

**Tiger
Analytics**

Data Engineering

Interview Questions



Ankita Gulati

Shubh Goyal



Job Details

- **Position:** Data Engineer
- **Experience:** 4+ years
- **Location:** Bangalore
- **Work mode:** Hybrid
- **Compensation:** ₹19–22 LPA
- **Total Rounds:** 3
- **Top Required Skills:**
 1. SQL
 2. PySpark / Python
 3. Cloud Data Engineering
 4. ETL / Data Modeling
 5. Big Data & Streaming
 6. System Design

Round 1

Core Technical Assessment

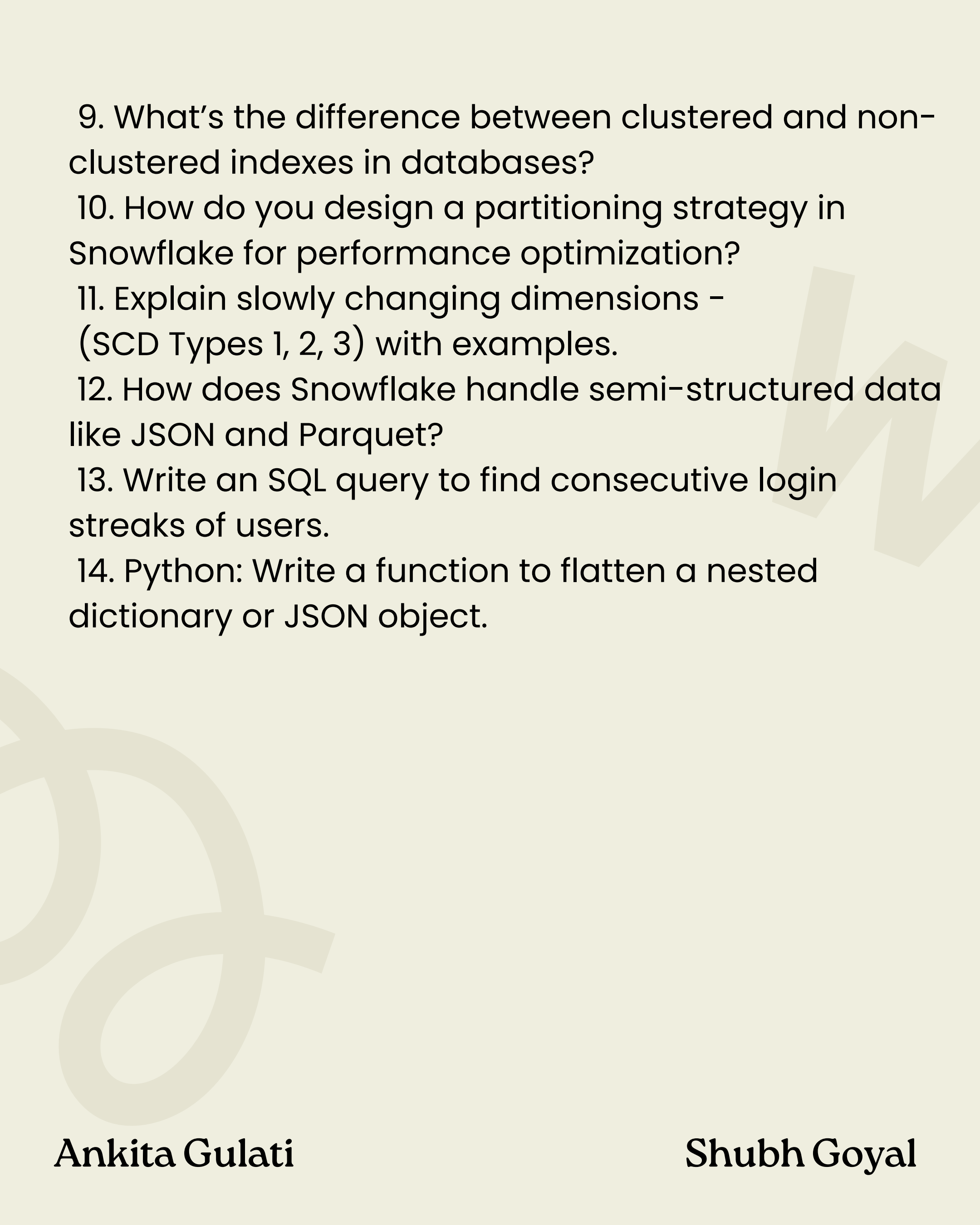
1. Solve 8 SQL problems (easy → medium). Topics include filtering, aggregation, subqueries, and joins.
2. Write an SQL query to find the N-th highest salary from an employee table.
3. Write an SQL query to identify the employee-manager hierarchy in an organization.
4. Explain the differences between `RANK()`, `DENSE_RANK()`, and `ROW_NUMBER()` with examples.
5. Write a query to get the first and last purchase date per customer from a transactions table.
6. Write a query to find duplicate rows in a dataset and count their frequency.
7. Write an SQL query to calculate the running total of sales per region using window functions.
8. Python: Write a function to count the frequency of words in a given list of strings.

9. Python: Write a function to reverse a string without using slicing.
10. SQL: Write a query to calculate the 7-day moving average of sales per product.
11. SQL: Given two tables (Orders, Shipments), identify orders that were placed but never shipped.
12. Python: Write a program to implement a custom iterator that generates prime numbers up to N.

Round 2

Techno-Managerial

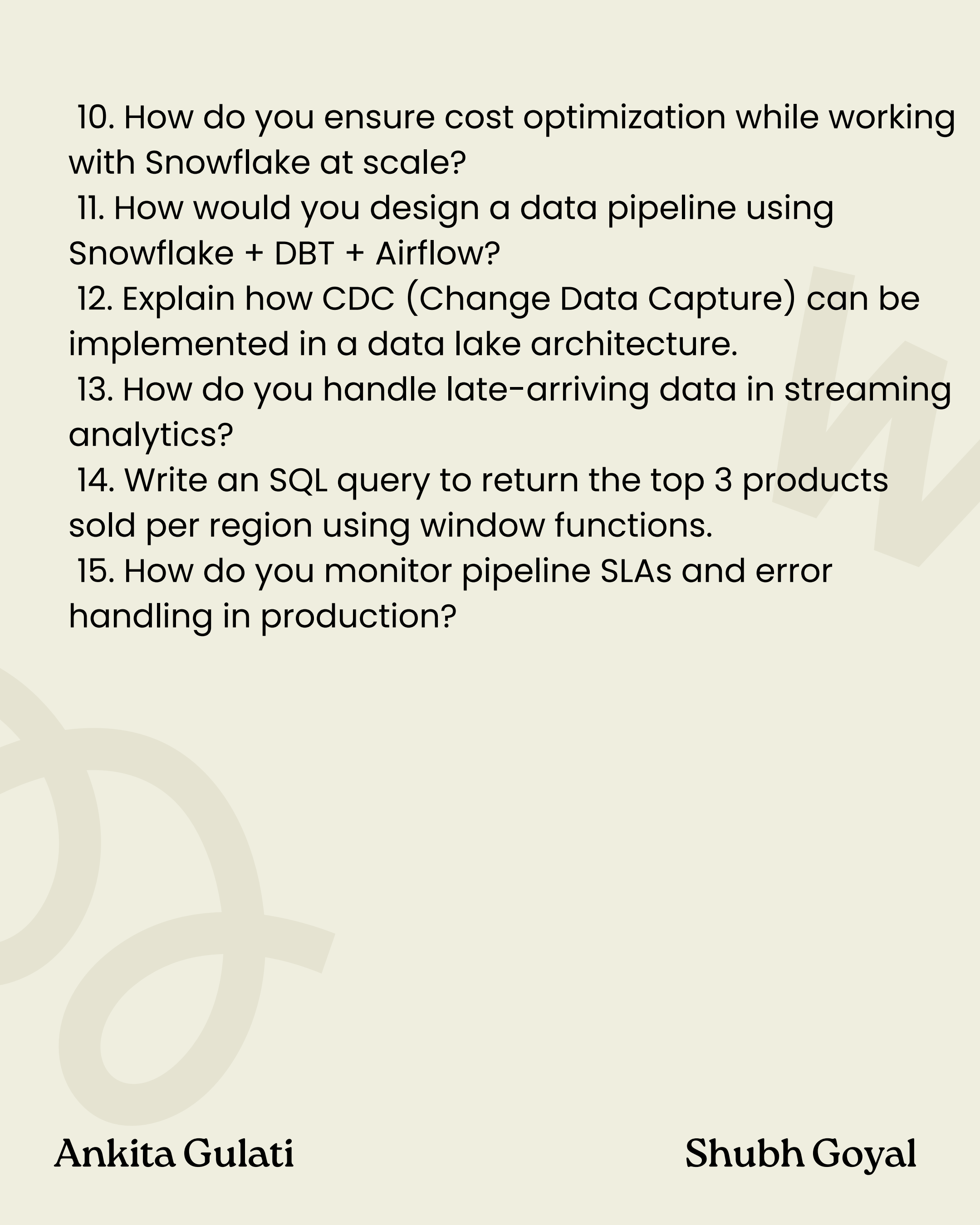
1. What are the core features of Snowflake, and how is it different from traditional warehouses?
2. Can you explain the Medallion Architecture (Bronze, Silver, Gold) and its use in data engineering pipelines?
3. How do you perform SQL query optimization in large-scale systems?
4. Write an SQL query involving a JOIN between two medium-sized tables, ensuring efficiency.
5. Python: Write a function to check if two strings are anagrams of each other.
6. Describe the architecture and implementation of your previous project.
7. What were the biggest challenges you faced in your project, and how did you overcome them?
8. How do you ensure data quality and validation in ETL pipelines?

- 
9. What's the difference between clustered and non-clustered indexes in databases?
 10. How do you design a partitioning strategy in Snowflake for performance optimization?
 11. Explain slowly changing dimensions - (SCD Types 1, 2, 3) with examples.
 12. How does Snowflake handle semi-structured data like JSON and Parquet?
 13. Write an SQL query to find consecutive login streaks of users.
 14. Python: Write a function to flatten a nested dictionary or JSON object.

Round 3

System Design + Conceptual

1. Explain the core principles of Data Warehousing (ETL vs ELT, staging layers, fact vs dimension tables).
2. What are the advantages of Snowflake's separation of storage and compute?
3. Explain how query caching works in Snowflake.
4. Discuss the architecture of your most impactful project and answer follow-up cross-questions.
5. How do you design a real-time data ingestion pipeline using Kafka and Snowflake?
6. How do you handle schema evolution in large data pipelines?
7. How would you design a pipeline that can scale from millions to billions of records per day?
8. What is backpressure in streaming systems and how do you handle it?
9. How would you design a data validation framework for incoming raw data before loading it into Snowflake?

- 
10. How do you ensure cost optimization while working with Snowflake at scale?
 11. How would you design a data pipeline using Snowflake + DBT + Airflow?
 12. Explain how CDC (Change Data Capture) can be implemented in a data lake architecture.
 13. How do you handle late-arriving data in streaming analytics?
 14. Write an SQL query to return the top 3 products sold per region using window functions.
 15. How do you monitor pipeline SLAs and error handling in production?

Thank You

**Best of luck with your
upcoming interviews
— you've got this!**

