

Data Engineer – Proposal Interview TD



01

02

03

Data Pipeline Design

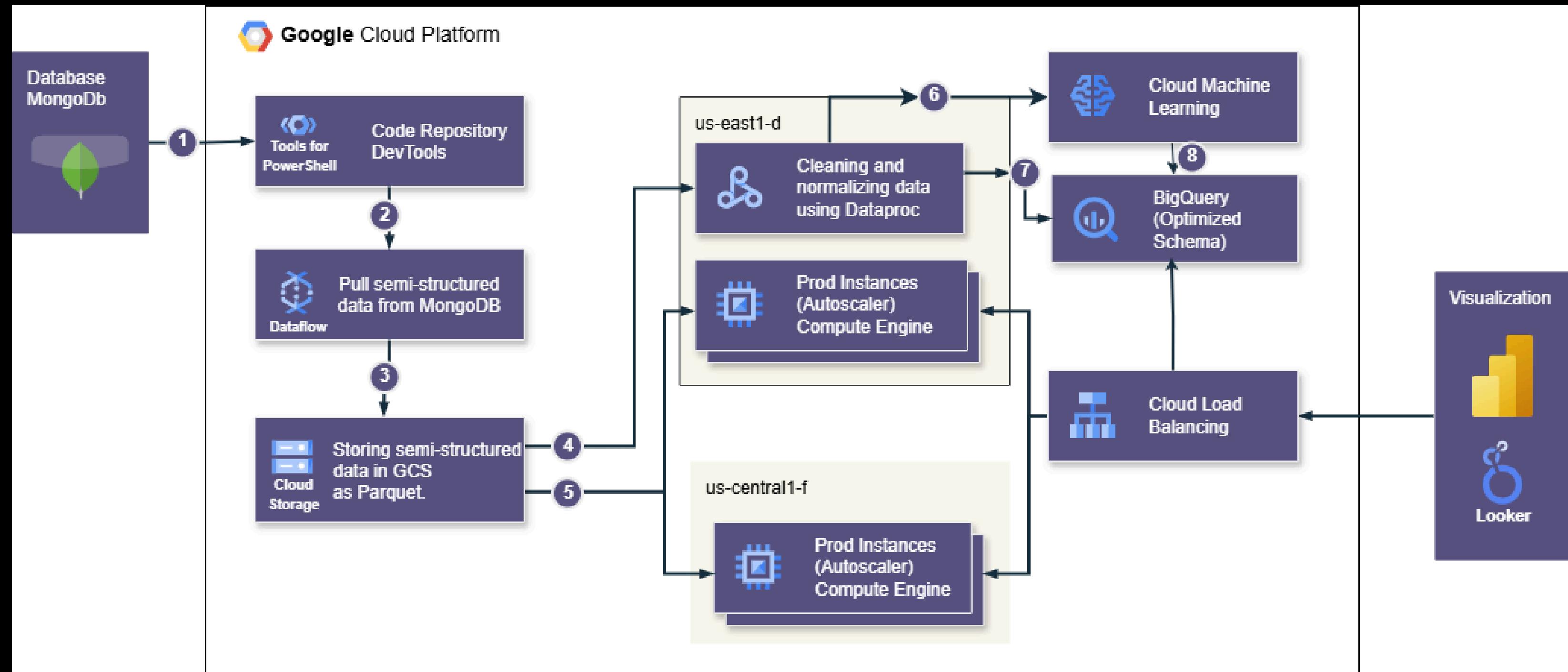
We would like to know your thought process and how you understand
& handle the high level thinking around data.

01

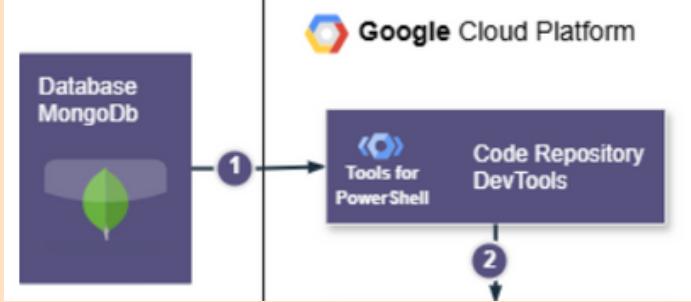
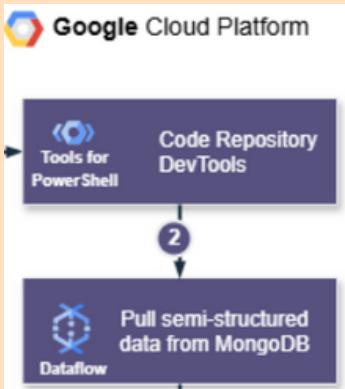
02

03

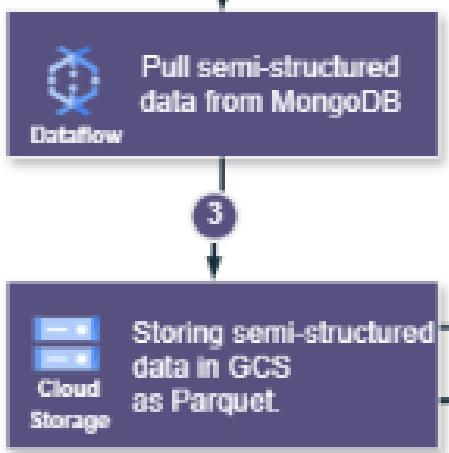
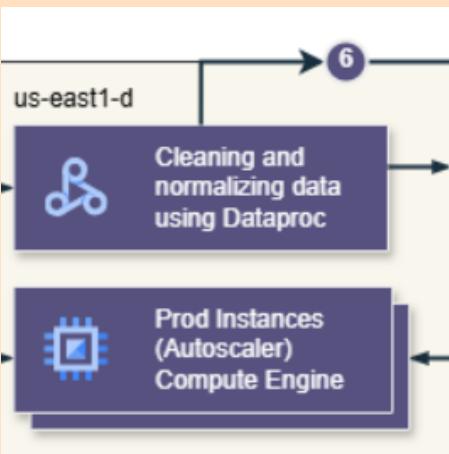
Cloud Architecture



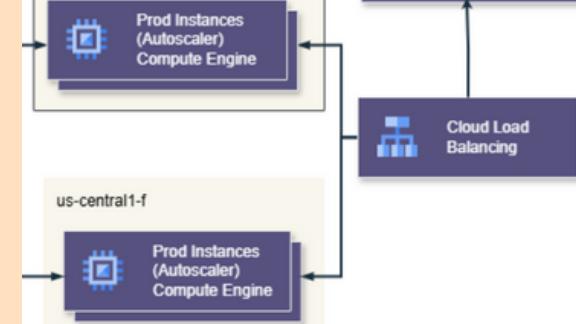
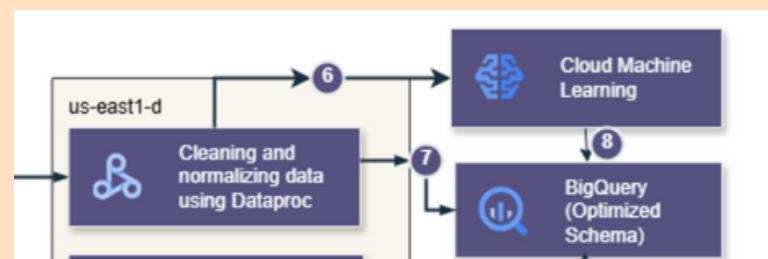
Process & Solution Table

Feature	Process	GCP Service
1.Extract	<ul style="list-style-type: none">PowerShell or automated scripts fetch semi-structured, denormalized data from MongoDB.Code is stored and versioned in a code repository.The script imports the following libraries <pre>1 #Assumption codeing 2 import apache_beam as beam 3 from apache_beam.options.pipeline_options import PipelineOptions 4 import pymongo 5 import json 6 import pandas as pd 7 import pyarrow.parquet as pq 8 from google.cloud import storage, bigquery 9 from pyspark.sql import SparkSession 10 from pyspark.sql.functions import col 11 12 # MongoDB Connection 13 MONGO_URI = "mongodb://your_mongo_uri_path" 14 MONGO_DB = "your_database" 15 MONGO_COLLECTION = "your_collection"</pre>	PowerShell, DevTools 
2.Ingest Data	<ul style="list-style-type: none">Extract semi-structured JSON data from MongoDB.Batch or streaming ingestion into Google Cloud.Handles schema mapping, deduplication, and basic transformations. <pre>17 # GCP Bucket & BigQuery Config 18 GCS_BUCKET = "your-bucket-name_path" 19 GCS_PARQUET_PATH = f"gs://{GCS_BUCKET}/data.parquet" 20 BIGQUERY_DATASET = "your_dataset" 21 BIGQUERY_TABLE = "your_table" 22 23 # Step 1: Extract - Read data from MongoDB 24 def read_from_mongo(): 25 client = pymongo.MongoClient(MONGO_URI) 26 db = client[MONGO_DB] 27 collection = db[MONGO_COLLECTION] 28 return list(collection.find({}))</pre>	Cloud Dataflow (Apache Beam) 

Process & Solution Table

Feature	Process	GCP Service
3.Store Raw Data	<ul style="list-style-type: none">The extracted data is stored in GCS as Parquet format (optimized for analytics).This allows efficient querying & downstream processing. <pre>30 # Step 2: Transform - Convert to DataFrame & Store as Parquet 31 def transform_and_store_parquet(data): 32 df = pd.DataFrame(data) 33 table = pa.Table.from_pandas(df) 34 pq.write_table(table, "/tmp/data.parquet") 35 36 # Upload to GCS 37 storage_client = storage.Client() 38 bucket = storage_client.bucket(GCS_BUCKET) 39 blob = bucket.blob("data.parquet") 40 blob.upload_from_filename("/tmp/data.parquet")</pre>	Cloud Storage (Parquet)  <pre>graph TD A[Pull semi-structured data from MongoDB] --> B((3)) B --> C[Storing semi-structured data in GCS as Parquet]</pre>
4.Transform Data	<ul style="list-style-type: none">Cleans & structures raw data (handling duplicates, missing values, normalization).Runs data transformations (flattening nested JSON, applying business logic).Prepares data for BigQuery. <pre>42 # Step 3: Load - Process Data using Dataproc & Load to BigQuery 43 def clean_data_with_spark(): 44 spark = SparkSession.builder.appName("Data Cleaning").getOrCreate() 45 df = spark.read.parquet(GCS_PARQUET_PATH) 46 47 # Cleaning Steps 48 cleaned_df = df.select(49 col("field1").alias("new_field1"), 50 col("field2").cast("integer"), 51 col("timestamp").cast("timestamp") 52) 53 54 # Save cleaned data back to GCS 55 cleaned_df.write.mode("overwrite").parquet(f"gs://{GCS_BUCKET}/cleaned_data.parquet")</pre>	Cloud Dataproc (Spark)  <pre>graph LR us-east1-d[us-east1-d] --> B((6)) B --> C[Cleaning and normalizing data using Dataproc] C --> D[Prod Instances (Autoscaler) Compute Engine]</pre>

Process & Solution Table

Feature	Process	GCP Service
5. Compute Resources	<ul style="list-style-type: none">Compute instances dynamically scale up/down based on the processing workload.Supports parallel execution of transformation jobs. <pre>68 # Run the Pipeline 69 if __name__ == "__main__": 70 data = read_from_mongo() 71 transform_and_store_parquet(data) 72 clean_data_with_spark() 73 load_to_bigquery()</pre>	Compute Engine (Autoscaler) 
6. Machine Learning (Optional)	<ul style="list-style-type: none">Prepares data for AI/ML model training.Generates predictions or recommendations for analytics.	Cloud Machine Learning 

Process & Solution Table

Feature	Process	GCP Service
7. Load to BigQuery	<ul style="list-style-type: none">Writes structured, partitioned, clustered tables for fast querying.Business-friendly schema for SQL-based users.Enables advanced analytics & reporting. <pre>57 def load_to_bigquery(): 58 client = bigquery.Client() 59 dataset_ref = client.dataset(BIGQUERY_DATASET) 60 table_ref = dataset_ref.table(BIGQUERY_TABLE) 61 62 job_config = bigquery.LoadJobConfig(source_format=bigquery.SourceFormat.PARQUET) 63 uri = f"gs://{GCS_BUCKET}/cleaned_data.parquet" 64 65 load_job = client.load_table_from_uri(uri, table_ref, job_config=job_config) 66 load_job.result()</pre>	BigQuery <pre>graph LR A[us-east1-d] --> B[Cloud Machine Learning] B --> C[Cloud Load Balancing] C --> D[Cleaning and normalizing data using Dataproc] D --> E[Prod Instances (Autoscaler) Compute Engine] E --> F[BigQuery (Optimized Schema)]</pre>
8. Visualization	<ul style="list-style-type: none">Connects BigQuery data to Looker dashboards.Uses Cloud Load Balancing for performance optimization.Enables data-driven decision-making.	Looker, Cloud Load Balancer <pre>graph LR A[Cloud Load Balancing] --> B[Looker] B --> C[Visualization]</pre>

Text Sanitizer

We would like to know how you design the application and how you implement the code based on the given requirements.

01

02

03

Process & Solution Table

Feature	Process	Solution
Upload Input.txt	<p>Step 1: Import necessary libraries Step 2: Read file contents</p> <ul style="list-style-type: none">• Open the uploaded file in read mode ("r"), ensuring encoding is set as "utf-8".• Read and print the content to verify the input.	<pre>from google.colab import files uploaded = files.upload() Choose Files No file chosen Saving Input.txt to Input.txt with open("Input.txt", "r", encoding="utf-8") as file: content = file.read() print(content) Hello, This is a sample text file. It contains multiple lines, different spacing, and some tabs. This should be sanitized correctly! Let's see how well it works. Additional lines are added to test the functionality. Tabs should be replaced with underscores. Line breaks should remain as they are. This is a longer paragraph to check how the text sanitizer handles large chunks of text. It should properly process and clean all input data efficiently, regardless of size. Another test line with multiple words and spaces to see how well extra spaces are handled. Finally, some numbers and symbols: 1234567890 !@#\$%^&()</pre>
Sanitize text	<p>Step 3: Sanitize text</p> <ul style="list-style-type: none">• Initialization• Create an instance of the TextSanitizer class. <p>Step 4: Read Text (read_file)</p> <ul style="list-style-type: none">• Open the provided file and read its content.• Store the text content into an internal variable.	<pre>import string import json class TextSanitizer: """Class for reading, sanitizing, and analyzing text.""" def __init__(self, file_path, output_file="output.json"): """ Initialize the TextSanitizer with a file path. :param file_path: Path to the input text file. :param output_file: Path to the output JSON file. """ self.file_path = file_path self.output_file = output_file self.text = "" self.sanitized_text = "" self.statistics = {} def read_file(self): """Reads input text from a file.""" try: with open(self.file_path, 'r', encoding='utf-8') as file: self.text = file.read() except FileNotFoundError: print(f"Error: File '{self.file_path}' not found.") exit(1) def sanitize_text(self): """Sanitizes the text by converting it to lowercase and replacing tabs.""" self.sanitized_text = self.text.lower().replace('\t', ' ') def generate_statistics(self):</pre>

Process & Solution Table

Feature	Process	Solution
Sanitize text	<p>Step 5: Sanitize Text</p> <ul style="list-style-type: none">• Replace tabs (\t) with underscores (_).• Convert text to lowercase (implied by counting only lowercase characters later).• Retain line breaks as they are. <p>Step 6: Analyze Text</p> <ul style="list-style-type: none">• Calculate the frequency of each alphabetic character (a to z).• Store these counts in a dictionary. <p>Step 7: Output Results</p> <ul style="list-style-type: none">• Save sanitized text and character frequency statistics into a JSON-formatted output file (output.json by default).• Print sanitized text to the console.• Print frequency of each alphabet character.	<pre>def output_result(self): """Prints sanitized text and statistics or writes to a file.""" output_data = { "sanitized_text": self.sanitized_text, "statistics": self.statistics } # Save output to a JSON file with open(self.output_file, 'w', encoding='utf-8') as file: json.dump(output_data, file, indent=4) print(f"Output saved to '{self.output_file}'") # Print results print("\n• **Sanitized Text:**") print(self.sanitized_text) print("\n• **Character Frequency:**") for char, count in self.statistics.items(): print(f"{char}: {count}") def process(self): """Runs the complete text processing workflow.""" self.read_file() self.sanitize_text() self.generate_statistics() self.output_result() # Run the sanitizer on the uploaded Input.txt file if __name__ == "__main__": file_path = "Input.txt" # Adjust if necessary sanitizer = TextSanitizer(file_path) sanitizer.process()</pre> <p>Output saved to 'output.json'</p> <p>• **Sanitized Text:** hello, this is a sample text file. it contains multiple lines, different spacing, and some tabs. this should be sanitized correctly! let's see how well it works. additional lines are added to test the functionality. tabs should be replaced with underscores. line breaks should remain as they are. this is a longer paragraph to check how the text sanitizer handles large chunks of text. it should properly process and clean all input data efficiently, regardless of size. another test line with multiple words and spaces to see how well extra spaces are handled. finally, some numbers and symbols: 1234567890 !@#\$%^&()</p> <p>• **Character Frequency:** a: 40 b: 7 c: 15 d: 23 e: 61 f: 9 g: 5 h: 22 i: 34 j: 0 k: 4 l: 38 m: 8</p>

SQL

We would like to know your thought process and how you implement SQL based on the given requirements.

01

02

03

Process & Solution Table

Feature	Process	Solution																		
SQL Step 1: Aggregate Sales Data	<ul style="list-style-type: none">Retrieve data from three tables:<ul style="list-style-type: none">Sales_Transaction – transaction records with details like quantities and retail prices.Product – provides information on product IDs, product names, and class IDs.Product_class – includes details about product class IDs and their names.Join these tables using the product IDs and transaction IDs to connect transactions with product details, creating a comprehensive dataset containing each sale's product and class information.	<pre>e Blame 37 lines (37 loc) · 985 Bytes 1 WITH SalesData AS (2 SELECT 3 p.product_id, 4 p.product_name, 5 pc.product_class_id, 6 pc.product_class_name, 7 SUM(s.quantity) AS total_quantity, 8 SUM(s.quantity * s.retail_price) AS sales_value 9 FROM Sales_Transaction s 10 JOIN Product p ON s.product_id = p.product_id 11 JOIN Product_Class pc ON p.product_class_id = pc.product_class_id 12 GROUP BY p.product_id, p.product_name, pc.product_class_id, pc.product_class_name</pre> <p>Sales Transaction</p> <table><thead><tr><th>transaction_id</th><th>product_id</th><th>quantity</th></tr></thead><tbody><tr><td>1</td><td>1</td><td>5</td></tr><tr><td>1</td><td>2</td><td>7</td></tr><tr><td>2</td><td>3</td><td>1</td></tr><tr><td>3</td><td>2</td><td>3</td></tr><tr><td>..</td><td>..</td><td>..</td></tr></tbody></table>	transaction_id	product_id	quantity	1	1	5	1	2	7	2	3	1	3	2	3
transaction_id	product_id	quantity																		
1	1	5																		
1	2	7																		
2	3	1																		
3	2	3																		
..																		

Process & Solution Table

Feature	Process	Solution																														
<p>SQL Step 2: Rank Products within Each Product Class</p>	<ul style="list-style-type: none">Calculate key metrics for each product:<ul style="list-style-type: none">Total Quantity Sold ($\text{SUM}(s.\text{quantity})$): the total units sold for each product.Total Sales Value ($\text{SUM}(s.\text{quantity} * s.\text{retail_price})$): the total revenue generated by each product.Group these aggregations by product details (<code>product_id</code>, <code>product_name</code>) and class details (<code>product_class_id</code>, <code>product_class_name</code>), producing a summarized view of sales performance per product.	<pre>15 SELECT 16 product_class_id, 17 product_class_name, 18 product_id, 19 product_name, 20 total_quantity, 21 sales_value, 22 RANK() OVER (23 PARTITION BY product_class_id 24 ORDER BY sales_value DESC, total_quantity ASC 25) AS rank 26 FROM SalesData 27)</pre> <p>Product</p> <table><thead><tr><th>product_id</th><th>product_name</th><th>retail_price</th><th>product_class_id</th></tr></thead><tbody><tr><td>1</td><td>aa</td><td>10</td><td>1</td></tr><tr><td>2</td><td>bb</td><td>20</td><td>1</td></tr><tr><td>3</td><td>cc</td><td>30</td><td>2</td></tr><tr><td>..</td><td>..</td><td>..</td><td>..</td></tr></tbody></table> <p>Product Class</p> <table><thead><tr><th>product_class_id</th><th>product_class_name</th></tr></thead><tbody><tr><td>1</td><td>Class A</td></tr><tr><td>2</td><td>Class B</td></tr><tr><td>3</td><td>Class C</td></tr><tr><td>..</td><td>..</td></tr></tbody></table>	product_id	product_name	retail_price	product_class_id	1	aa	10	1	2	bb	20	1	3	cc	30	2	product_class_id	product_class_name	1	Class A	2	Class B	3	Class C
product_id	product_name	retail_price	product_class_id																													
1	aa	10	1																													
2	bb	20	1																													
3	cc	30	2																													
..																													
product_class_id	product_class_name																															
1	Class A																															
2	Class B																															
3	Class C																															
..	..																															

Process & Solution Table

Feature	Process	Solution																								
SQL Step 3: Filter Top 2 Ranked Products	<ul style="list-style-type: none">Apply a ranking algorithm (RANK()) partitioned by product_class_id.Within each product class, products are sorted first by sales value (sales_value DESC) to prioritize products generating the highest revenue.If two products have the same sales value, the secondary criterion of total quantity (total_quantity ASC) is used to break ties, ranking products with fewer quantities sold higher (a typical business decision might be prioritizing high-value sales even if lower in quantity).	<pre>28 SELECT 29 product_class_name, 30 product_id, 31 product_name, 32 total_quantity, 33 sales_value, 34 rank 35 FROM RankedSales 36 WHERE rank <= 2 37 ORDER BY product_class_id, rank;</pre> <p>Expected Output</p> <table><thead><tr><th>product_class_name</th><th>rank</th><th>product_name</th><th>sal</th></tr></thead><tbody><tr><td>Class A</td><td>1</td><td>aa</td><td>1234</td></tr><tr><td>Class A</td><td>2</td><td>bb</td><td>9999</td></tr><tr><td>Class B</td><td>1</td><td>cc</td><td>2500</td></tr><tr><td>Class B</td><td>2</td><td>dd</td><td>2500</td></tr><tr><td>..</td><td>..</td><td>..</td><td>..</td></tr></tbody></table>	product_class_name	rank	product_name	sal	Class A	1	aa	1234	Class A	2	bb	9999	Class B	1	cc	2500	Class B	2	dd	2500
product_class_name	rank	product_name	sal																							
Class A	1	aa	1234																							
Class A	2	bb	9999																							
Class B	1	cc	2500																							
Class B	2	dd	2500																							
..																							

Thank You

FOR YOUR ATTENTION

01

02

03