Q1) You are hired as a data engineer for ShopSmart, a national retail chain that operates 100+ stores and an online e-commerce platform. ShopSmart wants to build a central analytics warehouse to analyze sales performance, customer behavior, and inventory trends across multiple channels.

# **1** Identify Fact and Dimension Tables

#### **Fact Tables**

<b>Fact Table</b>	Measures / Metrics
Sales_Fact	<pre>sales_amount, quantity_sold, discount_applied, profit</pre>
Inventory_Fact	stock_on_hand, reorder_level, units_sold
Promotion_Fact	<pre>promotion_discount, promotion_revenue, product_count</pre>

#### **Dimension Tables**

Dimension Table	Attributes
Customer_Dim	<pre>customer_id, first_name, last_name, gender, birth_date, loyalty_level, city, state, region</pre>
Product_Dim	<pre>product_id, product_name, category, subcategory, brand, supplier_id, price</pre>
Store_Dim	store_id, store_name, city, state, region, store_manager
Time_Dim	date_key, date, day, week, month, quarter, year, holiday_flag
Promotion_Dim	<pre>promotion_id, promotion_name, start_date, end_date, discount_percent</pre>

# **2** Star Schema Design

**Star Schema:** Fact tables in the center, directly connected to dimension tables.

```
Customer_Dim

|
Product_Dim - Sales_Fact - Time_Dim
|
|
Store_Dim
```

- Fact table: Sales\_Fact
- Dimensions: Customer Dim, Product Dim, Store Dim, Time Dim
- Advantages:
  - Simple structure for fast querying
  - o Denormalized dimensions reduce joins

# **3** Snowflake Schema Design

Snowflake Schema: Dimension tables normalized into sub-dimensions.

```
Customer_Dim

|
|
|
| Sales_Fact
| | \
|
| Product_Dim Store_Dim Time_Dim
|
| Category_Dim
```

- Product Dim normalized to Category Dim
- Store Dim could normalize city/state/region
- Advantages:
  - o Reduces data redundancy
  - Smaller storage footprint
- Disadvantages:
  - $\circ$  More joins  $\rightarrow$  slower query performance

# Slowly Changing Dimensions (SCD)

- Customer Dim: Address or loyalty level may change → Type 2 SCD
- Product Dim: Product category updates → Type 2 SCD
- This ensures **historical analytics** are accurate.

# **5** Example Queries

#### **Total Sales per Region per Month**

```
t.year,
    t.month,
    s.region,
    SUM(f.sales_amount) AS total_sales
FROM Sales_Fact f
JOIN Time_Dim t ON f.time_key = t.date_key
JOIN Store_Dim s ON f.store_key = s.store_id
GROUP BY t.year, t.month, s.region
ORDER BY t.year, t.month, s.region;
```

#### **Top 5 Products by Revenue**

```
SELECT
    p.product_name,
    SUM(f.sales_amount) AS revenue
FROM Sales_Fact f
JOIN Product_Dim p ON f.product_key = p.product_id
GROUP BY p.product_name
ORDER BY revenue DESC
LIMIT 5;
```

#### **Customer Retention Analysis**

```
SELECT
     c.loyalty_level,
     COUNT(DISTINCT f.customer_key) AS active_customers
FROM Sales_Fact f
JOIN Customer_Dim c ON f.customer_key = c.customer_id
WHERE f.time_key BETWEEN '2025-01-01' AND '2025-12-31'
GROUP BY c.loyalty_level;
```

## **6** Justification of Schema Choice

Schema	<b>Use Case</b>	Pros	Cons
Star	Fast reporting, dashboards	Simple queries, good for BI	Slight data redundancy
Snowflake	Normalized data warehouse	Saves storage, reduces redundancy	Complex queries, slower joins

Q2 You are a data engineer for QuickEats, an online food delivery platform operating in multiple cities. QuickEats collects and processes data from multiple sources. Currently, the system struggles with scalability, real-time processing, and analytics performance. Suggest a suitable model.

## **1** Identify Fact and Dimension Tables

#### **✓** Fact Tables

Fact Table Measures / Metrics

Orders Fact: order\_amount, delivery\_fee, discount\_amount, total\_payment

delivery\_time, order\_status

Delivery Fact: delivery\_time, distance\_traveled, pickup\_time, drop\_time,

delivery rating

Payment Fact: payment amount, tax amount, commission, payment status

App\_Events\_Fact: click\_count, session\_duration, device\_type, action\_type

#### **Dimension Tables**

**Dimension Table** Attributes

Customer Dim: customer id, name, phone, email, signup date, city, loyalty level

Restaurant Dim: restaurant\_id, restaurant\_name, cuisine\_type, city, rating,

partner\_since

DeliveryAgent Dim: agent id, agent name, vehicle type, city, experience level

Time Dim: time id, date, day, week, month, quarter, year

Location Dim: location id, city, state, region

PaymentMethod Dim: payment method id, method name, provider

## **2** Star Schema Design (Recommended for BI Dashboards)

#### **✓** Advantages:

- Fast performance for reporting.
- Simple joins.
- Best for dashboards (Power BI, Tableau, Looker).

## 3 Snowflake Schema Design (Normalized)

```
Customer_Dim

|
Restaurant_Dim - Orders_Fact - Time_Dim
|
DeliveryAgent_Dim
|
Location_Dim
|
Region_Dim
```

#### **V** Features:

- Normalized dimensions like Location → Region → Country.
- Reduces duplicate data.

## **Slowly Changing Dimensions (SCD)**

Dimension	Attribute Change Example	SCD Type
Customer_Dim	Loyalty level upgrade from Silver → Gold	Type 2
Restaurant_Dim	Updated rating or menu expansion	Type 2
DeliveryAgent_Dim	Vehicle type change	Type 2
Location_Dim	City name correction	Type 1

## **5** Example Queries

## **✓** Total Revenue by City

```
SELECT 1.city, SUM(o.total_payment) AS revenue
FROM Orders_Fact o
JOIN Location_Dim 1 ON o.location_key = 1.location_id
GROUP BY 1.city
ORDER BY revenue DESC;
```

## **✓** Top 5 Restaurants by Orders

SELECT r.restaurant\_name, COUNT(o.order\_id) AS total\_orders FROM Orders\_Fact o
JOIN Restaurant\_Dim r ON o.restaurant\_key = r.restaurant\_id
GROUP BY r.restaurant\_name
ORDER BY total\_orders DESC
LIMIT 5;

## Average Delivery Time by Agent

SELECT d.agent\_name, AVG(f.delivery\_time) AS avg\_delivery\_time
FROM Delivery\_Fact f
JOIN DeliveryAgent\_Dim d ON f.agent\_key = d.agent\_id
GROUP BY d.agent\_name
ORDER BY avg\_delivery\_time;

## **6** Suggested Modern Architecture for QuickEats

**Layer** Technology

Data Ingestion Kafka / AWS Kinesis (Real-time), Airbyte/Fivetran

(Batch)

Data Storage Data Lake (S3/Google Cloud), Warehouse

(Snowflake/BigQuery)

Processing Apache Spark, Flink (Real-time stream processing)

ETL/ELT dbt + Airflow

BI Dashboard Power BI, Tableau, Looker

Orchestration Apache Airflow

## **7** Justification of Data Model Choice

Schema	<b>Use Case</b>	Advantages	Disadvantages
Star	Sales dashboards, business reporting	Fast queries, easy to maintain	Some redundancy
Snowflake	Large dimensional data	Saves space, normalized	More joins
Data Lakehouse (Recommended)	Real-time + Analytics	Supports batch + streaming, scalable	Slightly complex setup

Q3) You are a data engineer for StreamFlix, a global video streaming platform (like Netflix). StreamFlix collects millions of events per day. The company wants to build a high-performance analytics warehouse to support:

- Real-time viewer engagement analytics
- Top trending videos per region
- Al models for recommendation engines

# **1** Identify Fact and Dimension Tables

#### **✓** Fact Tables

Fact Table Measures / Metrics

Viewership Fact: watch\_duration, progress\_percent, buffering\_time,

resolution played, watch status

StreamingSession\_Fact: session\_duration, device\_time\_spent, data\_consumed\_mb

Subscription\_Fact: subscription\_amount, discount, renewal\_status
Search Fact: total searches, click throughs, search time

Recommendation Fact: recommendation clicks, recommendation impressions

#### **Dimension Tables**

**Dimension Table** Attributes

User Dim: user\_id, name, gender, age\_group, country, subscription\_type,

join date

Video Dim: video\_id, title, genre, sub\_genre, language, release\_year,

maturity\_rating

Device\_Dim: device\_id, device\_type, os, app\_version

Time Dim: time id, date, hour, day, week, month, quarter, year

Location Dim: location id, country, region, city

SubscriptionPlan Dim: plan id, plan name, price, resolution limit, screens allowed

# **2** Star Schema Design (For Fast BI Reporting)

#### Advantages

- Optimized for query performance
- Simple structure for dashboards
- Ideal for daily analytics

# **3** Snowflake Schema Design (Normalized Model)

```
User_Dim

|
Video_Dim - Viewership_Fact - Time_Dim
|
Genre_Dim
|
Device_Dim
|
Location_Dim
|
Region_Dim
```

- ✓ Advantages: Reduces redundancy
- X More joins → Slight slower performance

# Handling Slowly Changing Dimensions (SCD)

Dimension	Change Example	SCD Type
User_Dim	Subscription changes from Basic → Premium	Type 2
Video_Dim	Video updated from SD $\rightarrow$ HD version	Type 2
Device_Dim	App version updates	Type 1
SubscriptionPlan_Dim	Plan pricing updates	Type 2

# **5** Example Analytical Queries

## Top Trending Videos Per Region

```
SELECT l.region, v.title, COUNT(f.video_id) AS total_views
FROM Viewership_Fact f

JOIN Video_Dim v ON f.video_key = v.video_id

JOIN Location_Dim l ON f.location_key = l.location_id

GROUP BY l.region, v.title

ORDER BY total_views DESC

LIMIT 10;
```

## Average Watch Time by Subscription Type

SELECT u.subscription\_type, AVG(f.watch\_duration) AS avg\_watch\_time
FROM Viewership\_Fact f
JOIN User\_Dim u ON f.user\_key = u.user\_id
GROUP BY u.subscription type;

## **Device Usage Analysis**

SELECT d.device\_type, COUNT(f.session\_id) AS sessions
FROM StreamingSession\_Fact f
JOIN Device\_Dim d ON f.device\_key = d.device\_id
GROUP BY d.device type;

# **6** Recommended Modern Architecture (Streaming + Batch)

Layer **Technology** Kafka / AWS Kinesis / **Event Ingestion** Pub/Sub Apache Spark Streaming / **Real-Time Processing** Flink Batch ETL Airflow + dbt Data Lake (S3/GCS) + Delta Storage Lake Snowflake / BigQuery / Warehouse Redshift Feast (for ML Feature Store Recommendations) Tableau / Power BI / Looker BI Layer

## **7** Justification of Schema Choice

Model	<b>Use Case</b>	Pros	Cons
Star Schema	Viewer reports, dashboards	Fast query, simple	Some redundancy
Snowflake	Large metadata handling	Efficient storage	Complex joins
Lakehouse ( 🔽	Real-time + ML use	Supports streaming +	Setup
Recommended)	cases	batch + AI	complexity