

DATA MODELING

BY - SHUBHAM WADEKAR

Copyright © Shubham Wadekar. All rights reserved.

This material is for personal use only. No part may be copied, shared, resold, or published without written permission. Unauthorized distribution is strictly prohibited and may result in action.

For permissions or licensing, contact: [shubham.p.wadekar@gmail.com]

Disclaimer: This content is for educational purposes. Accuracy is attempted but not guaranteed. No job outcome is promised.

Format of a Data Modeling Interview Round:

During a data modeling interview, the interviewer typically presents a real-world scenario or problem statement that requires a data model. For example, you might be tasked with modeling a situation for a retail store, university, e-commerce website, or a company.

What's Expected:

- **Problem Decomposition:** Break down the problem statement into discernible database entities (tables). Identify columns for each table such that they facilitate easy querying and data retrieval.
- **Dynamic Problem-Solving:** The interview might be interactive. While you might know some questions upfront, in many cases, the interviewer will pose new questions after you've designed your initial model. Sometimes, to answer a new question, you might need to tweak your model. This iterative process tests both your flexibility and comprehensive modeling skills.
- **SQL Queries:** Once your data model is established, be prepared to write SQL queries on top of these entities to answer specific business questions. For instance: "What is the sales amount for each customer till now?"

Tips for Preparation:

- **Practice Real-world Scenarios:** Before the interview, practice modeling various scenarios to get a feel for identifying entities, relationships, and potential challenges.
- **Understand Data Evolution:** Dive deeper into concepts like SCD, as changes and historical data management are often integral to real-world data scenarios.
- **Mock Interviews:** Practice with peers or platforms offering mock interviews. This familiarizes you with the format and enhances your confidence.

The data modeling interview is not just about creating a static model, but showcasing your ability to adjust, refine, and query that model effectively based on evolving requirements.

Data Modeling Concepts

Logical vs. Physical Data Models:

At its core, data modeling is the practice of translating complex real-world scenarios into structured formats for easier data management and processing. Logical data models represent an abstract portrayal of the organizational data, often without considering technological constraints, focusing on capturing information needs and relationships among various data points. Physical data models, on the other hand, delve deeper into the specifics of data storage and access, taking into consideration constraints and features of the database system like indexes, triggers, and stored procedures.

Entities, Attributes, and Relationships:

Entities represent the primary objects or concepts to be stored in a database, such as a “Customer” or “Product”. Each entity is characterized by its attributes, which are the data we want to store for that entity. For instance, a “Customer” entity might have attributes like “Name”, “Address”, and “Phone Number”. Relationships depict how entities are connected to one another. For example, a “Purchase” entity might link a “Customer” to the “Product” they bought, illustrating the relationship between those two entities.

Cardinality and Modality:

Cardinality and modality are concepts that further refine and describe the nature of relationships between entities. Cardinality represents the numerical aspects of the relationship, indicating how many instances of an entity relate to instances of another entity. It could be “one-to-one”, “one-to-many”, or “many-to-many”. Modality, on the other hand, specifies the necessity of the relationship, determining if an entity instance must participate in the relationship or if it’s optional.

ER Diagrams and Their Components:

Entity-Relationship (ER) diagrams provide a graphical representation of the logical data model. They visually depict entities, the attributes of these entities, and how they interrelate. Common components of ER diagrams include rectangles (representing entities), ovals (indicating attributes), diamonds (showing relationships), and connecting lines, often annotated with symbols to represent cardinality and modality constraints. By offering a clear visual layout, ER diagrams assist database designers in understanding, communicating, and implementing the structure of the data they’re working with.

What is an index and why is it important?

A database index is a data structure that improves the speed of data retrieval operations on a database table. Similar to how a book's index helps readers find specific content quickly without reading the entire book, a database index allows the database system to locate desired records without scanning the whole table. Without indexes, the database engine would need to go through each row in the table (a full table scan) to retrieve the relevant records, a process that can be extremely slow for large datasets.

Types of indexes: B-tree, Hash, Bitmap, Composite:

Different indexing methods cater to various query patterns and table structures. The **B-tree** index, with its tree-like structure where each node has a sorted list of keys, is the most common and is especially suited for range queries. **Hash** indexes use a hash function to distribute rows across a set number of buckets and are highly efficient for exact match queries. **Bitmap** indexes use bit vectors and are particularly useful for tables with low cardinality columns (i.e., columns with a limited number of distinct values). Finally, **Composite** indexes include multiple columns within a single index, optimizing queries that filter by those columns simultaneously.

Certainly! Each type of index comes with its own set of advantages and disadvantages, depending on the use case and the specific characteristics of the data.

B-tree:

Pros:

1. Range Queries: B-tree indexes are particularly well-suited for queries that retrieve a range of values since the keys are stored in a sorted manner.
2. Ordered Access: Enables efficient access to data in both ascending and descending order.
3. Dynamic: Can grow and shrink dynamically as data is added or removed.

Cons:

1. Space: Can consume more storage than some other types of indexes, especially if the table has a lot of columns or large values.
2. Maintenance Overhead: While B-trees dynamically adjust, the process of balancing the tree (especially in write-heavy scenarios) can introduce overhead.

Hash:

Pros:

1. Fast Lookups: Hash indexes provide very fast access for exact-match queries.
2. Space-Efficient: Typically, they're more space-efficient than B-trees for exact-match indexing.

Cons:

1. No Range Queries: They're not suited for range queries because hash indexes don't store keys in a sorted manner.
2. Collisions: Hashing can lead to collisions where different keys have the same hash value, requiring additional mechanisms to handle them.
3. Not Adaptive to Growing Data: The initial size of the hash table might not suit the data volume indefinitely, leading to rehashing.

Bitmap:

Pros:

1. Low Cardinality: Ideal for columns with a limited number of distinct values.
2. Combining Indexes: Bitmap indexes can be combined using bitwise operations, which is useful when querying on multiple indexed columns.
3. Space-Efficient: Typically use less space than other index types for low-cardinality columns.

Cons:

1. Write Operations: Bitmap indexes can be slower to update, making them less ideal for write-heavy scenarios.
2. Not Suitable for High Cardinality: For columns with many unique values, bitmap indexes might be inefficient compared to other index types.

Composite:

Pros:

1. Multiple Column Queries: Optimized for queries that filter or sort by multiple columns simultaneously.
2. Reduced Lookups: Can sometimes reduce the need for the database to access the actual table data if the indexed columns provide the necessary information.

Cons:

1. Size: Composite indexes are typically larger than single-column indexes, consuming more storage.
2. Maintenance Overhead: As with B-trees, there's an overhead in maintaining the index when performing write operations.
3. Order Sensitivity: The order of columns in a composite index is crucial. A query filtering by the second column in the index might not utilize the index if the first column isn't specified in the query.

Selecting the right index type depends on understanding the nature of the queries, the characteristics of the data, and the specific requirements of the application.

Benefits of indexing: speed, query performance:

The primary advantage of indexing is speed. With an appropriate index in place, data retrieval can be significantly faster as the system can skip large chunks of data. This speed-up is crucial for databases supporting applications where timely data retrieval is essential. Improved query performance means reduced response times, leading to more efficient applications, especially in scenarios with large data volumes or many concurrent users.

Overhead and trade-offs:

However, indexing is not without its overhead. While read operations become faster, write operations (like INSERT, UPDATE, or DELETE) might become slower because, in addition to updating the table, the associated indexes also need updates. Indexes also consume additional storage space. The decision to index and the choice of indexing method should consider the nature of the application: whether it's read-heavy or write-heavy. Over-indexing can lead to inefficiencies and performance issues, so careful selection and periodic reviews of indexes are essential. Properly maintained, the benefits of enhanced query performance generally outweigh the drawbacks, but as always, it's a balance.

Data Partitioning

Partitioning is a crucial technique in data warehousing for enhancing retrieval performance, managing large datasets, and ensuring scalability. Partitioning involves dividing a large table into smaller, more manageable pieces, yet treating these pieces as a single entity. The partitions can be based on certain column(s) values, such as date or region, depending on the access patterns and data distribution.

Benefits of Partitioning:

1. **Performance Enhancement:** Queries that filter by a partitioning key can access only the relevant partition(s) instead of scanning the entire table. This can significantly improve retrieval times.
2. **Manageability:** Maintaining smaller partitions can be more manageable than a vast, monolithic table. You can efficiently back up, restore, or even archive individual partitions based on need.
3. **Parallelism:** During query execution, multiple partitions can be read concurrently, leveraging parallel processing capabilities of modern systems.
4. **Maintenance Tasks:** Operations like indexing, backups, or even certain updates can be localized to affected partitions, ensuring faster operations and minimizing impact.

Types of Partitioning:

1. **Range Partitioning:** Data is partitioned based on ranges of a column's values, e.g., partitioning sales data by year or month.
2. **List Partitioning:** Data is partitioned based on a predefined list of column values. An example could be partitioning data based on regions such as "North", "South", "East", and "West".
3. **Hash Partitioning:** Data is divided based on a hash value of the partitioning column. This method aims to distribute data evenly across partitions but may not be as intuitive for query optimizations based on specific key values.
4. **Composite Partitioning:** Combining multiple partitioning strategies, such as range-list or range-hash.

Considerations:

1. **Choice of Partition Key:** It's crucial to choose the right column(s) for partitioning. This decision should be based on the primary access patterns, common query filters, and the dataset's nature.
2. **Partition Size:** Partitions should not be too granular, leading to excessive overhead, or too large, negating the benefits.
3. **Maintenance Overhead:** While partitioning can aid in data management, introducing many partitions can also increase complexity, especially in terms of managing schema changes or rebalancing partitions.

Partitioning in a distributed data warehouse like AWS Redshift offers a range of benefits tailored to the unique challenges and opportunities of distributed computing. Here are the primary advantages:

1. **Enhanced Query Performance:** In a distributed system, data is spread across multiple nodes. By partitioning data based on common query patterns, Redshift ensures that each node possesses as much relevant data as possible for a query, minimizing data shuffling between nodes. This co-location of related data is known as “data locality”, which greatly boosts query performance.

2. **Parallel Processing:** Distributed systems like Redshift inherently operate using parallel processing. When data is partitioned appropriately, each node can work on a chunk of data simultaneously, expediting aggregate and computation tasks.

3. **Efficient Data Loads:** Redshift utilizes a columnar storage format. When combined with partitioning, data can be loaded, transformed, and encoded more efficiently. This is especially true if data loads align with partition boundaries, such as loading data for a specific date or region.

4. **Storage Optimization:** Redshift can compress data more effectively when it's partitioned, as each partition often contains similar data. This not only saves on storage costs but also improves query performance since less data needs to be scanned and read into memory.

5. **Data Pruning:** Redshift can skip reading irrelevant partitions based on query predicates. For instance, if partitions are based on date and a query requests data for a specific month, Redshift will only scan partitions relevant to that month, ignoring the rest.

6. **Scalability and Elasticity:** As data grows, the need to scale the infrastructure becomes essential. In Redshift, adding more nodes is straightforward, and if data is partitioned, it can be redistributed across nodes with minimal reshuffling, ensuring consistent performance.

7. **Backup and Recovery:** Partitioned data can be backed up or restored more efficiently. If a particular node or partition faces an issue, it can be addressed in isolation without impacting the entire cluster.

8. **Improved Maintenance Tasks:** Tasks like Vacuuming (reclaiming space from deleted rows and resorting rows) and Analyzing (updating statistics) can be performed more efficiently on partitioned data. Operations can be localized to impacted partitions, saving computational resources and time.

In conclusion, partitioning in a distributed data warehouse like AWS Redshift amplifies the inherent advantages of distributed and columnar systems, ensuring fast, scalable, and cost-effective data operations.

Database Keys for Data Warehousing

Database keys play a crucial role in data warehousing. They help maintain relationships, ensure data integrity, and improve query performance.

1. **Primary Keys:** A primary key is a unique identifier for each row in a database table. It can be a single column or a combination of columns that guarantee the uniqueness of each record. Primary keys are essential for maintaining data integrity, preventing duplicate entries, and providing a unique reference point for other tables.

2. **Foreign Keys:** A foreign key is a column or set of columns that reference the primary key of another table. The purpose of foreign keys is to maintain the relationships between tables and ensure referential integrity. This means that a foreign key value must always match an existing primary key value in the related table or be NULL. Foreign keys help with data consistency, reduce redundancy, and make it easier to enforce constraints and retrieve related information.

Natural and Surrogate Keys

1. **Natural Keys:** Natural keys (also known as business keys or domain keys) are derived from the source system data and are used to identify a record based on its inherent attributes. They can be cryptic or easily understandable, depending on the nature of the data. Natural keys have their limitations, as they can be subject to change, and sometimes they may not guarantee uniqueness. However, keeping natural keys in dimension tables as secondary keys can be beneficial for traceability, data validation, and troubleshooting purposes.

For example, in a table storing employee details, the employee's social security number could serve as a natural key. Similarly, a customer's email address could be a natural key in a customer's table (assuming each customer has a unique email address)

2. **Surrogate Keys:** Surrogate keys are system-generated, unique identifiers that have no business meaning. They are used as primary keys in data warehousing to ensure data integrity and simplify the relationships between tables. Surrogate keys are usually auto-incremented numbers. Surrogate keys are preferred over natural keys for several reasons, including the ability to handle changes in the source data, improved performance, and the support for slowly changing dimensions.

For example, consider a table storing customer details. Even though each customer might have a unique email address, we might choose to use a surrogate key, such as CustomerID, to uniquely identify customers. So, a customer might be assigned an ID like 1, 2, 3, and so on, with no relation to the actual data about the customer.

Comparison:

Meaning: Natural keys have a business meaning, while surrogate keys do not.

Change: Natural keys can change (for example, a person might change their email address), while surrogate keys are

static and do not change once assigned.

Simplicity: Natural keys can be complex if they're composed of multiple attributes, while surrogate keys are simple

(usually just a number).

Performance: Surrogate keys can improve query performance because they're usually indexed and simpler to manage.

Data Anonymity: Surrogate keys provide a level of abstraction and data protection, particularly useful when sharing data without exposing sensitive information.

When designing a data warehouse, it is crucial to make the right choices regarding key usage. As a general guideline, use surrogate keys as primary and foreign keys for better data integrity and handling of data changes. Keep natural keys in dimension tables as secondary keys to maintain traceability and ease troubleshooting. Finally, discard or do not use natural keys in fact tables, as they can lead to redundancy and complexity in managing the data warehouse.

Fact and Dimension tables.

What is a Fact Table?

- A fact table stores the numbers or measurements about a business.
- Think of it as "what happened".
- Example: How many products were sold? What was the total amount?

Date	Product_ID	Store_ID	Units_Sold	Revenue
2024-07-01	101	1	10	500
2024-07-01	102	2	5	250

- Units_Sold and Revenue = these are facts (measurable data).
- Date, Product_ID, and Store_ID = these are keys that link to dimension tables.

Following are the three types of fact tables –

Transactional Fact Table: This is a very basic and fundamental view of corporate processes. It can be used to depict the occurrence of an event at any given time. The facts measure is only valid at that specific time and for that specific incident. "One row per line in a transaction," according to the grain associated with the transaction table. It typically comprises data at the detailed level, resulting in a huge number of dimensions linked with it. It captures the smallest or atomic level of dimension measurement. This allows the table to provide users with extensive dimensional grouping, roll-up, and drill-down reporting features. It's packed yet sparse at the same time. It can also be big at the same time, depending on the number of events (transactions) that have occurred.

Snapshot Fact Table: The snapshot depicts the condition of things at a specific point in time, sometimes known as a "picture of the moment." It usually contains a greater number of non-additive and semi-additive information. It aids in the examination of the company's overall performance at regular and predictable times. Unlike the transaction fact table, which adds a new row for each occurrence of an event, this represents the performance of an activity at the end of each day, week, month, or any other time interval. However, to retrieve the detailed data in the transaction fact table, snapshot fact tables or periodic snapshots rely on the transaction fact table. The periodic snapshot tables are typically large and take up a lot of space.

Accumulating Fact Table: These are used to depict the activity of any process with a well-defined beginning and end. Multiple data stamps are commonly found in accumulating snapshots, which reflect the predictable stages or events that occur over the course of a lifespan. There is sometimes an extra column with the date that indicates when the row was last updated.

What is a Dimension Table?

- A dimension table stores the details or description about things (like product, customer, store).
- Think of it as the "who, what, where, when" part of the data.

Product_ID	Product_Name	Category	Brand
101	iPhone 15	Mobile	Apple
102	Galaxy S24	Mobile	Samsung

You use this to add meaning to your data. So instead of just seeing "Product_ID = 101", you now know it's "iPhone 15".

Simple Analogy:

Imagine a sales report.

- Fact Table = The actual sales numbers (10 iPhones sold, made ₹5000).
- Dimension Table = The details about products, dates, or stores.

Feature	Fact Table	Dimension Table
Contains	Measurements / metrics	Descriptive attributes
Size	Usually large	Smaller
Examples	Sales, Revenue, Quantity	Product, Customer, Time, Store
Type of data	Numeric	Text / descriptive
Use in reporting	Shows what happened	Shows who, what, where, when

What is a Star Schema?

The Star Schema is a widely used data modeling technique in data warehousing and business intelligence (BI) applications. It is designed for optimized query performance in analytical processing by simplifying complex relational models into a denormalized structure.

1. Overview of Star Schema

- **Purpose:** Optimized for read-heavy analytical queries (OLAP - Online Analytical Processing).
- **Structure:** Resembles a star, with a central fact table connected to multiple dimension tables.
- **Denormalized:** Dimension tables contain redundant data to avoid complex joins.
- **Simple & Fast:** Designed for fast aggregation and reporting.

2. Components of Star Schema

A. Fact Table

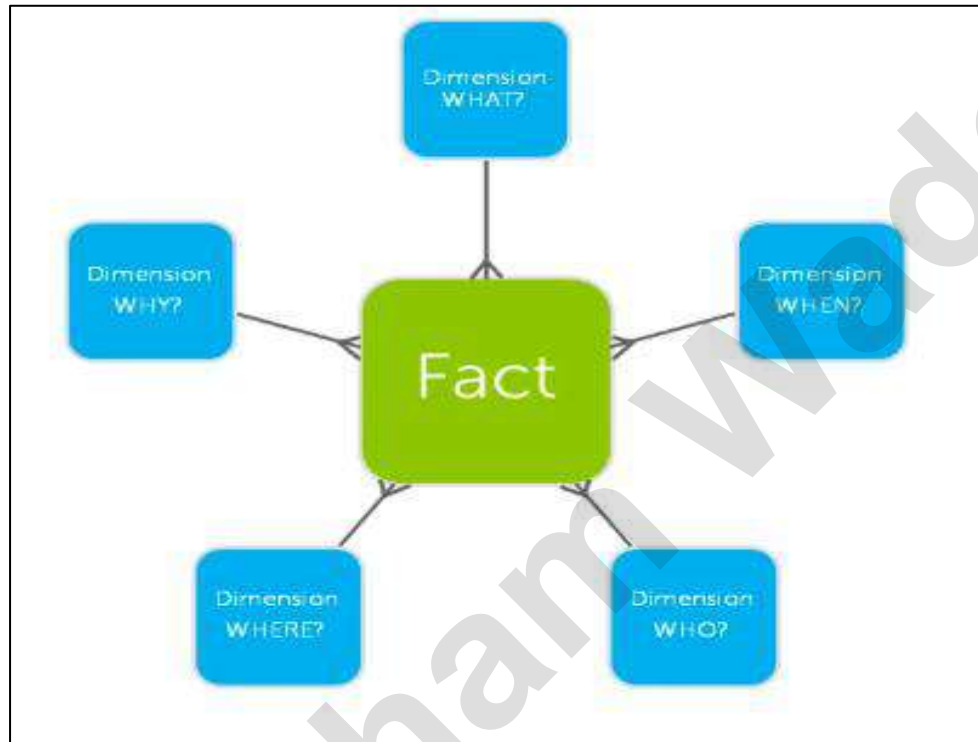
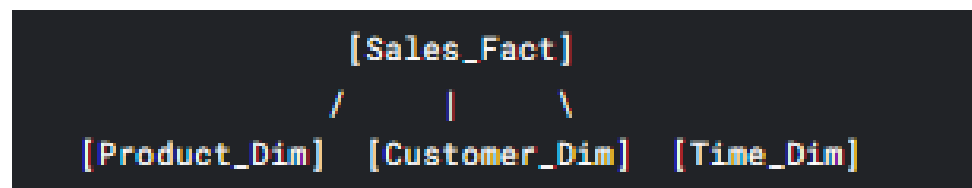
- The central table in the star schema.
- Contains quantitative data (metrics) that can be aggregated (e.g., sales amount, quantity sold).
- Foreign keys link to dimension tables.
- Example Fields
Sales_Fact: sale_id, product_id (FK), customer_id (FK), time_id (FK), amount, quantity.

B. Dimension Tables

- Surround the fact table like the points of a star.
- Contain descriptive attributes used for filtering and grouping.
- Highly denormalized (redundant data to avoid joins).
- Example Dimensions:
 - Product_Dim: product_id (PK), product_name, category, price.
 - Customer_Dim: customer_id (PK), name, address, city.
 - Time_Dim: time_id (PK), date, month, quarter, year.

3. Example of Star Schema

Sales Data Warehouse Example



- **Fact Table (Sales_Fact):**
 - sale_id, product_id, customer_id, time_id, amount, quantity.
- **Dimension Tables:**
 - Product_Dim: product_id, name, category.
 - Customer_Dim: customer_id, name, city.
 - Time_Dim: time_id, date, month, year.

4. Advantages of Star Schema

Fast Query Performance: Fewer joins = faster aggregations.

Simple to Understand: Clear structure for business users.

Optimized for Analytics: Ideal for BI tools (Power BI, Tableau).

Scalable: Handles large datasets efficiently.

Denormalized Structure: Reduces complex joins.

5. Disadvantages of Star Schema

Data Redundancy: Dimension tables store duplicate data.

Not Ideal for OLTP: Designed for analytics, not transactional systems.

Limited Flexibility: Changes in dimensions may require schema updates.

6. When to Use Star Schema?

Data Warehousing (e.g., sales, financial reporting).

Business Intelligence & Dashboards.

Aggregation-Heavy Queries (SUM, AVG, GROUP BY).

When Query Speed is More Important Than Storage Efficiency.

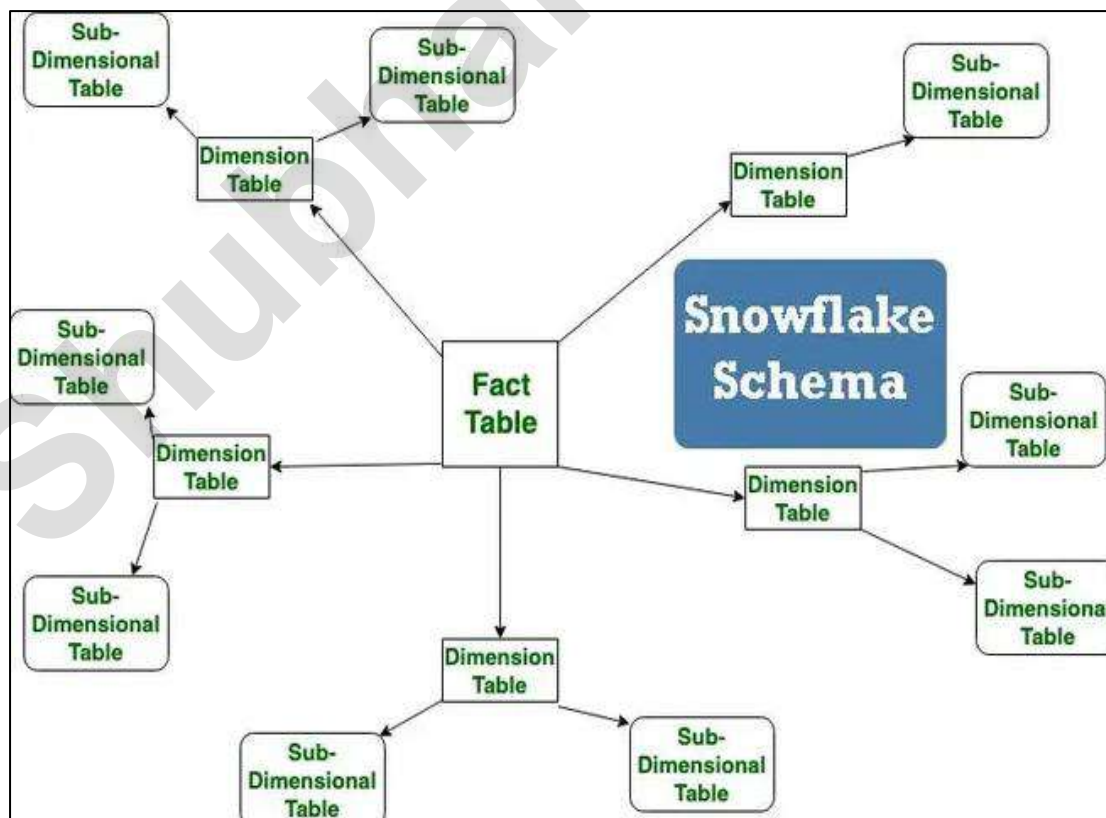
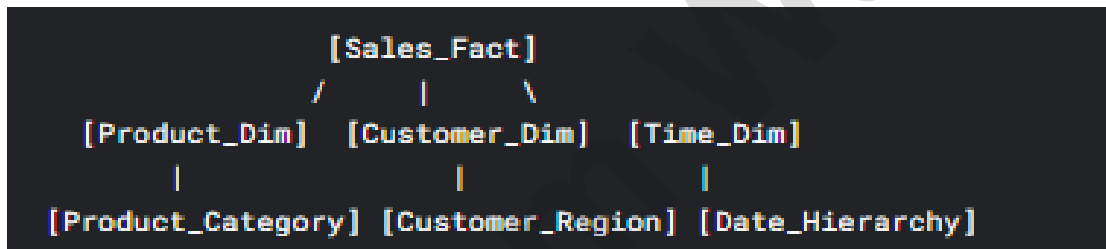
What is Snowflake Schema?

The Snowflake Schema is an extension of the Star Schema, but with normalized dimension tables to reduce redundancy.

Key Features:

- **Normalized Dimensions:** Dimension tables are split into multiple related tables (hierarchical structure).
- **Resembles a Snowflake:** Branches out from the central fact table.
- **Reduced Data Redundancy:** More efficient storage than Star Schema.
- **More Complex Joins:** Requires additional joins, which can impact performance.

Structure:



Example:

- **Fact Table (Sales_Fact):**
 - sale_id, product_id, customer_id, time_id, amount.
- **Dimension Tables:**
 - Product_Dim → Linked to Product_Category.
 - Customer_Dim → Linked to Customer_Region.
 - Time_Dim → Linked to Date_Hierarchy.

Advantages:

Reduced Storage (less redundant data).

Better for Data Integrity (normalized structure).

Easier Maintenance (changes affect fewer tables).

Disadvantages:

Slower Queries (more joins required).

More Complex Design (harder for business users to understand).

When to Use?

When storage efficiency is critical.

When dimensions have hierarchical relationships.

In data warehouses where normalization is preferred.

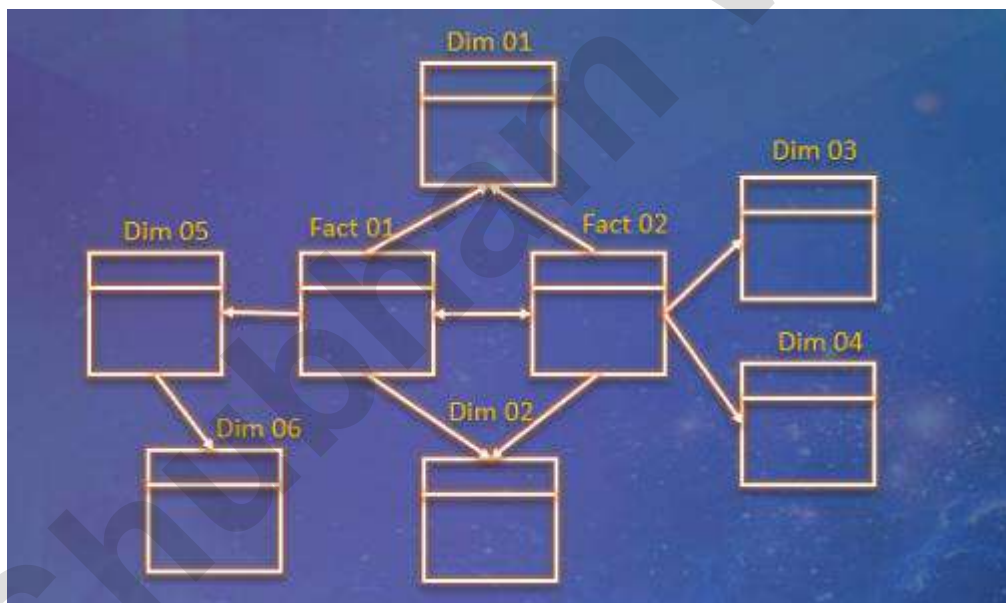
What is Galaxy Schema (Fact Constellation)?

The Galaxy Schema (or Fact Constellation) is a more complex model where multiple fact tables share dimension tables.

Key Features:

- **Multiple Fact Tables:** Different business processes share dimensions.
- **Shared Dimensions:** Avoids duplication (e.g., same Time_Dim for sales and inventory).
- **Highly Flexible:** Supports complex analytics across different business areas.
- **More Complex Than Star/Snowflake:** Requires careful design.

Structure:



Example:

- **Fact Tables:**
 - Sales_Fact (sales data).
 - Inventory_Fact (stock levels).
- **Shared Dimensions:**
 - Time_Dim (used by both fact tables).
 - Product_Dim (links sales and inventory).

Advantages:

Supports Complex Business Processes.
Reduces Redundancy (shared dimensions).
Flexible for Advanced Analytics.

Disadvantages:

Very Complex Design.
Slower Queries (multiple fact tables).
Harder to Maintain.

When to Use?

Enterprise data warehouses with multiple business processes.
When different fact tables need shared dimensions (e.g., sales + inventory).
Advanced analytics requiring cross-functional reporting.

Slowly Changing Dimensions

Slowly Changing Dimensions (SCD) refer to the techniques used in data warehousing to manage and track changes in dimension attributes over time. Given that certain dimension attributes, like a customer's address or a product's description, can change without the need for a new unique identifier, SCDs ensure that historical data remains consistent and accurate. Various SCD types provide strategies to handle these changes, such as overwriting old data, adding new records, or maintaining a history of changes, enabling analysts to observe data trends and changes across timeframes accurately.

Let's delve into each of the SCD types:

Type 0 - The passive method:

Type 0 treats dimensions as static and doesn't track changes. Once an attribute is set, it remains unchanged throughout the dimension's lifetime, regardless of changes in the source system.

Example: Assume a student's major at entry to college is stored. Even if the student changes majors later, the recorded major remains unchanged.

Before:	
Student ID	Major
001	Computer Science
After the student changes to "Mathematics":	
Student ID	Major
001	Computer Science

Type 1 - Overwriting the old value:

With Type 1, when a dimension attribute changes, the old value is simply overwritten with the new value. This method doesn't keep any history of old values.

Example: If a customer changes their address, the new address simply replaces the old one in the customer dimension table.

Before:	
Customer ID	Address
100	123 Main St.
After the customer moves to "456 Elm St.":	
Customer ID	Address
100	456 Elm St.

Type 2 - Creating a new additional record:

In Type 2, when an attribute changes, a new record is added to the dimension table, and the old record is marked as outdated. This approach maintains a full history of attribute values.

Example: When a product's price changes, a new product record is created with the updated price, while the old record is flagged as inactive or is given an end date.

Before:			
Product ID	Price	Active	
P001	\$10	Yes	
After a price change to \$12:			
Product ID	Price	Active	
P001	\$10	No	
P002	\$12	Yes	

Type 3 - Adding a new column:

Type 3 adds a new column to track changes. When an attribute changes, the current value is moved to this new column, and the original column is overwritten with the new value. This method only maintains the previous value.

Example: If an employee's role changes, their previous role is stored in a "Previous Role" column, and the main "Role" column is updated with the new role.

Before:		
Employee ID	Role	
E001	Developer	
After changing role to "Manager":		
Employee ID	Role	Previous Role
E001	Manager	Developer

Type 4 - Using historical table:

With Type 4, a separate history table is maintained. Whenever there's a change, the current state of the dimension is written into the history table before the change is applied to the dimension table.

Example: If a store's location changes, before updating the main table, the current store record is pushed into a "Store History" table, preserving all details prior to the change.

Before:

Store ID	Location
S001	Downtown

After relocating to "Uptown":

Main Table:

Store ID	Location
S001	Uptown

History Table:

Store ID	Location
S001	Downtown

Type 6 - Combine approaches of types 1,2,3 (1+2+3=6):

Type 6 is a hybrid approach, blending elements from Types 1, 2, and 3. It allows overwriting certain attributes (Type 1), adding new records for others (Type 2), and creating new columns for some (Type 3).

Example: Consider a salesperson's region. If they change regions, a new record is created (Type 2), their "Previous Region" is updated (Type 3), and some attributes like their contact number might just be overwritten if they get a new one (Type 1).

Before:

Salesperson ID	Region	Previous Region	Contact Number
SP001	East	None	555-1234

After changing to "West" region and new contact number "555-5678":

Salesperson ID	Region	Previous Region	Contact Number
SP001	East	None	555-1234
SP002	West	East	555-5678

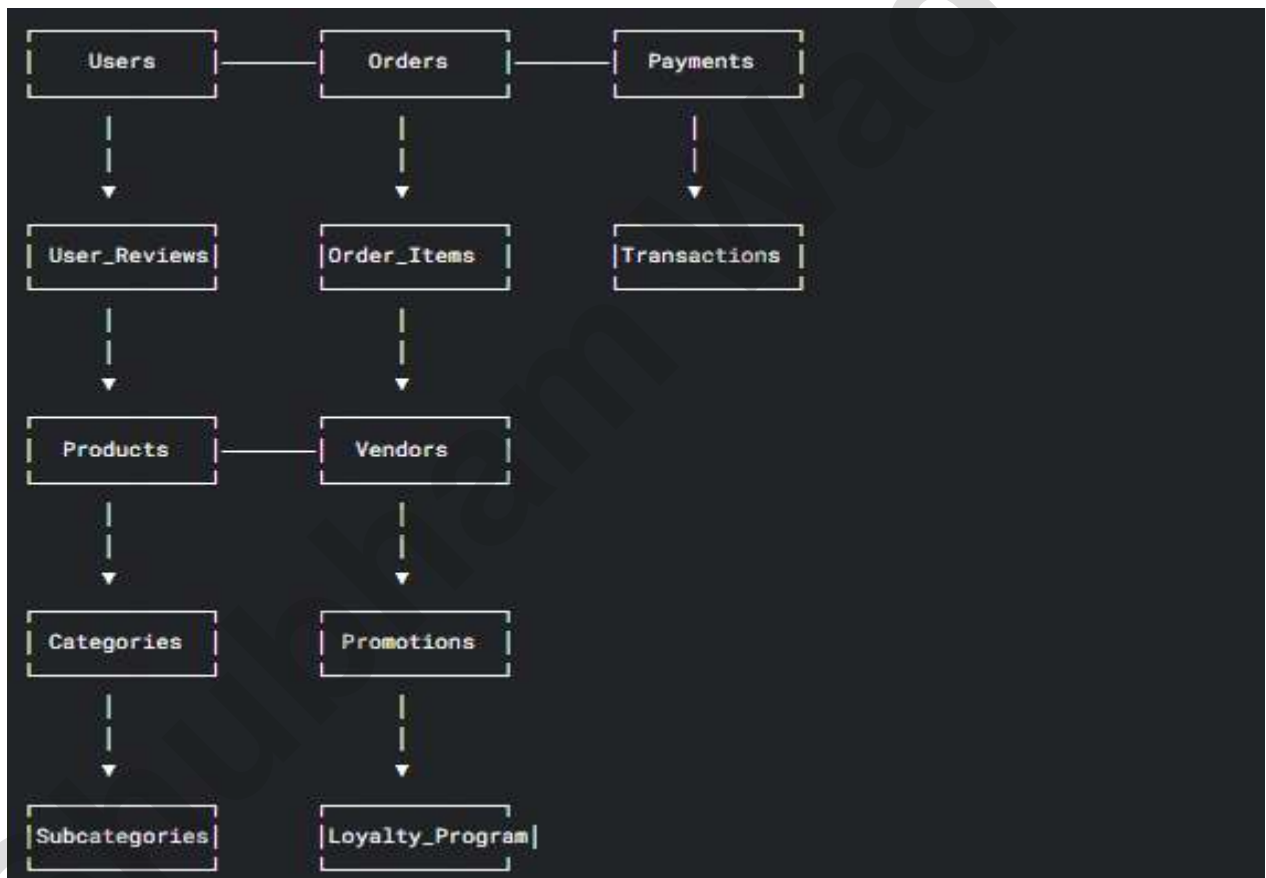
These techniques help businesses choose the most appropriate way to handle historical data, depending on their specific analytical and operational needs.

Design a data model for an e-commerce platform that supports product listings, user reviews, vendor profiles, and order processing. The platform should also cater to promotional offers and loyalty programs.

Overview

This data model supports core e-commerce operations including product management, user reviews, vendor profiles, order processing, promotions, and loyalty programs. It uses a hybrid approach combining relational and dimensional modeling for optimal transactional and analytical performance.

1. Core Entity-Relationship Diagram (ERD)



2. Detailed Tables Structure

A. User Management

```
CREATE TABLE Users (  
    user_id INT PRIMARY KEY,  
    username VARCHAR(50) UNIQUE,  
    email VARCHAR(100) UNIQUE,  
    password_hash VARCHAR(255),  
    first_name VARCHAR(50),  
    last_name VARCHAR(50),  
    registration_date TIMESTAMP,  
    last_login TIMESTAMP,  
    is_vendor BOOLEAN DEFAULT FALSE  
);
```

```
CREATE TABLE User_Addresses (  
    address_id INT PRIMARY KEY,  
    user_id INT REFERENCES Users(user_id),  
    address_type ENUM('home','work','other'),  
    street VARCHAR(100),  
    city VARCHAR(50),  
    state VARCHAR(50),  
    zip_code VARCHAR(20),  
    country VARCHAR(50),  
    is_default BOOLEAN DEFAULT FALSE  
);
```

B. Product Catalog

```
CREATE TABLE Categories (  
    category_id INT PRIMARY KEY,  
    name VARCHAR(50),  
    description TEXT,  
    parent_category_id INT REFERENCES Categories(category_id)  
);
```



```
CREATE TABLE Products (  
    product_id INT PRIMARY KEY,  
    vendor_id INT REFERENCES Vendors(vendor_id),  
    category_id INT REFERENCES Categories(category_id),  
    name VARCHAR(100),  
    description TEXT,  
    sku VARCHAR(50) UNIQUE,  
    price DECIMAL(10,2),  
    stock_quantity INT,  
    weight DECIMAL(10,2),  
    is_active BOOLEAN DEFAULT TRUE,  
    created_at TIMESTAMP,  
    updated_at TIMESTAMP  
);
```

```
CREATE TABLE Product_Images (  
    image_id INT PRIMARY KEY,  
    product_id INT REFERENCES Products(product_id),  
    image_url VARCHAR(255),  
    is_primary BOOLEAN DEFAULT FALSE,  
    sort_order INT  
);
```

C. Vendor Management

```
CREATE TABLE Vendors (  
    vendor_id INT PRIMARY KEY,  
    user_id INT REFERENCES Users(user_id),  
    company_name VARCHAR(100),  
    tax_id VARCHAR(50),  
    business_registration TEXT,  
    avg_rating DECIMAL(3,2),  
    total_products INT DEFAULT 0,  
    joined_date TIMESTAMP  
);
```

```
CREATE TABLE Vendor_Bank_Details (  
    bank_id INT PRIMARY KEY,  
    vendor_id INT REFERENCES Vendors(vendor_id),  
    account_name VARCHAR(100),  
    account_number VARCHAR(50),  
    bank_name VARCHAR(100),  
    branch_code VARCHAR(20),  
    is_verified BOOLEAN DEFAULT FALSE  
);
```

D. Order Processing

```
CREATE TABLE Orders (  
    order_id INT PRIMARY KEY,  
    user_id INT REFERENCES Users(user_id),  
    order_date TIMESTAMP,  
    status ENUM('pending','processing','shipped','delivered','cancelled','refunded'),  
    total_amount DECIMAL(12,2),  
    tax_amount DECIMAL(10,2),  
    shipping_amount DECIMAL(10,2),  
    discount_amount DECIMAL(10,2),  
    payment_method VARCHAR(50),  
    shipping_address_id INT REFERENCES User_Addresses(address_id),  
    tracking_number VARCHAR(100)  
);
```

```
CREATE TABLE Order_Items (  
    order_item_id INT PRIMARY KEY,  
    order_id INT REFERENCES Orders(order_id),  
    product_id INT REFERENCES Products(product_id),  
    quantity INT,  
    unit_price DECIMAL(10,2),  
    subtotal DECIMAL(12,2),  
    discount_applied DECIMAL(10,2)  
);
```

E. Reviews & Ratings

```
CREATE TABLE User_Reviews (  
    review_id INT PRIMARY KEY,  
    user_id INT REFERENCES Users(user_id),  
    product_id INT REFERENCES Products(product_id),  
    rating TINYINT CHECK (rating BETWEEN 1 AND 5),  
    review_text TEXT,  
    review_date TIMESTAMP,  
    is_verified_purchase BOOLEAN DEFAULT FALSE  
);
```

```
CREATE TABLE Vendor_Reviews (  
    vendor_review_id INT PRIMARY KEY,  
    user_id INT REFERENCES Users(user_id),  
    vendor_id INT REFERENCES Vendors(vendor_id),  
    rating TINYINT CHECK (rating BETWEEN 1 AND 5),  
    review_text TEXT,  
    review_date TIMESTAMP  
);
```

F. Promotions & Loyalty

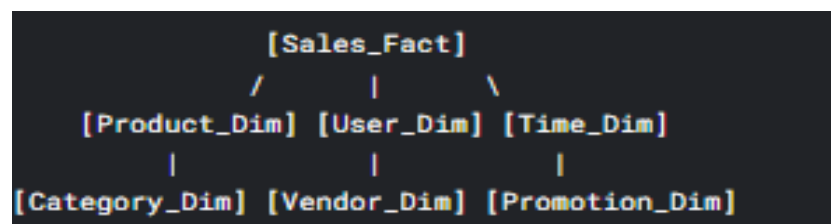
```
CREATE TABLE Promotions (  
    promotion_id INT PRIMARY KEY,  
    name VARCHAR(100),  
    description TEXT,  
    discount_type ENUM('percentage','fixed','buy_x_get_y'),  
    discount_value DECIMAL(10,2),  
    start_date TIMESTAMP,  
    end_date TIMESTAMP,  
    min_order_amount DECIMAL(10,2),  
    is_active BOOLEAN DEFAULT TRUE  
);
```

```
CREATE TABLE Product_Promotions (
    product_promo_id INT PRIMARY KEY,
    product_id INT REFERENCES Products(product_id),
    promotion_id INT REFERENCES Promotions(promotion_id),
    special_price DECIMAL(10,2)
);
```

```
CREATE TABLE Loyalty_Program (
    loyalty_id INT PRIMARY KEY,
    user_id INT REFERENCES Users(user_id),
    points_balance INT DEFAULT 0,
    tier ENUM('bronze','silver','gold','platinum'),
    join_date TIMESTAMP,
    last_activity_date TIMESTAMP
);
```

```
CREATE TABLE Loyalty_Transactions (
    transaction_id INT PRIMARY KEY,
    loyalty_id INT REFERENCES Loyalty_Program(loyalty_id),
    order_id INT REFERENCES Orders(order_id),
    points_earned INT,
    points_redeemed INT,
    transaction_date TIMESTAMP,
    description VARCHAR(255)
);
```

3. Analytical Schema (Star Schema for Reporting)



4. Key Business Rules

1. **Inventory Management:** Stock decreases when order is confirmed
2. **Review Validation:** Only verified purchasers can leave reviews
3. **Loyalty Points:** 1 point per \$1 spent, 100 points = \$1 discount
4. **Promotion Stacking:** Only one promotion applicable per product
5. **Vendor Payouts:** Processed weekly for delivered orders

5. Indexing Strategy

- Products: category_id, vendor_id, sku
- Orders: user_id, status, order_date
- User_Reviews: product_id, rating
- Promotions: is_active, end_date

This model supports:

- High-volume transactions
- Complex product relationships
- Multi-vendor operations
- Real-time analytics
- Personalized promotions

Comprehensive E-Commerce Data Model Explanation

1. Model Overview

This data model is designed for a scalable, multi-vendor e-commerce platform that handles:

- Product catalog management
- User and vendor interactions
- Order processing workflow
- Review and rating systems
- Promotional campaigns
- Loyalty programs

The architecture combines OLTP (transactional) and OLAP (analytical) approaches, using:

- Normalized relational tables for operational data
- Dimensional modeling (star schema) for business intelligence
- Hybrid indexing strategy for performance optimization

2. Detailed Component Analysis

A. User Management System

Core Tables:

1. **Users** (Central identity management)
 - Stores authentication credentials (with password hashing)
 - Tracks user type (customer/vendor) via `is_vendor` flag
 - Includes timestamps for lifecycle management
2. **User_Addresses** (Flexible location handling)
 - Supports multiple address types (home/work/other)
 - Enables default address selection for checkout
 - Geographically structured for shipping calculations

Key Features:

- **Vendor flag** triggers additional profile creation
- **Address normalization** supports international commerce
- **Registration tracking** for user acquisition analytics

B. Product Catalog Structure

Hierarchical Design:

Categories → Subcategories → Products → Product_Variants

Critical Tables:

1. **Categories**
 - Self-referential design for unlimited hierarchy depth
 - Enables faceted navigation (e.g., Electronics > Phones > Smartphones)

2. Products

- Contains both merchandising and logistical attributes
- SKU system supports inventory management
- Active/inactive flag for catalog management

3. Product_Images

- Supports multiple media assets per product
- Sort ordering for gallery display
- Primary image designation for thumbnails

Inventory Considerations:

- Real-time stock tracking via stock_quantity
- Weight data for shipping cost calculations
- Timestamps for product lifecycle analysis

C. Vendor Management Module

Dual-Profile System:

1. Vendors (Business entity)

- Links to user accounts (1:1 relationship)
- Stores legal and operational details
- Maintains performance metrics (avg_rating)

2. Vendor_Bank_Details

- Secure financial information storage
- Verification system for fraud prevention
- Supports multiple account configurations

Business Logic:

- Automatic vendor metrics calculation
- Separate from user profile for security
- Supports marketplace commission models

D. Order Processing Engine

Workflow States:

Pending → Processing → Shipped → Delivered

↓

↓

Cancelled Refunded

Transaction Components:

1. **Orders** (Master record)

- Comprehensive financial breakdown
- Shipping method integration
- Status-driven workflow triggers

2. **Order_Items**

- Snapshot of product details at purchase time
- Supports partial order operations
- Enables discount tracking per item

Financial Integrity:

- Immutable pricing records
- Tax and shipping separation
- Multi-address support (billing/shipping)

E. Review & Rating Ecosystem

Dual Feedback System:

1. **User_Reviews** (Product-centric)

- 5-star rating scale with textual feedback
- Verified purchase flag for authenticity
- Anti-spam measures via purchase validation

2. **Vendor_Reviews** (Seller-centric)

- Separate rating dimension
- Supports marketplace trust systems
- Timestamped for recent activity weighting

Quality Controls:

- Purchase verification requirement
- Rating constraints (1-5 range)
- Temporal analysis capabilities

F. Promotion & Loyalty Framework

Promotion Types Supported:

- Percentage discounts (20% off)
- Fixed amount reductions (\$10 off)
- BOGO deals (Buy 1 Get 1 Free)
- Tiered pricing (Spend \$100, get \$15 off)

Loyalty Program Features:

- Point accrual system (1:1 dollar ratio)
- Tiered membership benefits
- Transaction history tracking
- Reward redemption mechanics

Temporal Controls:

- Campaign scheduling (start/end dates)
- Minimum spend requirements
- Active/inactive toggle

3. Analytical Layer (Star Schema)

Fact Tables:

- **Sales_Fact:** Daily aggregated transactions
 - Measures: revenue, units sold, discounts
 - Dimensions: product, time, promotion

Dimension Tables:

1. **Product_Dim:** Slowly changing dimensions
 - Type 2 SCD for price history
 - Category hierarchies

- 2. **Time_Dim:** Granular date analytics
 - Fiscal periods
 - Holiday flags
- 3. **Promotion_Dim:** Campaign performance
 - Lift calculations
 - ROI analysis

BI Capabilities:

- Cohort analysis
- Promotion effectiveness
- Customer lifetime value
- Inventory turnover

4. Performance Optimization

Indexing Strategy:

Table	Indexed Columns	Type	Purpose
Products	category_id	B-tree	Category browsing
	vendor_id	B-tree	Vendor management
	sku	Hash	SKU lookups
Orders	user_id	B-tree	User history
	status + order_date	Composite	Order monitoring
User_Reviews	product_id + rating	Composite	Product analytics

Partitioning Approach:

- Orders by order_date (monthly ranges)
- User_Reviews by review_date (quarterly)

5. Business Rules Implementation

1. Inventory Management

```
CREATE TRIGGER update_inventory
AFTER INSERT ON Order_Items
FOR EACH ROW
BEGIN
    UPDATE Products
    SET stock_quantity = stock_quantity - NEW.quantity
    WHERE product_id = NEW.product_id;
END;
```

2. Review Validation

```
CREATE FUNCTION is_verified_purchase(user_id INT, product_id INT)
RETURNS BOOLEAN
BEGIN
    RETURN EXISTS (
        SELECT 1 FROM Order_Items oi
        JOIN Orders o ON oi.order_id = o.order_id
        WHERE o.user_id = user_id
        AND oi.product_id = product_id
        AND o.status = 'delivered'
    );
END;
```

3. Loyalty Points Accrual

```
CREATE PROCEDURE process_loyalty_points(IN order_id INT)
BEGIN
    INSERT INTO Loyalty_Transactions
    SELECT
        NULL,
        l.loyalty_id,
        order_id,
        FLOOR(o.total_amount),
        0,
        NOW(),
        CONCAT('Purchase #', order_id)
    FROM Orders o
    JOIN Loyalty_Program l ON o.user_id = l.user_id
    WHERE o.order_id = order_id;

    UPDATE Loyalty_Program lp
    JOIN Loyalty_Transactions lt ON lp.loyalty_id = lt.loyalty_id
    SET lp.points_balance = lp.points_balance + lt.points_earned
    WHERE lt.order_id = order_id;
END;
```

6. Scalability Considerations

Horizontal Scaling:

- Users/Orders: Shard by geographic region
- Products: Shard by category

Vertical Scaling:

- Product search: Elasticsearch integration
- Reviews: MongoDB for unstructured data

Cache Strategy:

- Redis for:
 - Product detail pages
 - Promotional pricing
 - Inventory status

7. Security Implementation

Data Protection:

- Payment info: PCI-compliant tokenization
- Passwords: Argon2 hashing
- PII: Encryption at rest

Access Control:

- RBAC matrix:
 - Customers: Read-only access to products
 - Vendors: CRUD on own products
 - Admins: Full access with audit trails

8. Extension Points

Future Enhancements:

1. Recommendation Engine

- Collaborative filtering tables
- User behavior tracking

2. Subscription Model

- Recurring billing tables
- Membership benefits

3. Multi-Warehouse

- Inventory location tracking
- Regional availability

This model provides a comprehensive foundation for a modern e-commerce platform, balancing operational efficiency with analytical depth while maintaining flexibility for future growth.