by Month = 'January' and Year = 2023.

3. Group by **Category** and calculate SUM of Sales_Amount and Quantity_Sold.

Result:

Category	Total_Sales	Total_Quantity
Electronics	\$700	3

Benefits of Dimensional Modeling

1. **Simplified Queries:** Intuitive structure for analytical queries.
2. **Performance:** Optimized for read-heavy operations.
3. **Scalability:** Handles large datasets effectively.
4. **Flexibility:** Supports drill-down, roll-up, and slicing and dicing.

By using dimensional modeling, **ShopSmart** can efficiently analyze sales trends, understand customer behavior, and optimize operations to improve business performance.

Snowflake Schema and Fact Constellation Example

Dimensional modeling often involves variations of the star schema for structuring data. These include the **snowflake schema** and the **fact constellation schema**, which accommodate more complex relationships between dimensions and facts.

1. Snowflake Schema

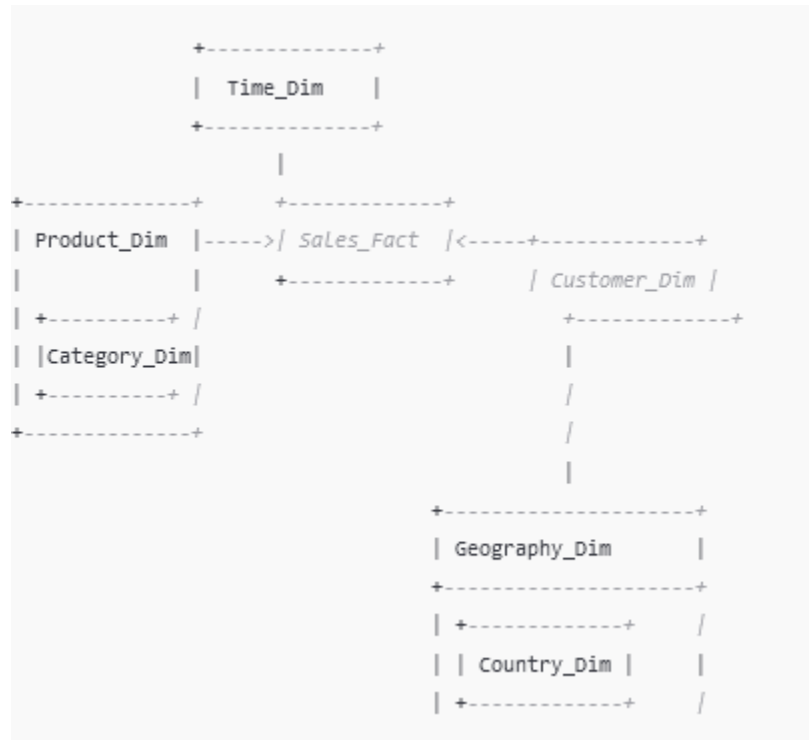
The **snowflake schema** normalizes the dimension tables into multiple related tables to reduce redundancy and save storage space. This design expands dimension tables into hierarchical sub-tables.

Characteristics:

- Normalized dimension tables (1NF or 2NF).
- Data redundancy is reduced but at the cost of more complex queries.
- Suitable for data warehouses with **large, complex dimensions**.

Illustration of Snowflake Schema

Here's a snowflake schema for a retail business:



Explanation:

- **Dimension Normalization:**
 - **Product_Dim** splits into **Category_Dim** for categories.
 - **Geography_Dim** splits into **Country_Dim** for country-level data.
- **Fact Table:** Central table (**Sales_Fact**) contains numerical data and foreign keys to all dimensions.

Example:

Sales_Fact

Sales_ID, Product_ID, Customer_ID, Date_ID, Store_ID, Sales_Amount, Quantity_Sold

Product_Dim

Product_ID, Name

Category_Dim

Category_ID, Cat_Name

Geography_Dim

Region_ID, City

Country_Dim

Country_ID, Country

2. Fact Constellation Schema

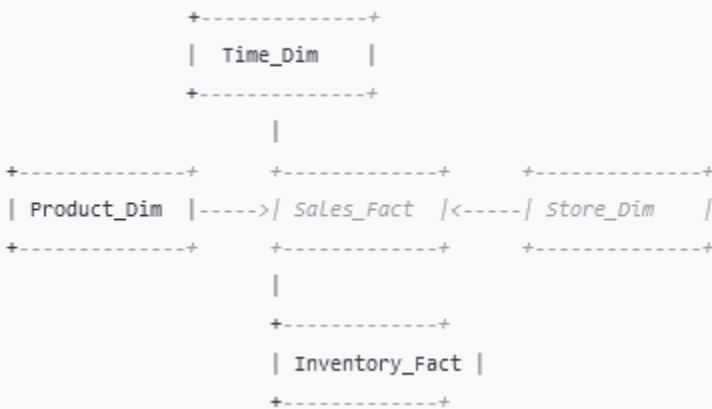
A **fact constellation schema** (also called a **galaxy schema**) is a more complex design that includes multiple fact tables sharing common dimension tables. This approach is suitable for modeling data marts and warehouses with interrelated subject areas.

Characteristics:

- **Multiple fact tables** to represent different business processes.
- **Shared dimensions** to connect fact tables.
- Best suited for **large-scale data warehouses** with interdependent analytical needs.

Illustration of Fact Constellation Schema

Here's a fact constellation schema for a retail business combining **Sales** and **Inventory** analysis:



Explanation:

- **Fact Tables:**
 - **Sales_Fact:** Tracks sales data (e.g., revenue, quantity sold).
 - **Inventory_Fact:** Tracks inventory data (e.g., stock levels, restocking events).
- **Shared Dimensions:**
 - **Product_Dim** and **Time_Dim** are shared across both fact tables.
 - **Store_Dim** links sales and inventory data.

Example:

Sales_Fact
Sales_ID, Product_ID, Store_ID, Date_ID, Sales_Amount, Quantity_Sold
Inventory_Fact
Inventory_ID, Product_ID, Store_ID, Date_ID, Stock_Level, Restock_Quantity
Shared Dimensions
Product_Dim, Store_Dim, Time_Dim

Comparison of Snowflake and Fact Constellation

Aspect	Snowflake Schema	Fact Constellation Schema
Complexity	Moderately complex (normalized).	High complexity (multiple fact tables).
Use Case	Large dimensions, normalized data.	Multiple subject areas (e.g., sales and inventory).
Query Performance	Slower than star schema.	Depends on the number of fact tables.
Data Redundancy	Low (due to normalization).	Medium (shared dimensions reduce redundancy).

When to Use

- **Snowflake Schema:** When dimension tables are large with hierarchies that can benefit from normalization.
- **Fact Constellation:** For interrelated analytical processes requiring multiple facts and shared dimensions.

Both schemas cater to specific needs in a data warehouse, enhancing **flexibility** and **scalability** for complex business environments.

Another Example of Fact Constellation Schema

A **fact constellation schema** (also known as a **galaxy schema**) is a data warehouse schema that involves multiple fact tables sharing common dimension tables. It provides a flexible framework for modeling complex business processes with interconnected analytical needs. This schema is well-suited for large-scale data warehouses that serve multiple subject areas, such as sales, inventory, and marketing.

Key Characteristics of Fact Constellation

1. **Multiple Fact Tables:** Each fact table represents a different business process (e.g., sales, inventory).
2. **Shared Dimensions:** Common dimension tables connect the fact tables.
3. **Granularity:** Fact tables may have different levels of granularity, depending on their purpose.
4. **Flexibility:** Supports cross-functional analyses by leveraging shared dimensions.

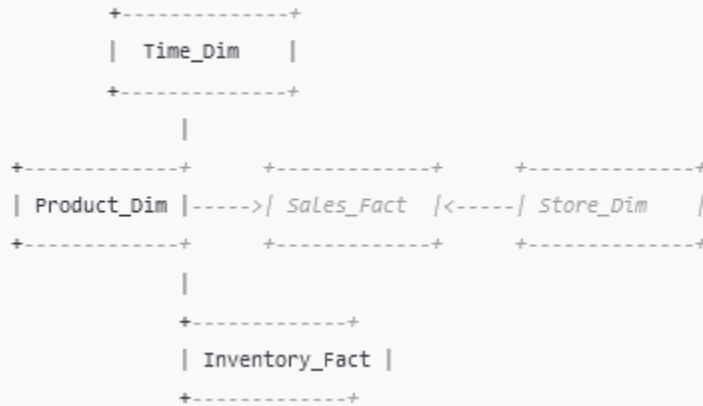
Example: Retail Business Case Study

A retail business wants to analyze:

- **Sales Performance:** Revenue, quantity sold, and discounts.

- **Inventory Management:** Stock levels, restocking events, and shrinkage.

Fact Constellation Schema Illustration



Schema Details

- Fact Tables:**
 - **Sales_Fact:** Tracks transactional sales data.
 - **Inventory_Fact:** Monitors inventory data.
- Shared Dimensions:**
 - **Product_Dim:** Details about products (e.g., name, category).
 - **Store_Dim:** Information about store locations (e.g., city, region).
 - **Time_Dim:** Temporal data (e.g., day, month, year).

Fact Tables

1. Sales_Fact

Attribute	Description
Sales_ID	Unique identifier for each transaction.
Product_ID	Foreign key to Product_Dim.
Store_ID	Foreign key to Store_Dim.
Date_ID	Foreign key to Time_Dim.
Sales_Amount	Total revenue for the transaction.
Quantity_Sold	Total quantity of items sold.

2. Example Data:

Sales_ID	Product_ID	Store_ID	Date_ID	Sales_Amount	Quantity_Sold
1	101	201	20230101	500	2
2	102	202	20230101	200	1

3. Inventory_Fact

Attribute	Description
Inventory_ID	Unique identifier for inventory record.
Product_ID	Foreign key to Product_Dim.
Store_ID	Foreign key to Store_Dim.
Date_ID	Foreign key to Time_Dim.
Stock_Level	Quantity of product in stock.
Restock_Quantity	Quantity restocked during the period.

4. Example Data:

Inventory_ID	Product_ID	Store_ID	Date_ID	Stock_Level	Restock_Quantity
1	101	201	20230101	50	10
2	102	202	20230101	30	5

Dimension Tables

1. Product_Dim

Attribute	Description
Product_ID	Unique identifier.
Product_Name	Name of the product.
Category	Product category.
Brand	Brand of the product.

2. Example:

Product_ID	Product_Name	Category	Brand
101	Laptop	Electronics	Dell

3. Store_Dim

Attribute	Description
Store_ID	Unique identifier.
Store_Name	Name of the store.
Location	Location of the store.

4. Example:

Store_ID	Store_Name	Location
201	ShopSmart NY	New York

5. Time_Dim

Attribute	Description
Date_ID	Unique identifier.
Day	Day of the transaction.
Month	Month of the transaction.
Year	Year of the transaction.

6. Example:

Date_ID	Day	Month	Year
20230101	Monday	January	2023

Analysis with Fact Constellation Schema

- Sales Analysis:**
 - Total sales by product category and store location.
 - Monthly sales trends.
- Inventory Analysis:**
 - Stock levels by product category.
 - Restocking patterns by store.
- Cross-Functional Analysis:**
 - Correlation between sales and restocking.
 - Identify inventory shortages for high-demand products.

Benefits of Fact Constellation Schema

Aspect	Description
Scalability	Supports multiple business processes.
Flexibility	Enables cross-functional analysis.
Data Sharing	Shared dimensions reduce redundancy.
Complexity	Better suited for complex data warehouses.

Note:

The **fact constellation schema** is ideal for retail businesses with interconnected analytics needs, such as combining sales and inventory data. By sharing dimension tables, this schema ensures efficient data organization and supports a wide range of queries for decision-making.