DATABRICKS SCENARIO BASED Q&A

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.

Scenario 1. How can Azure Databricks be used for real-time data processing?

Sure, I'll explain this based on my understanding and experience.

Azure Databricks supports real-time or near real-time data processing using something called Structured Streaming. This is very useful when we are dealing with data that keeps coming continuously, like logs, sensor readings, transaction records, or user activities. The idea is to process the data as it arrives instead of waiting for a daily or hourly batch.

Let me explain how I would approach this in a real-world scenario.

Let's say I am working on a system that receives transaction data from an e-commerce app. We want to detect fraudulent activity in near real-time.

To implement this, I would follow these steps:

Step 1: Identify the streaming source

First, we need to know where the data is coming from. Usually, real-time data is pushed into a streaming platform like Azure Event Hubs, Kafka, or IoT Hub.

Let's assume we are using Azure Event Hubs here.

Step 2: Set up the Databricks cluster

I would spin up a Databricks cluster and create a notebook where the streaming logic will be written. I'll also make sure that the cluster has the required libraries installed, like Spark Event Hubs connector.

Step 3: Connect to the Event Hub and read the streaming data

In the notebook, I will write code to read the data stream using Spark Structured Streaming. Here's an example code snippet I would use:

```
from pyspark.sql.types import StringType
from pyspark.sql.functions import col
event_hub_config = {
  'eventhubs.connectionString': '<connection_string>',
  'eventhubs.consumerGroup': '$Default'
}

# Reading data stream from Event Hub
raw_stream = spark.readStream \
  .format("eventhubs") \
  .options(**event_hub_config) \
```

.load()

```
# Convert binary 'body' to string

transactions = raw_stream.selectExpr("cast(body as string) as json_data")
```

Step 4: Parse the incoming data

Usually, the incoming data will be in JSON format. So, I will parse the JSON and extract fields like transaction amount, user ID, location, etc.

```
from pyspark.sql.functions import from_json
```

```
from pyspark.sql.types import StructType, StringType, DoubleType
```

```
schema = StructType() \
    .add("transaction_id", StringType()) \
    .add("user_id", StringType()) \
    .add("amount", DoubleType()) \
    .add("location", StringType())

parsed_stream = transactions \
    .withColumn("data", from_json(col("json_data"), schema)) \
    .select("data.*")
```

Step 5: Apply transformation or business logic

Let's say I want to flag any transaction where the amount is above 5000.

```
suspicious_txns = parsed_stream.filter(col("amount") > 5000)
```

I can also add more logic, like checking transaction patterns or comparing with user history.

Step 6: Write the output to a sink

After processing the data, we need to store or send it somewhere. We can write it to Azure Data Lake, a Delta table, or even send it to Power BI for visualization.

Here's how I would write the suspicious transactions to a Delta table:

```
query = suspicious_txns.writeStream \
.format("delta") \
.outputMode("append") \
.option("checkpointLocation", "/mnt/checkpoints/txns") \
.start("/mnt/output/suspicious_txns")
```

The checkpoint directory helps Spark remember what data it has already processed, so it doesn't reprocess in case of failures.

Step 7: Monitor the stream

Databricks provides a streaming UI where I can monitor the real-time job. It shows me the input rate, processing rate, and any errors. I always keep an eye on this to make sure the job is healthy.

Step 8: Handle failures and scaling

If data volume increases, Databricks can automatically scale the cluster to handle the load. I also make sure to implement error handling and retries wherever needed.

In summary, I would use Azure Databricks with Structured Streaming to read data from Event Hubs, process it in real-time using Spark DataFrames, apply logic like filtering or enrichment, and write the results to Delta Lake or other sinks. This allows me to detect events like fraud or system errors instantly, without waiting for batch jobs.

Scenario 2. How can Databricks help detect fraud by analyzing credit card transactions?

This is a very common use case in financial industries, and Databricks is a very good platform for solving this because of its ability to handle large volumes of data and perform both real-time and batch analytics.

If I were asked to design a fraud detection system for credit card transactions using Azure Databricks, I would explain it like this:

Step 1: Understand the data source

Credit card transactions come from different systems and can include details like:

- transaction ID
- card number or user ID
- amount
- timestamp
- location
- merchant name

This data might be stored in:

- A data lake (like Azure Data Lake or Blob Storage)
- Real-time event systems like Azure Event Hubs or Kafka
- Databases like SQL Server or Cosmos DB

Let's say we're getting transaction data both in batch (daily) and in real-time (through Event Hub).

Step 2: Ingest the data into Databricks

For real-time fraud detection, I would use Structured Streaming to read data from Azure Event Hubs.

```
# Example: Reading streaming data from Event Hub
event_hub_config = {
    "eventhubs.connectionString": "<your_connection_string>"
}
raw_stream = spark.readStream \
    .format("eventhubs") \
    .options(**event_hub_config) \
    .load()
```

For batch analysis, I would use scheduled jobs to read historical transaction data from storage.

historical_df = spark.read.format("delta").load("/mnt/data/transactions")

Step 3: Preprocess and clean the data

I would clean the incoming data to make sure it's usable:

- Remove duplicates
- Handle null values
- Convert timestamp and data types
- Parse any JSON or nested structures

Example:

from pyspark.sql.functions import from_json, col

schema = ... # define transaction schema

parsed_data = raw_stream.withColumn("jsonData", from_json(col("body").cast("string"),
schema)).select("jsonData.*")

Step 4: Define fraud detection rules

I would use a combination of rule-based logic and machine learning.

Rule-based checks (simple but effective):

- Transactions over a certain limit, say \$5000
- Transactions from different countries within 5 minutes
- Too many transactions in a short period (like 10 in 1 minute)

Example of flagging high-value transactions:

suspicious_df = parsed_data.filter(col("amount") > 5000)

Behavioral analysis:

We can detect if a transaction is happening in a location where the customer has never shopped before or at an unusual hour.

Step 5: Use historical data to train a machine learning model

In Databricks, we can use **MLflow** and **Spark MLlib** to train a classification model (like logistic regression or random forest) on past data.

For example:

- Input features: time of transaction, location, amount, device used
- Label: is_fraud (1 or 0)

I would prepare the data like this:

from pyspark.ml.feature import VectorAssembler

from pyspark.ml.classification import RandomForestClassifier

```
assembler = VectorAssembler(inputCols=["amount", "hour", "location_code"],
outputCol="features")
final_df = assembler.transform(training_df)

rf = RandomForestClassifier(featuresCol="features", labelCol="is_fraud")
model = rf.fit(final_df)

Then I would use this model in real-time to predict fraud on incoming transactions.
```

Step 6: Send alerts for suspicious transactions

If a transaction is marked as suspicious or fraud, I would write it to a separate Delta table or send an alert (email, API call, or Power BI dashboard).

```
suspicious_df.writeStream \
    .format("delta") \
    .option("checkpointLocation", "/mnt/checkpoints/fraud") \
    .start("/mnt/output/suspicious_transactions")
```

Step 7: Monitor and improve continuously

Fraud patterns change over time. So I would:

- Retrain models regularly
- Monitor false positives and false negatives
- Work with business analysts to fine-tune rules and thresholds

In summary, I would use Azure Databricks to read transaction data in real time, apply basic rules and ML models to detect unusual or suspicious behavior, and take quick action. Databricks makes this easier because it combines data engineering, streaming, and machine learning in one place, which is perfect for fraud detection use cases.

Scenario 3. How can you implement data partitioning in Databricks, and when should you repartition data?

Sure, I'll explain this step by step based on how I have done it in real projects.

Partitioning in Databricks is a way to split large datasets into smaller, manageable parts. This helps to make data processing faster and more efficient. There are two types of partitioning we usually talk about in Databricks:

- 1. Partitioning the data physically when storing it
- 2. Repartitioning data in memory during processing

Let me walk you through both.

Step 1: Partitioning while writing data (physical partitioning)

This means organizing your data into folders based on one or more columns. This is helpful when users are querying only a part of the data, like filtering by a specific date, region, or customer.

For example, I had a case where I was storing customer orders, and the data was growing very fast. To optimize performance, I decided to partition the data by year and month columns before writing to Delta format.

Here's the code I used:

```
df.write \
.format("delta") \
.partitionBy("year", "month") \
.mode("overwrite") \
.save("/mnt/data/customer_orders")
```

This will create subfolders like /year=2025/month=06/ and so on.

Now, when someone queries only June 2025 data, Spark will read just that folder, which makes the query much faster.

Step 2: Repartitioning data during processing (logical partitioning)

Sometimes, the data is not evenly distributed, or you're going to do operations like joins or groupBy. In those cases, I repartition the data to balance the work across all the Spark executors.

For example, if I'm joining two DataFrames on customer_id, I will repartition both DataFrames on that column:

```
df1 = df1.repartition("customer_id")
df2 = df2.repartition("customer_id")
```

This helps to avoid skew and ensures that all partitions do roughly equal work.

Also, before writing data, if I see too many small files are being created, I use coalesce() to reduce the number of partitions:

df.coalesce(5).write.format("delta").save("/mnt/final_output/")

This is useful especially when I'm writing data daily or hourly, and I don't want hundreds of tiny files.

Step 3: When should I repartition the data?

I usually decide to repartition in these situations:

- **Before a join**: To avoid data skew, especially if I'm joining on a large key like customer_id or region_id.
- Before groupBy or aggregation: If one partition has a lot of data and others have very little, repartitioning helps to distribute the load evenly.
- **Before writing**: To control the number of output files. Too many files slow down downstream processes and cost more in storage.
- After reading from a source: Sometimes the default number of partitions is too high or too low. I adjust it based on the size of data and the size of the cluster.

Step 4: How do I know if repartitioning is needed?

I check a few things:

- I use df.rdd.getNumPartitions() to see how many partitions the DataFrame has.
- I use the Spark UI to see if any stage is taking too long or if there's a data skew (some tasks take a long time, others are fast).
- I check if joins or aggregations are taking too long.
- I look at the size and count of files in the output folder.

If partitions are too big (like >1GB), or too small (like <10MB), I optimize them using repartition() or coalesce().

Step 5: Tips I follow

- Use repartition() when I want to increase the number of partitions.
- Use coalesce() when I want to decrease them without full shuffle.
- Don't over-partition small datasets. That can slow things down.
- Partition by columns that are frequently used in filters (like date or region).

In summary, I use data partitioning to speed up queries and manage large datasets better. I use partitionBy() when storing data and repartition() or coalesce() when processing data. I always look at the data size, distribution, and how it will be used downstream to decide when to apply partitioning or repartitioning.

Scenario 4. How do you approach debugging complex Spark applications in Databricks?

When working with large Spark applications in Databricks, debugging can get tricky, especially if the job involves joins, shuffles, and multiple transformations. So I always follow a structured and simple approach to debug and fix issues. Here's how I usually do it, step by step:

Step 1: Reproduce the issue using a small sample

The first thing I do is try to isolate the problem using a small sample of data. This helps me test quickly and avoids wasting cluster resources.

```
sample_df = full_df.limit(1000)
```

Then I run the transformation logic only on this small sample. If the problem still appears, I can focus on just this data without waiting for a full job to finish.

Step 2: Use display() or show() at every step

Instead of chaining all transformations in one cell, I break them down step by step.

For example:

```
step1_df = df.filter("amount > 100")
display(step1_df)
step2_df = step1_df.groupBy("region").count()
display(step2_df)
```

By checking the data at each stage, I can see where something unexpected happens (like nulls, wrong values, or empty DataFrames).

Step 3: Check schema and data types

Many bugs happen due to column name mismatches or wrong data types (like treating a string as an integer). So I always print schema before using the data.

```
df.printSchema()
```

```
df.select("id", "amount", "date").show(5)
```

If I'm doing a join, I confirm that both DataFrames have the join key in the same format (like both are StringType or IntegerType).

Step 4: Use try-except blocks to capture errors

For Python notebooks, I wrap risky logic in try-except blocks. This helps me log the error message clearly.

try:

```
result_df = df1.join(df2, "customer_id", "inner")
except Exception as e:
    print("Join failed:", str(e))
```

This makes it easier to know exactly which part of the code is failing.

Step 5: Check the Spark UI for performance issues

When the code runs but takes too long, I go to the Spark UI from the Databricks job or notebook run. There, I check:

- Are any stages taking too long?
- Are there tasks with too much data? (means skewed partitions)
- Are shuffles taking a lot of time?
- Are there many small or large partitions?

If I see one task taking 10 minutes and others only a few seconds, I know there's a data skew and I might need to repartition the data.

Step 6: Use explain() to see execution plan

When I do joins or aggregations, I always run .explain() on the DataFrame to see the physical plan.

df.explain(True)

This tells me if:

- There's a BroadcastHashJoin (good for small lookups)
- Or a ShuffleHashJoin (expensive, needs repartitioning)

If I want to force a broadcast join, I can use:

from pyspark.sql.functions import broadcast df.join(broadcast(dim_df), "id")

Step 7: Add logging or checkpoints in long workflows

For complex pipelines with multiple steps, I add log messages or print() statements in between.

Example:

```
print("Step 1 done: Filter applied")
print("Step 2 done: GroupBy complete")
```

Also, for streaming jobs, I use to make sure progress is saved and can be resumed.

Step 8: Handle nulls and bad data early

If I suspect the issue is with bad data, I add null checks or filters to prevent errors.

df = df.filter(df["customer_id"].isNotNull())

Sometimes bad records sneak in from source systems, so validating data early helps avoid crashes in later steps.

Step 9: Test joins carefully

Most of the issues I've faced are during joins. I test joins by:

- Checking the count before and after join
- Doing a left anti join to find missing records

df1.join(df2, "id", "left_anti").show()

This helps me confirm if any keys are missing.

Step 10: Use notebook history and revision tracking

In Databricks, I use the notebook revision history to roll back if a recent change breaks the code. It's like a backup system that helps a lot during debugging.

In summary, my approach is:

- Break the code into small steps
- Use sample data to test
- Check schema and nulls
- Use Spark UI and explain plans
- Add logs and try-excepts
- Investigate joins and shuffles carefully

This method helps me debug most Spark issues in Databricks without spending too much time or cluster resources.

Scenario 5. How do you build a CI/CD pipeline for Databricks notebooks?

When I build a CI/CD pipeline for Databricks notebooks, the main goal is to make sure the code is version controlled, tested, and deployed automatically from dev to test to production environments. This makes the development process more reliable and allows teams to collaborate safely.

Here's how I would explain my approach step by step in an interview:

Step 1: Store notebooks in a Git repository (like Azure DevOps or GitHub)

First, I make sure all my notebooks are stored in a Git repository. In Databricks, I can connect my workspace with Git using the built-in Git integration.

So, from the notebook menu, I select:

• Revision history → Git → Link notebook to Git repo

Now, any changes to the notebook can be tracked using Git. I can commit and push changes like normal code.

If I'm working with a team, we use **feature branches** and pull requests to review changes before merging into the main branch.

Step 2: Use Databricks Repos to sync with the Git repo

Inside Databricks, I use **Repos** to clone the Git repository directly into the workspace.

Workspace > Repos > Add Repo > Enter Git URL

This allows me to edit the code inside Databricks and commit back to Git from the notebook UI or using CLI.

Step 3: Add unit tests using Python and pytest

I write Python code and tests in .py files, not just notebooks. For testing, I use pytest to write unit tests for my data transformation functions.

For example, I may have:

transformations.py

```
def clean_data(df):
    return df.dropna()
```

test_transformations.py

```
def test_clean_data(spark_session):
    from transformations import clean_data
    input_df = spark_session.createDataFrame([(1, None), (2, "abc")], ["id", "name"])
    result_df = clean_data(input_df)
    assert result_df.count() == 1
```

This way, I can test Spark code automatically using pytest.

Step 4: Create a CI pipeline in Azure DevOps or GitHub Actions

Now I create a CI pipeline that runs every time I push code or create a pull request. In Azure DevOps, I create a azure-pipelines.yml file.

Basic example:

trigger:

- main

pool:

vmImage: ubuntu-latest

steps:

- task: UsePythonVersion@0

inputs:

versionSpec: '3.x'

- script: |

pip install -r requirements.txt

pytest tests/

displayName: 'Run Pytest Unit Tests'

This ensures that all code changes are tested before they are merged.

Step 5: Use Databricks CLI or REST API to deploy notebooks

For deployment, I use **Databricks CLI** in the CD pipeline. This allows me to upload notebooks, run jobs, or configure clusters automatically.

First, I configure the CLI:

databricks configure -- token

Then, I use CLI commands to deploy:

databricks workspace import_dir ./notebooks /Workspace/Production --overwrite

This uploads all notebooks to the production folder in Databricks.

Step 6: Set up jobs in Databricks to run the notebooks

After deploying, I trigger Databricks jobs either from the pipeline or schedule them to run daily/hourly.

Using REST API or CLI, I can trigger jobs like this:

databricks jobs run-now -- job-id 12345

Or schedule it in the UI.

Step 7: Use environment-specific configurations

In real scenarios, I separate dev, test, and prod environments. Each has its own:

- Cluster
- Storage paths
- Secrets (like database passwords)

I manage this using a config file or Databricks secrets. For example:

```
if env == "dev":
   path = "/mnt/dev_data/"
elif env == "prod":
   path = "/mnt/prod_data/"
```

This makes it easy to promote the same code to different environments without changing the logic.

Step 8: Use MLflow (if doing ML) to track and promote models

If the notebook involves machine learning, I use MLflow for model versioning and deployment. I track experiments in dev, and promote a model to production using MLflow registry.

In summary, my CI/CD process for Databricks includes:

- Storing notebooks and code in Git
- Writing unit tests using pytest
- · Creating pipelines to test and deploy code
- Using Databricks CLI to upload and trigger jobs
- Managing environments using configs or secrets

This helps make the development process faster, more reliable, and team-friendly.

Scenario 6. How do you choose between using Pandas, Dask, and PySpark in Databricks for analysis?

This is a very practical question, and I usually decide between Pandas, Dask, and PySpark based on the size of the data, performance needs, and team requirements. Let me explain how I choose, step by step, based on my real experience.

Step 1: Understand the differences first

- Pandas is best for small data (fits in memory). It's fast and easy to use but runs only on a single machine.
- Dask works like a distributed version of Pandas. It's better for larger data than Pandas, but it's still not as scalable as PySpark.
- PySpark is built for big data. It can handle huge datasets because it runs on a cluster.
 It's slower for small data but the best choice for distributed processing.

Step 2: Decide based on data size

 If the dataset is less than 1 GB and fits comfortably into the driver node's memory, I choose Pandas.

Example:

import pandas as pd

df = pd.read_csv("/dbfs/tmp/small_data.csv")

- If the data is a few GBs but I need more parallel processing than Pandas, I might consider Dask, especially outside Databricks.
- If the dataset is very large (like 50 GB, 100 GB, or more), or coming from distributed sources like Delta Lake or Parquet files on ADLS, I go with PySpark.

Example:

df = spark.read.format("delta").load("/mnt/data/transactions/")

Step 3: Think about the environment - Databricks is optimized for PySpark

Since Databricks is built on top of Apache Spark, I mostly use PySpark for large-scale transformations and aggregations. It's fully supported, and I can take advantage of:

- Spark optimizations like Catalyst and Tungsten
- Built-in connectors to data lakes, ADLS, Delta Lake
- Cluster power to process billions of records

Even if someone is used to Pandas, Databricks provides pandas API on Spark:

import pyspark.pandas as ps

df = ps.read_csv("/dbfs/tmp/medium_data.csv")

This lets users write Pandas-like code while still using Spark under the hood.

Step 4: Think about the type of work

• For exploratory data analysis (EDA) on small samples, I often use Pandas, especially for quick visualizations using matplotlib or seaborn.

Example:

sample_df = df.limit(1000).toPandas()

- For production pipelines, joins, aggregations, and transformations on large data, I always use PySpark.
- For machine learning, I use Pandas or PySpark ML, depending on the data size.

Step 5: Think about the team and maintainability

If the whole team is experienced with Spark and the data is large, then PySpark is the best fit.

If some users are data analysts or scientists familiar with Pandas, then I use pyspark.pandas or let them work on a smaller sampled dataset using to Pandas().

Step 6: Summary of how I choose

Data Size	Tool	When I Use It
< 1 GB	Pandas	Quick analysis, small CSVs, plotting
1–10 GB	Dask or pandas-on-Spark	s Slightly bigger data, but not huge
> 10 GB	PySpark	Production jobs, joins, aggregations
Mixed team users	s pyspark.pandas	When team prefers Pandas but data is big

In summary, in Databricks I mainly use PySpark because it's designed for distributed big data. I use Pandas only for small-scale analysis or quick data exploration. If someone in my team is used to Pandas, I recommend using pyspark.pandas to write Pandas-style code that runs on Spark. Dask is rare in Databricks because Spark is already more powerful and better integrated with the platform.

Scenario 7. How do you ensure data consistency and integrity when performing data transformations in Databricks?

When I work on data transformations in Databricks, especially for production pipelines, it's very important to make sure the data remains accurate, complete, and in the expected format. So I always follow a few simple but effective practices to maintain data consistency and data integrity. Let me explain how I do it step by step:

Step 1: Understand the source data before transforming

Before I do any transformation, I always check:

- Column names and data types
- Null values or missing data
- Duplicate records
- Unique keys or primary identifiers

For example:

```
df.printSchema()

df.select("customer_id").distinct().count()

df.select("customer_id").isNull().sum()
```

This gives me a clear understanding of what I'm dealing with, and I can plan accordingly.

Step 2: Apply validation rules before and after transformations

I add validation checks at different points in the pipeline.

For example, if I expect every transaction to have a positive amount, I add:

```
df = df.filter("amount > 0")
```

I also count records before and after to make sure data is not lost unexpectedly:

```
before_count = df.count()
after_transform_df = df_transformation(df)
after_count = after_transform_df.count()
if before_count != after_count:
    print("Warning: record count changed!")
```

Step 3: Use schema enforcement when writing data

When writing data to Delta tables, I enable schema enforcement to avoid writing unexpected or broken data.

df.write.format("delta").mode("append").option("mergeSchema",
"false").save("/mnt/data/sales")

This way, if the schema changes accidentally (like a new column appears), it will throw an error and stop the job. This protects the table structure.

Step 4: Use Delta Lake for ACID guarantees

In Databricks, I always use Delta Lake format instead of plain Parquet or CSV. Delta Lake gives me:

- ACID transactions (Atomicity, Consistency, Isolation, Durability)
- Schema enforcement
- Time travel for data recovery

So if anything goes wrong during a write, the data is not half-written. It's either fully written or not written at all.

Example:

df.write.format("delta").mode("overwrite").saveAsTable("sales_data")

Step 5: Use merge or upsert carefully to avoid duplicates

When updating existing data, I use MERGE statements with a clear matching condition to avoid duplicate or incorrect records.

MERGE INTO target_table AS target

USING updates_df AS source

ON target.id = source.id

WHEN MATCHED THEN UPDATE SET *

WHEN NOT MATCHED THEN INSERT *

I always check the number of matched and inserted rows using .count() before doing the actual merge.

Step 6: Handle nulls and data type mismatches early

Nulls and bad data types can silently break logic. So I handle them as early as possible:

df = df.fillna({"amount": 0, "category": "unknown"})

df = df.withColumn("amount", df["amount"].cast("double"))

This avoids unexpected failures later in joins or aggregations.

Step 7: Add checkpoints and logging in ETL pipelines

When the pipeline is big or streaming, I use checkpoints and log progress.

```
query = df.writeStream \
    .format("delta") \
    .option("checkpointLocation", "/mnt/checkpoints/") \
    .start("/mnt/output/")
```

This ensures that if the job fails, it can recover from the last good state and no data is lost or duplicated.

Step 8: Perform data quality checks after transformation

After transformation, I always run a few data quality tests, like:

- Check if any critical column has nulls
- Ensure unique keys are still unique
- Check totals or sums (e.g., total revenue is not negative)

Example:

```
assert df.select("order_id").distinct().count() == df.count()
assert df.filter("revenue < 0").count() == 0</pre>
```

Step 9: Version data and use time travel for safety

Delta Lake lets me track changes over time. I can always go back to an older version if something goes wrong.

```
SELECT * FROM sales_data VERSION AS OF 5
```

This gives a backup and recovery option to protect data integrity.

Step 10: Use unit tests for transformation logic

For reusable transformation functions, I write unit tests using pytest.

Example:

```
def test_drop_nulls():
  input_df = spark.createDataFrame([(1, None), (2, "abc")], ["id", "name"])
  result_df = drop_nulls(input_df, "name")
  assert result_df.count() == 1
```

This ensures that logic is working as expected.

In summary, I ensure data consistency and integrity in Databricks by: Validating source data Using Delta Lake for safety and ACID support Applying schema checks Handling nulls and types early Writing test cases and quality checks Monitoring the output carefully These steps help me build safe and reliable pipelines in Databricks.

Scenario 9. How do you ensure your Databricks code is production-grade and collaborative?

When I write code in Databricks, especially for production pipelines, I always follow some simple practices to make sure the code is clean, reusable, testable, and easy for other team members to understand or contribute to. Let me explain step by step how I make my Databricks code production-ready and collaborative.

Step 1: Organize code in separate notebooks or Python files

I try not to put everything in one big notebook. I break the code into smaller parts based on functionality, for example:

- One notebook for reading and validating data
- One for data transformation
- One for writing data to output
- One for job orchestration or scheduling

If it's reusable logic, I write it in .py files and keep them in a shared folder or repo. Then I import those modules in notebooks using %run or standard Python import.

Example:

```
# reusable_functions.py
def clean_data(df):
    return df.dropna()
In notebook:
from reusable_functions import clean_data
df = clean_data(df)
```

Step 2: Use Git for version control

I always connect my Databricks workspace to a Git repository like Azure DevOps or GitHub. That way, every change to notebooks or .py files is tracked. This helps in:

- Collaborating with team members
- Rolling back if something breaks
- Reviewing code through pull requests

In Databricks, I link notebooks with Git using the "Revision History → Git" option or use Repos to sync the whole repo inside the workspace.

Step 3: Follow coding best practices

I write clean, readable code:

- Use clear variable names
- Add comments to explain complex logic
- Avoid hardcoding paths, secrets, or column names
- Use functions or classes for repeated logic

Example:

```
def read_data(path, file_format="delta"):
  return spark.read.format(file_format).load(path)
```

I also avoid doing everything in one giant Spark command. Instead, I break down transformations step by step to make debugging easier.

Step 4: Add error handling and logging

To make code reliable, I always add try-except blocks and log errors.

Example:

```
try:
```

```
df = spark.read.format("csv").load("/mnt/data/")
except Exception as e:
  print(f"Error reading data: {str(e)}")
  raise
```

In production, I often use structured logging or tools like mlflow to track job runs and logs.

Step 5: Parameterize the notebook using widgets

Instead of hardcoding paths or dates, I use widgets so I can pass parameters dynamically.

```
dbutils.widgets.text("input_path", "")
input_path = dbutils.widgets.get("input_path")
```

This makes the same notebook reusable for different inputs.

Step 6: Write unit tests for functions

For transformation logic written in .py files, I write unit tests using pytest. I test things like:

- Null handling
- Schema changes
- Expected output

Example:

```
def test_clean_data():
  input_df = spark.createDataFrame([(1, None), (2, "xyz")], ["id", "name"])
  result_df = clean_data(input_df)
  assert result_df.count() == 1
```

This makes sure nothing breaks when changes are made later.

Step 7: Use jobs and workflows for production runs

For running notebooks in production, I use Databricks Jobs or Workflows. I can:

- Schedule them (daily, hourly)
- Set retry policies
- Chain multiple notebooks with dependencies
- Monitor status and get alerts if something fails

Example setup:

- Step 1: Read and validate data
- Step 2: Clean and transform
- Step 3: Write output to Delta
- Step 4: Send notification

Step 8: Use secrets and configs safely

I never hardcode credentials or keys. Instead, I store them in Databricks Secret Scope and access them like this:

```
dbutils.secrets.get(scope="my-secrets", key="sql-password")
```

For paths or environment settings, I use config files or environment variables to make the code environment-independent.

Step 9: Document the code and project

To make collaboration easier, I add documentation:

- A readme file in Git to explain the purpose, flow, and setup
- Markdown cells in notebooks to explain steps
- Inline comments for complex logic

This helps other developers quickly understand and contribute to the project.

Step 10: Review code through pull requests

Before merging changes, I always create a pull request and ask team members to review it. We check for:

- Code quality
- Test coverage
- Correct use of Spark
- Performance risks

This improves quality and brings fresh eyes to spot bugs.

In summary, to ensure Databricks code is production-grade and collaborative, I:

- Organize code in reusable modules
- Use Git for version control
- Write clean, tested, and parameterized code
- Use error handling and logging
- Automate runs with Databricks Jobs
- Keep secrets safe
- Document everything clearly

This makes the code reliable, easy to maintain, and team-friendly.

Scenario 10. How do you handle data versioning and reproducibility in Databricks projects?

When I work on Databricks projects, especially in a team or production setup, it's very important to keep track of which version of data was used, what code was applied, and make sure the same results can be reproduced later. So I follow a few simple but effective practices to manage data versioning and reproducibility. Here's how I do it step by step:

Step 1: Use Delta Lake for versioned data storage

The first thing I do is always write data in Delta format. Delta Lake automatically maintains a version history of each table.

So whenever I do an INSERT, UPDATE, or DELETE, a new version of the table is created behind the scenes.

Example:

df.write.format("delta").mode("overwrite").save("/mnt/data/customers")

Later, I can go back to an older version of the data using time travel:

Read version 3 of the table

df = spark.read.format("delta").option("versionAsOf", 3).load("/mnt/data/customers")

Or by timestamp:

df = spark.read.format("delta").option("timestampAsOf", "2024-12-01 10:00:00").load("/mnt/data/customers")

This helps me reproduce results exactly as they were at that time.

Step 2: Track code versions using Git

To make sure the code used for transformation is versioned, I always keep all notebooks and Python files in a Git repo (either GitHub, GitLab, or Azure DevOps).

Each commit and tag in Git represents a snapshot of the code at a certain point.

If I need to rerun a job from a week ago, I can simply check out the exact Git commit used that time and rerun it with the same data version.

In Databricks, I connect my workspace to Git by using the Repos feature or attaching notebooks directly to Git.

Step 3: Record metadata about data and code versions

When running production pipelines, I log the following information after each run:

- Data version used
- Input/output paths
- Git commit hash
- Date and time of run
- User or system that ran it

I usually store this in a Delta table called pipeline_run_log or push it to MLflow or another logging system.

Example:

```
run_log = [(job_id, "customers_table", 7, "abc123commit", "2025-07-19 10:00", "success")]
columns = ["job_id", "table", "data_version", "code_version", "run_time", "status"]
spark.createDataFrame(run_log,
columns).write.format("delta").mode("append").saveAsTable("pipeline_run_log")
```

This helps me go back later and know exactly what data and code were used for a particular run.

Step 4: Use notebooks with widgets for parameterization

I always design notebooks so that I can pass parameters like date, data version, or input path using widgets.

```
dbutils.widgets.text("data_version", "")
version = dbutils.widgets.get("data_version")
df = spark.read.format("delta").option("versionAsOf", int(version)).load("/mnt/data/customers")
```

This allows me to re-run the notebook with the same settings as before to reproduce the same result.

Step 5: Versioning models and experiments using MLflow

If the project includes machine learning, I use MLflow (which is built into Databricks) to track:

- Training data version
- Code version (Git hash)
- Model parameters
- Metrics and outputs

This way, I can reproduce the exact training run later.

```
import mlflow
with mlflow.start_run():
    mlflow.log_param("data_version", 7)
    mlflow.log_param("code_version", "abc123commit")
    mlflow.log_metric("accuracy", 0.91)
    mlflow.spark.log_model(model, "model")
```

MLflow also gives me the ability to register, version, and deploy models in a controlled way.

Step 6: Use structured folder naming or table naming conventions

To keep things clear and organized, I sometimes store versioned data in folders or table names like:

- /mnt/data/v1/customers/
- table_name_v3

But usually I prefer Delta Lake's built-in versioning, as it avoids creating too many tables or paths.

Step 7: Automate everything with Jobs or Workflows

In production, I use Databricks Jobs to automate pipeline runs, and each run is logged with metadata.

I tag each run with job ID and parameters so I can easily trace back what happened.

In summary, I ensure data versioning and reproducibility in Databricks by:

- Writing and reading data in Delta format with versioning
- Keeping all code in Git and using commit hashes
- Logging all metadata of each pipeline run
- Using widgets to rerun notebooks with same parameters
- Using MLflow to track model training details
- Structuring folders and table names for clarity
- Using Jobs/Workflows to automate and monitor everything

These practices help me and my team easily reproduce any result, debug issues, and maintain full transparency across the data pipeline.

Scenario 11. How do you handle personally identifiable information (PII) in Databricks to comply with privacy regulations?

When working with PII data in Databricks, I always follow strict security and privacy best practices to make sure we're complying with regulations like GDPR, HIPAA, and other data protection laws. The main goal is to protect sensitive information, avoid exposing personal data unnecessarily, and make sure only authorized users can access it. Here's how I handle PII step by step:

Step 1: Identify and classify PII data

First, I work with data owners and stakeholders to identify which columns are considered PII. Common examples include:

- Names
- Email addresses
- Phone numbers
- Social Security Numbers
- IP addresses
- Credit card numbers

Once I identify them, I tag or document them properly using a data catalog like Unity Catalog, or maintain metadata manually.

Step 2: Restrict access using Unity Catalog and role-based access control (RBAC)

To control who can access PII, I use Unity Catalog in Databricks. It allows me to:

- Assign fine-grained permissions at the table, column, or view level
- Create roles (like Data Analyst, Data Engineer, Compliance Officer)
- Make sure only authorized users can view or query sensitive columns

For example, if the column email is PII, I can hide it from general users:

GRANT SELECT ON TABLE customer_data TO `analyst_group`

REVOKE SELECT(email) ON TABLE customer_data FROM `analyst_group`

This way, users can query the table but won't see PII fields unless they're authorized.

Step 3: Use encryption for data at rest and in transit

In Databricks, data is already encrypted at rest and in transit by default. But I also make sure:

- Data stored in ADLS or Blob Storage is encrypted using customer-managed keys (CMK)
- I use secure endpoints (HTTPS) and mount points with access control
- Secrets like keys and tokens are stored using **Databricks Secret Scope**

For example, accessing credentials:

password = dbutils.secrets.get(scope="prod-secrets", key="sql-password")

Step 4: Mask or tokenize PII before using it in non-production environments

If I need to use PII data in development or testing, I never use the actual data. Instead, I:

- Mask the data by replacing values with random or generic values
- Tokenize it (replace with reversible tokens using a secure service)
- Hash sensitive fields if I only need uniqueness but not actual value

Example of simple masking:

from pyspark.sql.functions import lit

masked_df = df.withColumn("email", lit("masked@example.com"))

Or hashing:

from pyspark.sql.functions import sha2, col

hashed_df = df.withColumn("email_hash", sha2(col("email"), 256))

Step 5: Use views to separate PII from non-PII

If different teams need access to different levels of detail, I create views that only expose the necessary data.

Example:

CREATE VIEW customer_summary AS

SELECT customer_id, country, signup_date

FROM customer_data

So analysts can query useful data without seeing names, emails, or other PII.

Step 6: Audit access to sensitive data

I enable audit logging in Databricks so we can track:

- Who accessed which table
- When they accessed it
- What queries were run

This is important for compliance and internal investigations.

Step 7: Follow data retention and deletion policies

For compliance (like GDPR's "right to be forgotten"), I make sure we:

- Delete PII data when a user requests
- Remove records after a defined retention period (e.g., 90 days)
- Avoid keeping raw logs with PII longer than necessary

I use scheduled jobs or workflows to clean up data as per policy.

Step 8: Anonymize data when possible

If full identity is not needed, I anonymize the data permanently. This can include:

- Aggregating data (e.g., total sales by region)
- Removing direct identifiers
- Applying techniques like k-anonymity or differential privacy (in advanced use cases)

Example of aggregation:

df.groupBy("region").agg({"revenue": "sum"})

This removes the need for individual-level data.

Step 9: Document and train the team

Finally, I document:

- What data is PII
- How it should be handled
- Who has access
- What masking or encryption is applied

And I make sure the team follows data privacy best practices and understands their responsibilities.

In summary, to handle PII in Databricks and comply with privacy laws, I:

- Identify and tag PII fields
- Use Unity Catalog for access control
- Encrypt data and store secrets securely
- Mask or hash data in non-prod environments
- Use views to limit exposure
- Audit and log access
- Follow data retention and deletion policies
- Anonymize when possible
- Educate and document processes

These steps help ensure that we protect personal data, stay compliant, and build trust with customers and stakeholders.

Scenario 12. How do you manage and version control notebooks in Azure Databricks?

Managing and version-controlling notebooks is very important for any serious data project in Databricks, especially when multiple people are working together. It helps track changes, avoid conflicts, and roll back safely when needed. Here's how I handle notebook versioning in Databricks step by step.

Step 1: Use Git integration with Databricks

The best way to version control notebooks is to connect them with a Git provider like GitHub, Azure DevOps, or GitLab. In Databricks, I use the built-in Repos feature to do this.

Steps I follow:

- 1. In the left menu, I go to Repos → Click Add Repo.
- 2. I enter the Git URL (for example, GitHub repo URL).
- 3. Databricks clones the repo into my workspace.
- 4. I can now edit notebooks directly from the repo, commit, pull, and push changes from inside Databricks.

This way, all changes are version-controlled just like normal code.

Step 2: Work in feature branches

I never work directly on the main branch. Instead, I:

- Create a feature branch when I start working on something.
- Make changes to notebooks or .py files.
- Commit the changes with proper commit messages.

Example:

git checkout -b feature/add-customer-cleaning

Then after completing the task, I open a pull request to merge it back to the main branch after review.

Step 3: Use meaningful commit messages and pull requests

I make sure every change I make is clearly explained in the commit message or pull request.

For example:

vbnet

Added data validation logic to clean customer records with missing email or phone

This helps other team members understand what changed and why.

Step 4: Test locally if needed before committing

If the notebooks are using reusable functions or Spark logic, I sometimes convert important parts into .py modules and test them locally before committing.

This helps avoid pushing broken code.

Step 5: Collaborate using pull requests and code reviews

When a pull request is created, my teammates can:

- Review the changes
- Add comments
- Suggest improvements
- Approve before merge

This makes the project collaborative and avoids bugs in production code.

Step 6: Tag or release versions for stable checkpoints

When we hit a stable version of the pipeline or notebook, I create a release tag in Git. This acts like a snapshot of the project at that point in time.

Example:

git tag v1.0.0

git push origin v1.0.0

This makes it easy to roll back or re-run a specific version of the pipeline.

Step 7: Use %run to call shared code across notebooks

Instead of copying logic between notebooks, I keep common code in a shared notebook and call it using %run.

Example:

%run "./utils/common_functions"

This keeps the project clean and avoids duplicate logic.

Step 8: Use notebook revision history (optional backup)

Even if Git is not used, Databricks itself keeps a revision history of notebooks.

- I can go to Notebook → Revisions to view older versions
- I can revert to any earlier version if needed

But this is more of a backup feature and not as powerful as Git.

Step 9: Organize notebooks in a clear folder structure

This helps the team quickly find the right code without confusion.

Step 10: Keep configuration and secrets separate

I keep paths, parameters, and secrets out of the notebooks.

- Use widgets to pass runtime parameters.
- Use secret scopes for credentials.
- Store configurations in a JSON file or table.

This makes the notebooks reusable and safe for versioning.

In summary, I manage and version control notebooks in Azure Databricks by:

- Using Git integration (via Repos or direct Git links)
- Working in feature branches
- Writing clear commit messages and using pull requests
- Using tags to track stable versions
- Organizing code cleanly and calling shared logic with %run
- Keeping configs and secrets out of the code
- Optionally using Databricks' built-in revision history

These steps make the project more collaborative, organized, and easier to maintain across teams and environments.

Scenario 13. How do you manage job retries in Databricks?

Managing job retries in Databricks is important because sometimes jobs fail due to temporary issues like network errors, cluster delays, or data not being available yet. Instead of failing the whole pipeline, we can automatically retry the job a few times before finally failing it. This improves the stability and reliability of the data pipeline.

Here's how I manage job retries step by step:

Step 1: Configure retries in the job settings

When I create a job in Databricks (either through the UI or using the REST API), there's a built-in option to set retry behavior.

In the Databricks UI:

- 1. I go to Jobs → Click Create Job or edit an existing one.
- 2. In the task section, I scroll down to Advanced options.
- 3. Under Retries, I set:

Max Retries – how many times the job should retry if it fails (e.g., 3).

Retry Interval – time between retries (in seconds or minutes).

This helps Databricks automatically re-run the job if it fails due to temporary issues.

Step 2: Understand when retries work best

Retries are helpful when:

- There is a temporary data availability issue (e.g., waiting for an upstream process to finish).
- · There are cluster startup delays.
- A job fails due to a transient error like a connection timeout or driver node restart.

But retries won't help if the job fails due to:

- Bad logic in the notebook (e.g., syntax error)
- Corrupt data
- Permanent infrastructure issues

That's why I monitor the reason for failure before deciding the number of retries.

Step 3: Use conditional logic inside notebooks (optional)

In some cases, instead of letting the whole job retry, I add **custom retry logic** inside my notebook for specific steps like API calls or external service access.

```
Example using Python:

import time

max_retries = 3

for attempt in range(max_retries):

try:

# example: reading from a flaky API or temporary DB

df = spark.read.format("jdbc").load(...)

break

except Exception as e:

print(f"Attempt {attempt+1} failed: {e}")

time.sleep(10)

else:

raise Exception("All retries failed")
```

This helps control retries at a smaller level instead of the entire job.

Step 4: Monitor job failures and retries using logs

After enabling retries, I always make sure to monitor job runs in the Databricks Jobs UI or via alerts.

- I can see which runs were retried and whether the retry was successful.
- If I use Databricks REST API, I can fetch job run history programmatically.
- I also set up email or webhook alerts for failures, so I get notified if all retries fail.

Step 5: Use retry logic in workflows with dependencies

If I have a multi-task workflow, I can also configure retries for each individual task.

Example:

- Task A: Ingest data (retry 3 times)
- Task B: Clean data (retry 2 times)
- Task C: Train model (no retry)

This way, different parts of the pipeline can have their own retry logic, depending on how critical or error-prone they are.

Step 6: Use alerting and error-handling for the final failure

If all retries fail, I make sure the pipeline sends a notification or writes the failure log somewhere (like in a Delta table or error log file), so I or my team can investigate later.

In summary, to manage job retries in Databricks, I:

- Set retries and retry intervals in the job/task settings
- Use them mainly for temporary or flaky errors
- Optionally write custom retry logic in notebooks for more control
- Monitor job history and failures using the UI or APIs
- Use email/webhook alerts to track retry failures
- · Customize retry behavior for each task in a multi-step workflow

Scenario 14. How do you manage secrets securely in Azure Databricks when connecting to external sources?

When connecting to external sources like databases, APIs, or storage accounts, I make sure that all sensitive credentials (like usernames, passwords, access keys, or tokens) are stored securely. In Databricks, the best way to manage secrets is by using Databricks Secret Scopes. I never hardcode secrets directly in notebooks or code.

Here's how I manage secrets securely in Databricks step by step:

Step 1: Create a secret scope

A secret scope is like a container to hold all my secrets.

If I'm using Azure Databricks with Azure Key Vault (which is very common in enterprise settings), I usually create a Databricks-backed scope or link it to Azure Key Vault.

Option 1: Databricks-backed scope

From Databricks CLI:

databricks secrets create-scope --scope my-secret-scope

Option 2: Azure Key Vault-backed scope

If I want to connect Databricks to Azure Key Vault directly:

- 1. First, I create an Azure Key Vault in the Azure portal.
- 2. I add secrets in the Key Vault.
- 3. Then I use the Azure CLI or ARM template to link the Key Vault to Databricks using:

az databricks secret-scope create \

- --name my-kv-scope \
- --scope-backend-type AZURE_KEYVAULT \
- --resource-id <Azure_KeyVault_Resource_ID> \
- --dns-name <Azure_KeyVault_DNS_Name>

This allows Databricks to fetch secrets securely from Azure Key Vault at runtime.

Step 2: Add secrets to the scope

If I'm using a Databricks-managed scope, I can add secrets like this:

databricks secrets put --scope my-secret-scope --key db-password

It will prompt me to paste the password securely.

Step 3: Access secrets securely inside notebooks

Once the secret is stored, I do not write it directly in code. Instead, I access it using the dbutils.secrets.get() function.

Example:

```
db_password = dbutils.secrets.get(scope="my-secret-scope", key="db-password")
Now I can use this variable to connect to a database:
jdbc_url = "jdbc:sqlserver://myserver.database.windows.net:1433;database=mydb"
properties = {
    "user": "myuser",
    "password": db_password,
    "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
}
df = spark.read.jdbc(url=jdbc_url, table="customers", properties=properties)
```

Step 4: Secure access to secrets with ACLs

I make sure that only authorized users or service principals can access the secret scope.

In Databricks, I use Access Control Lists (ACLs) to control who can:

This way, the password is never shown in logs or stored in code.

- Read secrets
- Write secrets
- Manage scopes

Example (Databricks CLI):

databricks secrets put-acl --scope my-secret-scope --principal alice@example.com --permission READ

Step 5: Use environment variables for automation or CI/CD

When working with CI/CD pipelines (like Azure DevOps), I use environment variables in the pipeline, and read them into my Databricks job using secret scopes or cluster environment variables.

This ensures no secrets are stored in Git or hardcoded anywhere.

Step 6: Never print secrets in logs

Sometimes people accidentally do this:

print(dbutils.secrets.get(scope="my-scope", key="api-token"))

I always avoid this, because it will show the secret in notebook output or logs. Instead, I only use secrets in secure variables and use them as inputs to other functions.

Step 7: Rotate secrets regularly

For long-term projects, I also set reminders or policies to rotate secrets (especially database passwords or access tokens) every few months. If I'm using Azure Key Vault, this rotation can even be automated.

In summary, to manage secrets securely in Databricks when connecting to external sources, I:

- Create a secret scope (Databricks-managed or Azure Key Vault-backed)
- Store secrets securely using CLI or Azure portal
- Access secrets using dbutils.secrets.get() in notebooks
- Apply fine-grained access control (ACLs) to restrict who can access them
- Avoid printing secrets in logs
- Use secrets in CI/CD pipelines via environment variables
- Rotate secrets regularly for security

By following this approach, I make sure sensitive credentials are protected and not exposed anywhere in code or logs.

Scenario 15. How do you mask sensitive data (e.g., PII) before sharing it across teams in Databricks?

When dealing with sensitive data like names, phone numbers, email addresses, or social security numbers, it's very important to protect that data before sharing it with other teams for analysis or reporting. In Databricks, I use different techniques to mask or anonymize PII depending on the use case. Below is the step-by-step approach I follow:

Step 1: Identify sensitive data fields

First, I go through the dataset and identify columns that contain personally identifiable information (PII). These can be:

- Name
- Email
- Phone number
- Address
- National ID or SSN
- Bank account number

I usually tag or document these columns in a data dictionary or within notebook comments.

Step 2: Use column masking techniques in Spark

After identifying the sensitive fields, I mask them using simple transformations. Some common methods I use:

1. Partial masking (show only part of the value)

from pyspark.sql.functions import concat, lit, substring

```
\label{eq:df_def} $$ df-with Column("masked_email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("masked_email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("masked_email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("masked_email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("masked_email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("masked_email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("***@***.com"))) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("***@***.com")) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com")) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com")) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com")) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com")) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com")) $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com") $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com") $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("*****@****.com") $$ df = df-with Column("email", concat(substring("email", 1, 3), lit("****@****.com") $$ df = df-with Column("email", 1, 3), lit("****@*
```

df = df.withColumn("masked_phone", concat(lit("XXX-XXX-"), substring("phone", -4, 4)))

2. Hashing sensitive data

Hashing converts the data into an irreversible format, useful for anonymization.

from pyspark.sql.functions import sha2, col

```
df = df.withColumn("hashed_ssn", sha2(col("ssn"), 256))
```

Now even if someone sees the data, they can't reverse it back to the original SSN.

3. Nullify or remove columns if not needed

Sometimes, if a team doesn't need the PII at all, I just remove or set it to null.

```
df = df.drop("ssn", "name")
```

or

from pyspark.sql.functions import lit

```
df = df.withColumn("ssn", lit(None))
```

Step 3: Tokenization (if re-identification is needed)

If we need to link the masked data back to the original source later (like in customer support), I use tokenization. I store a mapping table like this:

customer_id original_ssn tokenized_ssn

```
1001 123-45-6789 abcde12345
```

This mapping table is stored in a secure, restricted access table, while the rest of the team only sees the token.

Step 4: Apply data masking using Delta Live Tables or views

Sometimes I create views that show masked data instead of raw tables.

CREATE OR REPLACE VIEW masked_customer_data AS

SELECT

customer_id,

CONCAT(SUBSTRING(name, 1, 1), '***') AS masked_name,

CONCAT('***-***-', SUBSTRING(phone, -4, 4)) AS masked_phone

FROM raw_customer_data;

Now other teams can query this view safely without accessing the raw PII.

Step 5: Use Unity Catalog for fine-grained access control

In Unity Catalog, I can create dynamic views or apply row- or column-level security, so that sensitive columns are hidden or masked based on user roles.

For example, only users in the pii_readers group can see the real name, others see NULL.

CREATE OR REPLACE VIEW secure_view AS

SELECT

customer_id,

CASE

WHEN is_member('pii_readers') THEN name

ELSE NULL

END AS name.

email

FROM raw_table;

This ensures masking is done automatically based on who is accessing the data.

Step 6: Document and validate masking logic

Before sharing data with any team, I:

- Document which columns were masked and how
- Review the logic with data governance or security teams (if required)
- Validate that no PII remains in the shared dataset

Step 7: Automate masking as part of the pipeline

If I'm building an ETL pipeline in Databricks, I include the masking logic as part of the final transformation step before saving the output.

masked_df.write.format("delta").save("/mnt/shared/masked_customer_data")

Now the shared table always contains masked values, and the raw data remains protected.

In summary, to mask sensitive data in Databricks before sharing:

- I identify all PII fields first
- Then I mask them using partial masking, hashing, nulling, or tokenization
- I can also use views or Delta Live Tables to serve masked versions of data
- For advanced control, I apply dynamic views using Unity Catalog
- I make sure masking is tested, automated, and documented properly

This keeps sensitive data safe while still allowing teams to do useful analysis.

Scenario 16. How do you pass parameters to a Databricks notebook from a pipeline or API?

In real-world projects, I often need to pass parameters into a Databricks notebook. For example, I might want to pass a date, table name, or file path when running a job. This is useful when the same notebook is reused in multiple pipelines or environments.

Here's how I handle passing parameters to a Databricks notebook step by step:

Step 1: Use widgets in the target notebook to receive parameters

In Databricks, I use dbutils.widgets to define and read parameters. These widgets act like input fields that receive values when the notebook is run from a pipeline, a job, or an API.

At the top of the notebook, I define the widgets like this:

```
dbutils.widgets.text("input_date", "")
```

input_date = dbutils.widgets.get("input_date")

dbutils.widgets.text("table_name", "")

table_name = dbutils.widgets.get("table_name")

This makes the notebook ready to accept those values.

Step 2: Pass parameters from a Databricks Job

If I'm running the notebook from a scheduled or triggered job, I can pass the parameters through the Job configuration.

In the Databricks UI:

- 1. Go to Jobs → Create or edit a job
- 2. In the Tasks section, I choose the notebook
- 3. Scroll to Parameters, and set key-value pairs like:
 - input_date = 2025-07-01
 - o table_name = sales_data

Now when the job runs, it will automatically pass these values to the notebook, and they will be captured using dbutils.widgets.get().

Step 3: Pass parameters from Azure Data Factory

If I'm using Azure Data Factory (ADF) to run a Databricks notebook as part of a pipeline, I can pass parameters from the ADF pipeline like this:

- 1. In the ADF pipeline, add a Databricks Notebook activity
- 2. Under Base Parameters, define the key-value pairs:

```
{
  "input_date": "2025-07-01",
  "table_name": "sales_data"
}
```

These values are sent to the notebook during execution, and inside the notebook I use dbutils.widgets.get() to access them.

Step 4: Pass parameters using the Databricks REST API

If I'm triggering notebook execution programmatically using the **REST API**, I use the 2.1/jobs/run-now endpoint and pass the parameters as part of the request.

```
Example request body (in JSON):
{
 "job_id": 123,
 "notebook_params": {
  "input_date": "2025-07-01",
 "table_name": "sales_data"
}
Or if I'm using the Jobs API to create and run a one-time job:
 "run_name": "Run Notebook with Parameters",
 "notebook_task": {
  "notebook_path": "/Shared/MyNotebook",
  "base_parameters": {
  "input_date": "2025-07-01",
  "table_name": "sales_data"
 }
},
"existing_cluster_id": "abc123"
}
```

This allows full control when triggering notebooks from custom apps or tools.

Step 5: Use default values for testing

In my notebooks, I often add default values to make local testing easier:

input_date = dbutils.widgets.get("input_date") or "2025-01-01"

table_name = dbutils.widgets.get("table_name") or "default_table"

This way, if I run the notebook manually without a pipeline, it still works with default values.

Step 6: Optional - Convert widgets to variables in SQL cells

If I want to use the parameters in SQL cells inside the notebook, I use this syntax:

-- Use the widget value as a SQL variable

%sql

SELECT * FROM \${table_name} WHERE transaction_date = '\${input_date}'

Databricks automatically replaces \${table_name} and \${input_date} with the actual values passed through widgets.

In summary, to pass parameters to a Databricks notebook:

- I define widgets at the top using dbutils.widgets.text()
- I pass values via:

Job configuration

Azure Data Factory

REST API

- I safely retrieve values using dbutils.widgets.get()
- I optionally add defaults for manual testing
- I can use the values inside Python or SQL cells
- This makes my notebooks reusable, dynamic, and production-ready

This parameterization is a best practice when working with automated pipelines and helps make the notebook flexible for multiple use cases.

Scenario 17. How do you scale up and down clusters automatically based on workload in Databricks?

In Databricks, automatic cluster scaling is essential for optimizing resource usage and cost management. To ensure that workloads are processed efficiently without over-provisioning resources, Databricks provides several mechanisms for scaling clusters automatically based on workload. Here's how I manage this process:

Step 1: Use Autoscaling with Clusters

Databricks supports autoscaling clusters, where the number of worker nodes is adjusted dynamically based on the workload. The cluster scales up when there's high demand (e.g., more data or more tasks to process) and scales down when the workload reduces.

How to configure autoscaling:

- 1. Create a new cluster in the Databricks UI.
- 2. Under Cluster Mode, choose Standard or High Concurrency (depending on the workload requirements).
- 3. In the Autoscaling section, enable the autoscaling option:

Min Workers: The minimum number of worker nodes that Databricks should keep running.

Max Workers: The maximum number of worker nodes that the cluster can scale to based on workload.

For example, if I expect low traffic most of the time but want the ability to scale when necessary, I might set:

- Min Workers: 2 (to ensure the cluster is always ready to handle small jobs)
- Max Workers: 10 (to handle high workloads during peak times)

Step 2: Adjust Cluster Type Based on Jobs

For different types of workloads (batch vs. streaming), I can scale clusters differently:

- Batch Jobs: For non-continuous workloads (like ETL jobs), I can use standard clusters with autoscaling, ensuring that the cluster scales based on the number of tasks or data volume.
- Streaming Jobs: For streaming workloads (e.g., processing real-time data), I can configure a high-concurrency cluster that can also scale based on the incoming data stream. This ensures optimal performance without wasting resources.

Step 3: Set up Databricks Jobs with Dynamic Cluster Allocation

When scheduling jobs in Databricks (for example, in a Databricks Job), I can set the cluster configuration to use job clusters that automatically scale.

- Job clusters are clusters that are spun up and scaled according to the job's needs. This ensures that resources are allocated only when required.
- Databricks will manage the cluster lifecycle automatically, creating it when the job starts and terminating it when the job finishes.

Example Job Cluster Setup:

- 1. Go to Jobs and create a new job.
- 2. Under Cluster, choose New cluster or Existing cluster.
- 3. Select Autoscaling and set the min and max worker nodes based on the job's anticipated load.

Databricks will automatically determine the resources needed, scaling the cluster dynamically while the job is running.

Step 4: Use Cluster Policy to Control Scaling Behavior

To ensure that the cluster does not scale beyond a certain point (e.g., to avoid cost overruns), I can define cluster policies in Databricks.

- 1. I create a cluster policy that restricts the number of workers or limits other settings like the instance types that can be used.
- 2. This ensures that autoscaling stays within the budget and doesn't over-provision resources.

For example, I can limit the maximum number of nodes that a cluster can scale to, ensuring I don't exceed resource limits for cost management.

Step 5: Monitor and Adjust Based on Usage

Once autoscaling is configured, it's important to continuously monitor the cluster's performance and adjust the scaling parameters as needed. Databricks provides built-in monitoring tools for tracking cluster resource usage:

- 1. **Cluster Metrics**: I can monitor CPU, memory, and disk utilization in real-time. If I notice a pattern where the cluster is consistently running at maximum capacity, I may increase the max worker nodes limit.
- 2. **Job Metrics**: I can also analyze the job performance and understand if the job duration or resource usage changes significantly over time.

Based on this data, I can tweak autoscaling parameters for optimization.

Step 6: Scaling Using Databricks REST API

If I want to scale clusters programmatically, I can use the Databricks REST API to resize the cluster or adjust worker nodes dynamically. For instance, the 2.0/clusters/edit API can be used to modify cluster settings based on workload triggers.

Here's how I might use the API to scale a cluster:

```
import requests
url = 'https://<databricks-instance>/api/2.0/clusters/edit'
data = {
    "cluster_id": "<cluster-id>",
    "num_workers": 10 # Set to desired number of workers based on workload
}
response = requests.post(url, json=data, headers={"Authorization": "Bearer <token>"})
```

In Summary:

To scale clusters up and down automatically based on workload in Databricks:

- 1. Enable autoscaling in the cluster configuration, setting min and max worker nodes.
- 2. Use job clusters that spin up and down as jobs are executed, automatically adjusting to workload demands.
- 3. Set up cluster policies to control the scaling limits and prevent over-provisioning.
- 4. Monitor cluster and job metrics to fine-tune scaling based on observed performance.
- 5. Optionally, use the Databricks REST API to resize clusters programmatically for advanced scaling needs.

By using these approaches, I ensure that Databricks clusters are optimally scaled based on workload, saving costs during low-demand periods while ensuring high performance when the workload increases.

Scenario 18. How do you secure your Databricks workspace and prevent unauthorized access?

To secure the Databricks workspace and prevent unauthorized access, I follow a layered security approach. This includes setting up proper authentication, role-based access control, network security, data encryption, and audit logging. Below is how I would explain it step by step during an interview:

Step 1: Enable Single Sign-On (SSO) and Azure Active Directory (AAD) Integration

The first thing I do is integrate Databricks with Azure Active Directory (AAD). This ensures that only authorized users within our organization can access the workspace.

- We configure Single Sign-On (SSO) so that users log in using their company credentials.
- This also helps enforce Multi-Factor Authentication (MFA) through AAD policies.

So even if someone gets hold of a username/password, they can't log in without the second factor.

Step 2: Use Role-Based Access Control (RBAC) for Users and Groups

I avoid giving wide access to everyone. Instead, I:

- Assign users to AAD security groups based on their roles (like data engineer, data scientist, analyst).
- In Databricks, I map these groups to specific workspace permissions and cluster permissions.

For example:

- Data Engineers can create and manage clusters.
- Analysts can run notebooks but not manage clusters.
- Only Admins can create or delete secret scopes.

This minimizes the risk of unauthorized actions within the workspace.

Step 3: Secure Clusters and Jobs

For each cluster, I define who can:

- Start, stop, or edit the cluster
- Attach notebooks or jobs to it

In Databricks, I make use of Cluster Access Control to:

- Prevent users from launching expensive clusters
- Avoid users running sensitive notebooks on shared clusters

I also use Job Access Control to restrict who can view or run scheduled jobs.

Step 4: Secure Data Access with Unity Catalog or Table ACLs

To prevent unauthorized access to data:

- I enable Unity Catalog to manage table-level and column-level permissions across all workspaces.
- I define fine-grained access control at the database, table, or even column level.

Example:

- Only the finance team can view salary columns.
- Data scientists can access anonymized views, not raw data.

If Unity Catalog isn't enabled, I use Table Access Control (TAC) for managing Hive metastore access in legacy environments.

Step 5: Use Secret Scopes for Managing Credentials

To avoid hardcoding sensitive information like database passwords, I:

- Store secrets in Databricks secret scopes or use Azure Key Vault-backed scopes
- Access them securely inside notebooks using dbutils.secrets.get()

This way, users can run code without ever seeing the underlying credentials.

Step 6: Configure IP Access Lists (IP Whitelisting)

To restrict who can access the workspace, I:

- Set up IP access lists so that only requests from specific IP ranges (like our corporate VPN) are allowed.
- This helps block external or unauthorized users from logging in even if they have credentials.

This feature is available in the Premium and Enterprise tiers.

Step 7: Network Security and Private Link

For maximum security, I configure:

- VNet injection to deploy Databricks in a private Azure Virtual Network
- Azure Private Link to securely connect to resources like storage, databases, and Key Vault without going over the public internet

This ensures that data in transit is protected from man-in-the-middle attacks.

Step 8: Encryption for Data at Rest and In Transit

Databricks automatically encrypts:

- Data at rest using Azure Storage encryption (with Microsoft-managed or customermanaged keys)
- Data in transit using TLS/SSL for communication

For sensitive use cases, I enable Customer-Managed Keys (CMK) to have full control over encryption keys.

Step 9: Enable Audit Logging

To track any suspicious or unauthorized activity:

- I enable audit logs in the Databricks workspace and send them to Azure Log Analytics or Azure Monitor
- These logs include login attempts, notebook runs, permission changes, job execution, etc.

We then set up alerts on unusual activity like failed login attempts or permission escalations.

Step 10: Follow Least Privilege Principle

Throughout the setup:

- I always follow the principle of least privilege, giving users and service principals only the minimum required permissions.
- This reduces the attack surface and helps contain damage if something goes wrong.

In Summary:

To secure a Databricks workspace and prevent unauthorized access, I:

- 1. Integrate with Azure Active Directory for SSO and MFA
- 2. Use role-based access control (RBAC) to manage user roles and permissions
- 3. Control access to clusters, notebooks, and jobs
- 4. Manage data access using Unity Catalog or Table ACLs
- 5. Store credentials securely using secret scopes
- 6. Restrict access using IP access lists and network rules
- 7. Use Private Link and VNet for secure data connections
- 8. Ensure encryption of data at rest and in transit
- 9. Enable audit logging for monitoring and investigation
- 10. Follow the least privilege principle across all services

These steps help protect sensitive data, reduce risks, and comply with enterprise and regulatory security requirements.

Scenario 19. How do you use Databricks Repos for version control and collaboration?

In Databricks, I use Repos to manage notebooks just like code files in a Git repository. It allows me and my team to collaborate better, track changes, and follow proper version control using tools like GitHub, Azure Repos, or GitLab. Below is how I would explain it in an interview, step by step in simple language:

Step 1: Connect Databricks to a Git Provider

Before using Repos, I first connect Databricks to a Git system. In most of my projects, I use either GitHub or Azure DevOps Repos.

- I go to User Settings in Databricks
- Then under Git Integration, I choose the provider (like GitHub)
- I enter a personal access token to authenticate

This step allows me to pull and push code between Databricks and the Git system.

Step 2: Create or Clone a Repo in Databricks

After Git is connected, I can start using Databricks Repos:

- I click on Repos from the sidebar
- Then click Add Repo
- I enter the Git URL of the repository (for example, https://github.com/org/project.git)
- Databricks clones the Git repo and shows the notebooks, Python scripts, and other files inside my workspace

If the repo already exists in GitHub, this pulls all the code into Databricks so I can work with it directly.

Step 3: Work on Notebooks or Files Collaboratively

Once the repo is in Databricks:

- I open the notebooks or .py files
- I write or modify code directly inside Databricks
- These files are linked to the Git branch, so I can commit changes back to Git

This helps multiple team members work on the same repo, just like how developers work with regular Git workflows.

Step 4: Use Git Commands Within Databricks Repos

Databricks provides built-in Git options. I can:

- Pull the latest changes from the remote repo
- Commit and push my local changes to Git
- Switch branches (e.g., move from main to feature/mybranch)
- Create branches directly from the Databricks UI

For example:

- I work on a feature in a notebook
- I switch to a new branch called feature-data-cleanup
- After testing the code, I commit and push it
- Then I raise a pull request in GitHub for review

This helps maintain code quality and supports code reviews.

Step 5: Resolve Conflicts and Keep History

When two people edit the same notebook, Git may show a conflict.

- To avoid this, we usually pull the latest changes before editing
- If there's a conflict, we can resolve it either using Git commands or from the GitHub interface

Also, because everything is in Git, I can:

- See the history of each file
- Roll back to a previous version
- Track who made what change

This gives full transparency and safety during collaboration.

Step 6: Use Repos for CI/CD Integration

I also use Databricks Repos as part of the CI/CD pipeline. For example:

- A developer pushes code to GitHub
- A GitHub Action is triggered
- The CI pipeline runs unit tests or deploys the notebooks to a staging workspace

Since the code lives in Git and not just inside Databricks, it becomes easy to integrate with automation and testing tools.

Step 7: Use Folders for Better Organization

Within a Databricks Repo, I structure folders like this:

/src --> Python modules or reusable functions

/notebooks --> Feature notebooks for development

/tests --> Unit test files

/config --> YAML or JSON configs

This makes the repo more organized and developer-friendly.

Step 8: Use Collaboration Best Practices

To make the best use of Repos in a team setting, I follow a few habits:

- Create a new Git branch for every new task or bug fix
- Use meaningful commit messages
- Review each other's code via pull requests
- Add comments to notebooks for clarity
- Use Git tags or releases for tracking versions in production

In Summary:

To use Databricks Repos for version control and collaboration, I:

- 1. Connect Databricks to GitHub or Azure DevOps
- 2. Clone the Git repo into Databricks using Repos
- 3. Work on notebooks and scripts directly from the repo
- 4. Pull, commit, and push changes using Git commands in the UI
- 5. Use branches to work safely and avoid conflicts
- 6. Integrate with CI/CD tools for testing and deployment
- 7. Follow Git best practices to work smoothly with the team

Using Databricks Repos helps us treat our notebooks like proper code, collaborate better, and build more reliable and maintainable data pipelines.

Scenario 20. How do you use Databricks to build a recommendation engine?

In an interview, I would explain that Databricks is a great platform for building recommendation engines because it supports large-scale data processing with Spark, provides built-in machine learning libraries, and integrates well with notebooks for collaboration.

Here is how I would build a recommendation engine in Databricks, step by step, using very simple language:

Step 1: Understand the goal of the recommendation engine

Before starting, I make sure I understand what kind of recommendation system is needed. For example:

- Product recommendation (e.g., recommend similar products to users)
- Movie recommendation (e.g., like Netflix)
- Content recommendation (e.g., suggest articles or videos)

There are two common types:

- 1. Collaborative filtering recommends items based on user behavior (e.g., what other similar users liked)
- 2. Content-based filtering recommends items similar to what the user liked in the past

In Databricks, I mostly use collaborative filtering with Spark's built-in ALS (Alternating Least Squares) algorithm.

Step 2: Ingest and prepare the data

I load the data first. Usually, I have data like:

- User ID
- Item ID (e.g., product, movie)
- Rating or interaction (e.g., 5 stars, click, purchase)

Example:

from pyspark.sql.functions import col

```
df = spark.read.csv("/mnt/data/user_item_ratings.csv", header=True, inferSchema=True)
ratings_df = df.select(
    col("user_id").cast("int"),
    col("item_id").cast("int"),
    col("rating").cast("float")
)
```

I make sure the data has no missing or incorrect values and is in the right format for the model.

Step 3: Train the recommendation model using ALS

Databricks has Spark MLlib, which includes the ALS model, a popular choice for collaborative filtering.

from pyspark.ml.recommendation import ALS

```
als = ALS(
    userCol="user_id",
    itemCol="item_id",
    ratingCol="rating",
    coldStartStrategy="drop", # to handle missing predictions
    nonnegative=True,
    implicitPrefs=False,
    rank=10,
    maxIter=10,
    regParam=0.1
)
```

This model learns user and item preferences and can predict how much a user might like an item they haven't interacted with yet.

Step 4: Generate recommendations

```
After training the model, I can now generate top-N recommendations for each user: user_recs = model.recommendForAllUsers(5) # top 5 recommendations per user user_recs.show()

Or, I can recommend items for a specific user: user_id = 123

user_df = spark.createDataFrame([(user_id,)], ["user_id"])

model.recommendForUserSubset(user_df, 5).show()
```

Step 5: Evaluate the model (optional)

If I have a test dataset, I can evaluate the model accuracy using metrics like **RMSE** (Root Mean Square Error):

```
from pyspark.ml.evaluation import RegressionEvaluator
predictions = model.transform(test_df)
evaluator = RegressionEvaluator(
    metricName="rmse",
    labelCol="rating",
    predictionCol="prediction"
)
rmse = evaluator.evaluate(predictions)
print(f"Root-mean-square error = {rmse}")
Lower RMSE means better predictions.
```

Step 6: Store the recommendations

I can save the output to a Delta table, a database, or a storage location like Azure Data Lake: user_recs.write.format("delta").mode("overwrite").save("/mnt/output/recommendations")

This allows other teams or applications (like web or mobile apps) to use the recommendations.

Step 7: Automate the pipeline using jobs

To make it run automatically:

- I schedule this notebook as a job in Databricks.
- The job can run daily or hourly, depending on how often new data comes in.
- I also include alerts in case the job fails.

Step 8: Monitor and retrain regularly

I keep track of how well the recommendations are performing (through feedback or click data), and retrain the model regularly with fresh data to keep it accurate.

Bonus: Handle real-time recommendations (if needed)

For near real-time use cases:

- I can use Databricks Structured Streaming to process user activity as it happens.
- Combine streaming data with a pre-trained model for generating fresh recommendations on the fly.

In summary:

To build a recommendation engine in Databricks, I:

- 1. Ingest user-item interaction data (ratings, clicks, purchases)
- 2. Clean and prepare the data
- 3. Train a collaborative filtering model using ALS
- 4. Generate top-N recommendations for each user
- 5. Save the output for use by other systems
- 6. Schedule and monitor the process as a Databricks job
- 7. Retrain the model regularly to improve accuracy

Databricks helps by providing scalable Spark clusters, built-in ML tools, and easy job orchestration, which makes it ideal for building large-scale, reliable recommendation systems.

Scenario 21. How do you use Databricks to count how many times the word "the" appears in a book?

If I get this question in an interview, I would explain that this is a classic example of a word count problem, and Databricks with PySpark is perfect for this because it can process large text files efficiently.

Here's how I would solve it step by step in Databricks using very simple words:

Step 1: Upload the book into Databricks

First, I would upload the text file of the book (let's say it's a .txt file like book.txt) into DBFS (Databricks File System).

- I go to the Data tab → click Add Data → upload the file.
- It gets saved to something like: /dbfs/FileStore/books/book.txt

Step 2: Read the book as a text file using Spark

I use PySpark to read the text file:

book_path = "/FileStore/books/book.txt"

text_df = spark.read.text(book_path)

This will read each line of the book into a DataFrame, where each row contains a line of text.

Step 3: Convert all words to lowercase

To count the word "the" accurately, I make sure everything is in lowercase (so that "The", "THE", and "the" are all treated the same):

from pyspark.sql.functions import lower, col

text_df = text_df.select(lower(col("value")).alias("line"))

Step 4: Split lines into words

Now I split each line into words using split:

from pyspark.sql.functions import split, explode

words_df = text_df.select(explode(split(col("line"), "\\s+")).alias("word"))

This breaks each line into individual words and creates one row per word.

Step 5: Filter only the word "the"

Then I filter the words to keep only those that exactly match "the":

the_words_df = words_df.filter(col("word") == "the")

Step 6: Count how many times it appears

Finally, I count how many times the word "the" appears:

```
the_count = the_words_df.count()
```

print(f"The word 'the' appears {the_count} times in the book.")

Optional: Remove punctuation

Sometimes words like "the," or "the." can get counted separately because of punctuation. So, to make the count more accurate, I can remove punctuation using a regular expression:

from pyspark.sql.functions import regexp_replace

```
clean_df = text_df.select(regexp_replace(col("line"), "[^a-zA-Z\\s]", "").alias("line"))
words_df = clean_df.select(explode(split(col("line"), "\\s+")).alias("word"))
```

the_words_df = words_df.filter(col("word") == "the")

the_count = the_words_df.count()

This removes commas, periods, and other symbols before splitting into words.

In summary:

To count the number of times "the" appears in a book using Databricks, I:

- 1. Upload the book text file to DBFS
- 2. Read it using spark.read.text()
- 3. Convert all lines to lowercase
- 4. Split lines into words using split and explode
- 5. Remove punctuation if needed
- 6. Filter the word "the"
- 7. Use .count() to get the total number

This is a simple and efficient way to process even large books using Spark in Databricks.

Scenario 22. How do you use Databricks to improve the accuracy of machine learning models?

If I get this question in an interview, I would explain that Databricks provides a lot of tools to build and train machine learning models at scale. But more importantly, it gives me the environment and flexibility to keep improving model accuracy over time.

Here's how I would answer this step by step in simple language:

Step 1: Start with good quality and clean data

The first and most important thing I focus on is data quality. Databricks helps here with Spark for big data processing and Delta Lake for managing versions of data.

- I remove duplicates and missing values
- I handle outliers or incorrect entries
- I fill or impute missing values where needed

Example:

clean_df = raw_df.dropna().dropDuplicates()

If the data is bad, even the best algorithm won't give good accuracy.

Step 2: Feature engineering to make data more meaningful

I use Databricks notebooks to perform feature engineering, which means creating new columns (features) from raw data to help the model understand patterns better.

Some things I usually do:

- Convert dates into day, month, weekday
- Convert text into numerical features using TF-IDF or Word2Vec
- Normalize or scale numeric columns
- Use one-hot encoding for categorical variables

Example:

from pyspark.ml.feature import OneHotEncoder, StringIndexer

indexer = StringIndexer(inputCol="category", outputCol="categoryIndex")

encoder = OneHotEncoder(inputCols=["categoryIndex"], outputCols=["categoryVec"])

These steps often lead to a big improvement in accuracy.

Step 3: Try different algorithms and compare

Databricks makes it easy to try many ML algorithms. I usually start with a few and compare them:

- Logistic Regression
- Decision Trees / Random Forest
- Gradient Boosted Trees (like XGBoost or LightGBM)
- Neural Networks (if using deep learning frameworks)

Example:

```
from pyspark.ml.classification import RandomForestClassifier

rf = RandomForestClassifier(labelCol="label", featuresCol="features")

model = rf.fit(train_data)
```

Then I compare results using metrics like accuracy, F1 score, or AUC.

Step 4: Use hyperparameter tuning

To make the model more accurate, I tune its settings (like number of trees, max depth, learning rate, etc). Databricks supports hyperparameter tuning using CrossValidator or TrainValidationSplit.

Example:

from pyspark.ml.tuning import ParamGridBuilder, CrossValidator

```
paramGrid = ParamGridBuilder() \
    .addGrid(rf.maxDepth, [5, 10]) \
    .addGrid(rf.numTrees, [20, 50]) \
    .build()
```

```
crossval = CrossValidator(estimator=rf,
estimatorParamMaps=paramGrid,
evaluator=evaluator,
numFolds=3)
```

This helps in finding the best combination of settings for the model.

Step 5: Use MLflow for tracking and comparing experiments

Databricks has built-in support for MLflow, which I use to:

- Track all my experiments
- Log model metrics, parameters, and code
- Compare models side-by-side
- Restore the best model anytime

Example:

```
import mlflow.spark
```

```
with mlflow.start_run():
   model = rf.fit(train_data)
   mlflow.spark.log_model(model, "rf-model")
   mlflow.log_metric("accuracy", accuracy)
```

This helps me keep track of which changes improved the accuracy.

Step 6: Try ensemble models or stacking

When I want even better results, I combine multiple models using ensemble methods:

- Bagging (e.g., Random Forest)
- Boosting (e.g., XGBoost)
- Stacking (combine outputs from different models)

In Databricks, I can run different models in parallel using multiple notebooks or clusters.

Step 7: Balance the dataset (if it's imbalanced)

For classification problems, if one class appears a lot more than the other (like 90% vs 10%), the model accuracy can be misleading.

So I use techniques like:

- Undersampling or oversampling
- SMOTE (Synthetic Minority Over-sampling Technique)

Databricks lets me use libraries like imbalanced-learn to handle this.

Step 8: Monitor and retrain regularly

Finally, I make sure the model is not just accurate during training, but also stays accurate over time.

- I set up a Databricks job to retrain the model weekly or monthly
- I monitor its accuracy using test data or user feedback
- If the accuracy drops, I go back, retrain, and improve it again

In summary:

To improve ML model accuracy in Databricks, I follow these steps:

- 1. Clean and prepare high-quality data
- 2. Do smart feature engineering
- 3. Try and compare different algorithms
- 4. Use hyperparameter tuning with CrossValidator
- 5. Track everything using MLflow
- 6. Try ensemble models or stacking
- 7. Handle imbalanced datasets
- 8. Retrain the model regularly and monitor its performance

Databricks supports all of this with scalable compute, built-in ML tools, and collaboration features, which makes it very helpful for improving model accuracy.

Scenario 23. How do you use Databricks to perform A/B testing?

If I'm asked this question in an interview, I would explain that A/B testing is used to compare two versions of something (like a feature, model, or webpage) to see which one performs better. Databricks is very helpful in performing A/B testing because it supports large-scale data processing, data visualization, and statistical analysis all in one place.

Here's how I would perform A/B testing in Databricks, step by step, using simple and clear language:

Step 1: Understand the goal of the A/B test

First, I clearly define what I'm testing. For example:

- Version A: Existing product recommendation model
- Version B: New product recommendation model

The goal could be: which version results in more purchases or user engagement.

Step 2: Collect A/B test data

Usually, the application or website already assigns users randomly into Group A or Group B.

So, I expect a dataset like this:

user_id group conversion clicks time_spent

101	Α	1	4	120
102	В	0	2	90

This data might be stored in a Delta table, a database, or in raw files. I load it in Databricks like this:

df = spark.read.format("delta").load("/mnt/data/ab_test_results")

Step 3: Validate the groups are randomized and balanced

Before analyzing results, I check if the A and B groups are similar in size and characteristics (like age, location, device type) to make sure the test is fair.

df.groupBy("group").count().show()

If the distribution is uneven, it could bias the results.

Step 4: Calculate key metrics for each group

I calculate the average of key metrics (e.g., conversion rate, click-through rate) for each group.

Example:

```
df.groupBy("group").agg(
    avg("conversion").alias("avg_conversion"),
    avg("clicks").alias("avg_clicks"),
    avg("time_spent").alias("avg_time")
).show()
```

This gives me a basic comparison between A and B.

Step 5: Perform statistical testing

To make sure the difference in results is not just random, I perform a statistical test like a t-test or chi-square test.

Since Spark doesn't support these directly, I often convert data to Pandas for this step:

```
pandas_df = df.toPandas()
from scipy.stats import ttest_ind
group_a = pandas_df[pandas_df['group'] == 'A']['conversion']
group_b = pandas_df[pandas_df['group'] == 'B']['conversion']
t_stat, p_value = ttest_ind(group_a, group_b)
print(f"T-statistic = {t_stat}, p-value = {p_value}")
```

If the p-value is less than 0.05, the difference is considered statistically significant.

Step 6: Visualize the results

Databricks notebooks support graphs and charts, which I use to visualize:

- Conversion rates per group
- Distribution of metrics
- Confidence intervals

```
Example with Matplotlib:
```

```
import matplotlib.pyplot as plt
plt.bar(['Group A', 'Group B'], [group_a.mean(), group_b.mean()])
plt.title("Conversion Rate Comparison")
plt.ylabel("Conversion Rate")
plt.show()
```

This helps make the results easier to understand and explain.

Step 7: Interpret the results

I interpret results based on:

- Which group performed better?
- Is the result statistically significant?
- What action should be taken (e.g., roll out version B)?

If Group B had a higher conversion rate and the p-value is low, I can say it's safe to move forward with version B.

Step 8: Share findings with stakeholders

Finally, I create a Databricks dashboard or export the notebook results to share with product teams, management, or clients. Databricks lets me display metrics, charts, and even export the notebook as PDF.

Step 9: Automate future A/B tests (optional)

If we do A/B tests regularly, I set up:

- A job to run analysis daily/weekly
- Save metrics and results into Delta tables
- Alert if a version significantly outperforms the other

In summary:

To perform A/B testing in Databricks, I:

- 1. Load and clean A/B test data
- 2. Make sure the groups are balanced
- 3. Calculate metrics like conversion rate
- 4. Use statistical tests (like t-test) to compare A vs B
- 5. Visualize results using charts
- 6. Make decisions based on p-values and metric differences
- 7. Share results through dashboards or reports
- 8. Automate future tests if needed

Databricks gives me the power to handle big datasets, run analysis at scale, use Python and statistics, and collaborate with teams in one place, which makes it a great platform for A/B testing.

Scenario 24. How do you use Databricks to find the tallest building in NYC from a dataset?

If I'm asked this in an interview, I would explain that Databricks can handle structured data using Spark DataFrames, and I can easily use filtering, sorting, and aggregation functions to find the tallest building.

Here's how I would approach it step by step, in very simple and clear terms:

Step 1: Load the dataset into Databricks

Let's assume I've been given a dataset of buildings. It may be in a CSV, JSON, or Delta file.

Here's how I load it into a DataFrame:

df = spark.read.option("header", True).csv("/FileStore/data/nyc_buildings.csv")

Or if it's a Delta table:

df = spark.read.format("delta").load("/mnt/data/nyc_buildings")

Step 2: Explore the dataset to understand the structure

I use display(df) or df.printSchema() to check what columns are available.

Let's say I see columns like:

- building_name
- city
- height_feet

Step 3: Filter only buildings from NYC

To make sure I'm only working with buildings in New York City, I filter the data:

```
nyc_df = df.filter(df.city == "New York")
```

If the city column contains variations like "NYC", "New York City", or "Manhattan", I might use like or isin():

nyc_df = df.filter(df.city.isin("New York", "New York City", "NYC", "Manhattan"))

Step 4: Convert height to numeric (if needed)

Sometimes height values are stored as strings. So I convert the height_feet column to a numeric type like float:

from pyspark.sql.functions import col

nyc_df = nyc_df.withColumn("height_feet", col("height_feet").cast("float"))

Step 5: Find the tallest building

```
Now I sort the buildings by height and take the first one (the tallest):
```

tallest_building = nyc_df.orderBy(col("height_feet").desc()).limit(1)

Or use max() if I just want the tallest height:

from pyspark.sql.functions import max

nyc_df.select(max("height_feet")).show()

To also get the building name and height together:

display(tallest)

Step 6: Display or export the result

Finally, I can display the result in a notebook using display() or save it to a table or CSV if needed:

tallest.write.mode("overwrite").csv("/FileStore/results/tallest_building_nyc.csv")

In summary:

To find the tallest building in NYC using Databricks, I:

- 1. Load the dataset into a DataFrame
- 2. Filter the data for buildings in New York City
- 3. Make sure height is in numeric format
- 4. Sort by height in descending order
- 5. Pick the top row (the tallest building)
- 6. Display or export the result

This is a simple example of how Databricks makes it easy to explore structured data and answer real-world questions using Spark and Python.

Scenario 25. How do you troubleshoot "out of memory" errors in Databricks notebooks?

If I get an "out of memory" error in Databricks during an interview scenario, I would explain that it usually means Spark is trying to load or process more data than the driver or executor memory can handle. I follow a step-by-step process to troubleshoot and fix it, using simple strategies to reduce memory usage or increase resources.

Here's how I handle this, step by step, in plain and simple words:

Step 1: Understand where the error is happening

First, I check whether the memory issue is happening on the driver or executor.

- If the notebook crashes or restarts, it's likely a driver memory issue.
- If the job fails with an error like "GC overhead limit exceeded" or "Java heap space", it's more likely an executor memory issue.

I check the Spark UI from the "Driver Logs" or "Spark Jobs" tab in the notebook to see the full error message and which task failed.

Step 2: Reduce the amount of data being processed

If I'm reading a large file or table, I try to read only the needed columns or rows instead of loading everything.

Example:

Instead of this:

df = spark.read.parquet("/mnt/data/large_table")

Do this:

df = spark.read.parquet("/mnt/data/large_table").select("id", "timestamp", "amount")

If I'm doing a join or aggregation, I also try to **filter early** to reduce data before heavy processing:

df = df.filter(df["timestamp"] >= "2024-01-01")

Step 3: Use .persist() or .cache() wisely

If I'm caching large datasets, I check whether I really need it. Caching too much data fills up memory quickly.

I either:

- Avoid using .cache() if it's not reused often
- Use .persist(StorageLevel.DISK_ONLY) if memory is tight

Example:

from pyspark import StorageLevel

df.persist(StorageLevel.DISK_ONLY)

Step 4: Repartition the data

Sometimes, Spark uses very large partitions (for example, one partition with 1 GB of data), which overloads a single executor. I fix this by using .repartition() or .coalesce() to control partition sizes.

Example:

Increase partitions for parallelism

df = df.repartition(200)

Or reduce partitions before writing

df = df.coalesce(10)

I usually check the number of partitions with:

df.rdd.getNumPartitions()

Step 5: Avoid collecting large data to the driver

One common cause of out-of-memory is calling .collect() or display() on a huge dataset.

Bad:

This brings all data into the driver's memory (can crash the notebook)

df.collect()

Better:

Use .limit() first

df.limit(100).toPandas()

Step 6: Increase cluster resources (as last option)

If none of the optimizations help, I scale the cluster by:

- Increasing the driver and worker memory (using a larger node type)
- Increasing the number of executors or using autoscaling

In the cluster settings:

- I can choose r5d.4xlarge instead of Standard_DS3_v2 for more memory
- Enable autoscaling to let Databricks add more nodes when needed

Step 7: Use Delta Lake to optimize file size and reads

If the data is stored in Delta format, I use **optimize** and **ZORDER** to reduce the amount of data Spark needs to read.

Example:

OPTIMIZE my_table ZORDER BY (user_id)

This helps read only the relevant data and avoid loading huge amounts into memory.

Step 8: Monitor memory usage in the Spark UI

I go to the Spark UI > Executors tab to see how much memory each executor is using. This helps me know if:

- Memory is being spilled to disk (slow but avoids crashing)
- Tasks are failing repeatedly on certain nodes

In summary:

To troubleshoot out-of-memory errors in Databricks notebooks, I follow these steps:

- 1. Identify if it's a driver or executor memory issue
- 2. Read only the needed data and filter early
- 3. Avoid unnecessary caching or use disk storage
- 4. Repartition data to reduce skew or big partitions
- 5. Avoid .collect() or .toPandas() on large datasets
- 6. Scale the cluster size or use autoscaling
- 7. Optimize Delta tables with OPTIMIZE and ZORDER
- 8. Use Spark UI to monitor memory usage

By combining these techniques, I can usually fix out-of-memory issues in a clean and efficient way without wasting resources.

Scenario 26. How do you use Databricks to find common records (joins) between two datasets?

If I'm asked this in an interview, I would say that finding common records between two datasets is done using joins, and Databricks supports various types of joins like inner join, left join, right join, and full join using PySpark or SQL. The most common one for finding matching records in both datasets is the inner join.

Here's how I would explain it step by step using very simple words:

Step 1: Load both datasets into DataFrames

Let's say I have two datasets:

- Dataset A: customers who signed up
- Dataset B: customers who made a purchase

I load both files (CSV or Delta or any format) into Spark DataFrames.

Load Dataset A (signups)

signups = spark.read.option("header", True).csv("/mnt/data/signups.csv")

Load Dataset B (purchases)

purchases = spark.read.option("header", True).csv("/mnt/data/purchases.csv")

Step 2: Check the common key between datasets

Usually, I need a key column to match records between two datasets. It could be:

- customer_id
- email
- order_id

I inspect both schemas using:

signups.printSchema()

purchases.printSchema()

Let's say both datasets have a column called customer_id.

Step 3: Perform an inner join to find common records

I use an inner join to get records that exist in **both** datasets (i.e., customers who signed up *and* made a purchase).

common_df = signups.join(purchases, on="customer_id", how="inner")

If the column names are different, I match them explicitly:

common_df = signups.join(purchases, signups["customer_id"] == purchases["cust_id"],

"inner")

Step 4: View the common records

To preview the results:

display(common_df)

Or to show the first 5 records:

common_df.show(5)

Step 5: Select only the needed columns (optional)

If I want to show only some fields from both datasets:

result = common_df.select("customer_id", "signups.signup_date", "purchases.amount")

Step 6: Save the common records (optional)

If I want to save the result for further analysis:

result.write.mode("overwrite").parquet("/mnt/output/common_customers")

Step 7: Use SQL for the same logic (optional)

If I prefer SQL, I register both DataFrames as temporary views:

signups.createOrReplaceTempView("signups")

purchases.createOrReplaceTempView("purchases")

Then I run an SQL query:

SELECT *

FROM signups s

JOIN purchases p

ON s.customer_id = p.customer_id

In summary:

To find common records between two datasets in Databricks:

- 1. Load both datasets as Spark DataFrames
- 2. Identify the key column used for joining
- 3. Use join() with "inner" type to find matches
- 4. Filter/select only the needed columns
- 5. Use display() or .show() to view results
- 6. Save or export the final result if needed
- 7. Optionally use SQL for better readability

This approach works with large datasets efficiently since Spark runs the join operation in parallel across the cluster.

Scenario 27. How do you design a pipeline in Databricks to process data from Azure Data Lake and load it into Azure SQL Database?

If I'm asked this in an interview, I would say that I usually build this kind of pipeline in Databricks using PySpark to read data from Azure Data Lake, clean or transform it, and then write it to Azure SQL Database. I will explain the full pipeline step-by-step in very simple words.

Step 1: Set up the environment and configurations

Before I begin, I make sure:

- Azure Data Lake is connected using a service principal or managed identity
- Azure SQL Database connection details are available (server, database, username, password)
- I have created secrets to store credentials securely in Databricks

Step 2: Mount Azure Data Lake to Databricks (if needed)

```
Sometimes, I mount the ADLS Gen2 path to make file access easier.

configs = {

"fs.azure.account.auth.type": "OAuth",

"fs.azure.account.oauth.provider.type":

"org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider",

"fs.azure.account.oauth2.client.id": dbutils.secrets.get(scope="my-scope", key="client-id"),

"fs.azure.account.oauth2.client.secret": dbutils.secrets.get(scope="my-scope", key="client-secret"),

"fs.azure.account.oauth2.client.endpoint": "https://login.microsoftonline.com/<tenant-id>/oauth2/token"
}
```

```
dbutils.fs.mount(
    source = "abfss://<container-name>@<storage-account-name>.dfs.core.windows.net/",
    mount_point = "/mnt/data",
    extra_configs = configs)
```

But if mounting is not allowed, I directly access the file path using spark.read.

Step 3: Read data from Azure Data Lake

```
I read the raw data from ADLS using Spark. It can be in CSV, Parquet, Delta, or JSON format.

df = spark.read.option("header", True).csv("/mnt/data/sales_data.csv")

If it's a Delta table:

df = spark.read.format("delta").load("/mnt/data/delta/sales_data")
```

Step 4: Clean or transform the data

Next, I apply transformations like filtering, joining, renaming columns, converting formats, etc.

Example:

from pyspark.sql.functions import col, to_date

Step 5: Prepare Azure SQL connection details

```
I store credentials in a secret scope and pass them securely to Spark.
```

```
jdbc_url = "jdbc:sqlserver://<server-name>.database.windows.net:1433;database=<db-
name>"
jdbc_username = dbutils.secrets.get(scope="my-scope", key="sql-user")
jdbc_password = dbutils.secrets.get(scope="my-scope", key="sql-password")
connection_properties = {
    "user": jdbc_username,
    "password": jdbc_password,
    "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
```

Step 6: Write transformed data to Azure SQL Database

I write the final cleaned data into a table in Azure SQL. I can use append or overwrite mode depending on the use case.

```
df_cleaned.write.jdbc(
    url=jdbc_url,
    table="dbo.processed_sales",
    mode="overwrite",
    properties=connection_properties
```

Step 7: Automate the pipeline using Databricks Jobs

Once the notebook is working:

- 1. I schedule it as a Databricks Job
- 2. I configure a cluster for execution
- 3. I set up a schedule (like every night at 1 AM)
- 4. I set alerts for success/failure

Step 8: Add logging and error handling (optional but useful)

I add logs to track progress:

print("Reading completed")

print("Row count: ", df_cleaned.count())

And I may use try-except blocks to catch errors and send email or Slack alerts.

In summary:

To design a pipeline in Databricks that reads data from Azure Data Lake and writes to Azure SQL Database, I:

- 1. Set up access to Azure Data Lake and secrets
- 2. Read the source data using Spark
- 3. Apply necessary data cleaning or transformation
- 4. Securely connect to Azure SQL using JDBC
- 5. Write the transformed data into the SQL table
- 6. Automate the notebook using a scheduled Databricks Job
- 7. Add logging and error handling for reliability

This is a common real-world scenario, and I try to make the pipeline modular, secure, and easy to maintain.

Scenario 28. How do you implement a Medallion Architecture in Databricks?

If I'm asked this in an interview, I would say that Medallion Architecture is a simple and effective way to organize data in different layers to ensure data quality, structure, and reusability. It breaks down data processing into three stages:

• Bronze: raw data

Silver: cleaned and filtered data

Gold: business-ready, aggregated data

Databricks is a perfect platform to implement this because of its support for Delta Lake, which brings ACID transactions, versioning, and scalable processing.

Here's how I would explain it step-by-step:

Step 1: Understand the role of each layer

• Bronze Layer (Raw Data):

This is where I store the raw data exactly how it comes in no filters, no changes. This helps with traceability and debugging.

• Silver Layer (Cleaned Data):

This is where I clean, filter, and transform the data for example, fixing nulls, converting data types, joining reference tables.

Gold Layer (Business Data):

This is the final, polished data ready for dashboards, machine learning, and reports. Usually includes aggregations or calculated fields.

Step 2: Create the Delta Lake folders or tables for each layer

I make sure to define three separate storage locations or Delta tables:

bronze_path = "/mnt/datalake/bronze"

silver_path = "/mnt/datalake/silver"

gold_path = "/mnt/datalake/gold"

I store each table in its own folder so it's easy to manage.

Step 3: Ingest data into the Bronze layer

Let's say I'm reading a JSON file from Azure Data Lake:

raw_df = spark.read.option("multiline", True).json("/mnt/raw/sales.json")

raw_df.write.format("delta").mode("overwrite").save(bronze_path + "/sales")

This way, I save the unprocessed data in Delta format.

Step 4: Process and write to Silver layer

```
I read from the Bronze layer, clean the data, and write it into Silver:

bronze_df = spark.read.format("delta").load(bronze_path + "/sales")

silver_df = bronze_df.filter("amount > 0") \
    .withColumnRenamed("cust_id", "customer_id") \
    .dropDuplicates()

silver_df.write.format("delta").mode("overwrite").save(silver_path + "/cleaned_sales")
```

silver_df.write.format("delta").mode("overwrite").save(silver_path + "/cleaned_sales"

Now I have validated data ready for analytics.

Step 5: Aggregate and write to Gold layer

Step 6: Schedule each layer using Databricks Jobs

I typically create a separate notebook for each layer:

- Notebook 1: Bronze ingestion
- Notebook 2: Silver cleaning
- Notebook 3: Gold aggregation

Then I schedule them in Databricks Jobs, making sure they run in the correct order, or use Job workflows with task dependencies.

Step 7: Add data validation and logging (optional but important)

```
To make the pipeline robust, I add row counts, schema checks, and logs:

print("Bronze row count:", raw_df.count())

print("Silver row count:", silver_df.count())

If needed, I add alerts using email or webhooks if row count suddenly drops or data fails
```

validation.

Step 8: Track data lineage using Unity Catalog or tables

If Unity Catalog is enabled, I register each Delta table in a schema so teams can discover and use them easily.

CREATE TABLE bronze.sales_data

USING DELTA

LOCATION '/mnt/datalake/bronze/sales'

In summary:

To implement Medallion Architecture in Databricks, I:

- 1. Define three layers: Bronze, Silver, Gold
- 2. Ingest raw data into the Bronze layer as-is
- 3. Clean and transform the data in the Silver layer
- 4. Create business-ready aggregated data in the Gold layer
- 5. Use Delta Lake for ACID and performance
- 6. Automate everything using Databricks Jobs
- 7. Add logging, validation, and alerting for reliability

This structure helps keep the data clean, scalable, and easy to debug or reuse across multiple teams.

Scenario 29. How do you implement data lineage tracking in Databricks?

If I'm asked this question in an interview, I would explain that data lineage means being able to trace the flow of data from its source, through all the transformations, and finally to the destination. In Databricks, there are both manual and automated ways to implement data lineage depending on whether Unity Catalog is enabled or not.

Here is how I would answer step by step, using simple words:

Step 1: Use Unity Catalog for built-in lineage (if available)

If Unity Catalog is enabled in the workspace, I prefer using it because Databricks automatically captures lineage for:

- SQL queries
- Delta Live Tables
- Workflows
- Notebook operations

How it works:

Whenever I run a SQL query or notebook that reads from or writes to a table registered in Unity Catalog, Databricks automatically logs:

- Which table or file was used as input
- What logic was applied (joins, filters, etc.)
- What table or view was created or updated

I can go to Data > Lineage tab in the Databricks UI and visually see where the data came from and where it went.

This makes it super easy to track lineage without writing any extra code.

Step 2: Register all Delta tables in Unity Catalog (for better tracking)

To make this work, I register all my Delta tables in Unity Catalog. For example:

CREATE TABLE my_catalog.my_schema.bronze_sales

USING DELTA

LOCATION '/mnt/datalake/bronze/sales_data'

Now every time I query or update this table, the lineage is captured.

Step 3: Use Delta Live Tables (DLT) for real-time lineage

If I'm using Delta Live Tables, lineage is tracked automatically. I just define a DLT pipeline with decorators like @dlt.table and Databricks shows a **visual flow chart** of the entire pipeline.

```
Example:
import dlt

from pyspark.sql.functions import *

@dlt.table

def bronze_sales():
    return spark.read.format("json").load("/mnt/raw/sales.json")

@dlt.table

def silver_sales():
    df = dlt.read("bronze_sales")
```

After running the pipeline, I go to the DLT UI and see exactly which table depends on which, and what transformations are applied. This is very useful for debugging and compliance.

Step 4: For manual lineage (if Unity Catalog is not used)

If Unity Catalog is not available, I use manual logging:

return df.filter(col("amount") > 0)

- I keep a metadata log table where I log each transformation.
- I use comments in notebooks to document where data is coming from and going.
- I tag output files or Delta tables with metadata like source, job name, and timestamp.

Example log entry:

```
log_df = spark.createDataFrame([{
    "source_table": "raw.sales_json",
    "target_table": "silver.cleaned_sales",
    "job_name": "clean_sales_notebook",
    "timestamp": current_timestamp()
}])
log_df.write.mode("append").saveAsTable("lineage.audit_log")
```

Step 5: Use third-party tools (optional)

In some projects, I also connect Databricks to tools like:

- Apache Atlas
- Purview (Azure Data Catalog)
- Collibra or Alation

These tools can fetch metadata from Unity Catalog or scan Delta tables and show visual lineage.

In summary:

To track data lineage in Databricks, I use the following approaches:

- 1. Unity Catalog: built-in, automatic lineage tracking for SQL, notebooks, and jobs
- 2. Delta Live Tables: visual, automatic lineage for pipelines
- 3. Register Delta tables in Unity Catalog to make lineage visible
- 4. Manual logging in audit tables when automation is not available
- 5. Third-party tools like Purview or Apache Atlas for extended cataloging

This helps ensure full traceability from raw data to final output, which is very important for debugging, audits, and data governance.

Scenario 30. How do you migrate notebooks from one Databricks workspace to another?

If I'm asked this question in an interview, I would say that there are a few different ways to migrate notebooks from one Databricks workspace to another, depending on whether I want to do it manually or automate it. The goal is to move the .dbc or .ipynb files, along with any associated dependencies, in a smooth and consistent way.

Here is how I would explain it step-by-step in simple words:

Step 1: Decide the method of migration

There are mainly three ways to migrate notebooks:

- 1. Manual Export and Import (using UI)
- 2. Using Databricks CLI (Command Line)
- 3. Using Repos (if both workspaces are connected to Git)

Method 1: Manual Export and Import (quick and easy)

This is good for small-scale migration.

Step 1.1: Export the notebook from the source workspace

- I go to the workspace folder.
- Right-click on the notebook.
- Choose Export → choose format (HTML, IPython, or DBC).

I usually choose DBC if I want to export multiple notebooks.

Or IPython (.ipynb) for one-by-one migration.

This downloads the file to my local machine.

Step 1.2: Import the notebook into the target workspace

- I log in to the destination workspace.
- Click on the Workspace section.
- Choose Import.
- Upload the .dbc or .ipynb file.
- Select the location to save it (e.g., under /Shared or /Users).

This imports the notebook and I can start using it right away.

Method 2: Using Databricks CLI (for automation and scripting)

This method is useful when I want to automate the migration.

Step 2.1: Install the Databricks CLI

On my local machine, I install it using:

pip install databricks-cli

Then configure it for both source and target workspaces:

databricks configure --profile source

databricks configure --profile target

Each command asks for a host URL and a personal access token from both workspaces.

Step 2.2: Export the notebook using CLI

databricks workspace export_dir /Workspace/MyProject ./myproject_backup --profile source

This command downloads all notebooks from the folder /Workspace/MyProject to a local folder.

Step 2.3: Import into target workspace

databricks workspace import_dir ./myproject_backup /Workspace/MyProject --profile target This uploads the same files into the second workspace.

Method 3: Using Git Integration with Repos

This is the cleanest and most collaborative way.

Step 3.1: Connect notebooks to a Git repo in the source workspace

- I go to the Repos section in Databricks.
- Click Add Repo and connect it to my GitHub/Azure DevOps repo.
- Commit and push all notebooks to Git.

Step 3.2: Clone the same repo in the destination workspace

- Go to Repos in the second workspace.
- Click Add Repo and enter the same Git URL.
- Now I have the same notebooks in both places, and I can sync changes by pushing and pulling.

This also helps in version control and team collaboration.

In summary:

To migrate notebooks between two Databricks workspaces, I:

- 1. Use Export/Import from UI for small quick jobs.
- 2. Use Databricks CLI for scripted or bulk migrations.
- 3. Use Git (Repos) for ongoing collaboration and version control.

If the notebooks depend on libraries, I also make sure to install the same libraries or use the same cluster setup in the new workspace for consistency.

Scenario 31. How do you optimize a Spark job with heavy shuffle operations?

If I'm asked this in an interview, I would explain that shuffle is one of the most expensive operations in Spark, and it can slow down the entire job if not handled properly. Shuffle happens when Spark needs to move data between partitions across the cluster usually during operations like groupBy, join, distinct, or repartition.

When I see that my Spark job has heavy shuffle (especially in Databricks), I follow these steps to optimize it:

Step 1: Understand where and why shuffle is happening

I start by identifying where shuffle is being triggered in my code. Common causes include:

- groupBy, groupByKey
- join without proper partitioning
- distinct
- orderBy, sortBy
- repartition

To detect shuffle, I use:

- Spark UI: After running a job, I go to the Spark UI and look at the stages. If I see high shuffle read/write or spill to disk, that's a sign of heavy shuffle.
- Query plan: I use .explain() on my DataFrame to see if there's an exchange (which means shuffle).

Step 2: Reduce data skew (if shuffle is caused by uneven data)

If one key has too much data (data skew), it causes some partitions to become huge while others stay small. This causes slowdowns.

To fix it:

• I add salt keys to break the skewed key into smaller groups.

Example:

from pyspark.sql.functions import col, lit, concat, rand

df = df.withColumn("salted_key", concat(col("key"), lit("_"), (rand() * 10).cast("int")))

This way, the key "A" becomes "A_1", "A_2", ..., distributing the load better across partitions.

Step 3: Use broadcast joins when one dataset is small

Joins are one of the biggest causes of shuffle. If one side of the join is small (say a few MBs), I broadcast it:

from pyspark.sql.functions import broadcast

result = large_df.join(broadcast(small_df), "id")

This avoids shuffle by sending the small table to all nodes.

Step 4: Repartition wisely (not too much, not too little)

Sometimes I repartition data to control the number of partitions before a join or groupBy.

- If partitions are too many, it increases shuffle and overhead.
- If too few, not all cores are used and some tasks take too long.

I try to keep partition size between 100MB to 200MB.

Example:

df = df.repartition(200, "id") # repartition by key before a join

Or I use coalesce() to reduce partitions after shuffle:

df = df.coalesce(50)

Step 5: Cache intermediate results to avoid re-computation

If I'm using the same transformed data multiple times, I cache it to avoid repeating shuffle:

df = df.cache()

df.count() # triggers the caching

This helps when I'm using the same data in multiple downstream operations.

Step 6: Avoid wide transformations when possible

Wide transformations (like groupByKey) cause a lot of shuffle. I try to use better alternatives:

- Use reduceByKey or aggregateByKey instead of groupByKey (in RDD).
- Use mapPartitions if I can do operations without needing a full shuffle.

Step 7: Use better partitioning strategy for joins

If both datasets are big, I repartition both using the same key before joining:

df1 = df1.repartition("id")

df2 = df2.repartition("id")

result = df1.join(df2, "id")

This ensures the same keys go to the same partition, reducing shuffle across the network.

Step 8: Monitor the shuffle using Spark UI

In Databricks, I always check:

- Shuffle Read/Write size in stages
- Task duration: skewed tasks run much longer
- Shuffle spill (memory/disk): if it's too high, I increase executor memory

Step 9: Tune Spark configs for shuffle performance

If needed, I also tune configs:

spark.conf.set("spark.sql.shuffle.partitions", "200")

spark.conf.set("spark.sql.adaptive.enabled", "true")

spark.conf.set("spark.sql.adaptive.shuffle.targetPostShuffleInputSize", "134217728") # 128MB

- I enable adaptive query execution (AQE) to automatically optimize shuffle.
- AQE can coalesce partitions and switch to broadcast joins automatically.

In summary:

To optimize a Spark job with heavy shuffle in Databricks, I:

- 1. Use Spark UI to find where shuffle is happening
- 2. Fix data skew by adding salt keys
- 3. Use broadcast joins when possible
- 4. Repartition data based on key and cluster size
- 5. Cache intermediate results to reduce shuffle
- 6. Avoid wide transformations and use efficient aggregations
- 7. Use adaptive query execution for auto-optimization

This keeps my job fast, resource-efficient, and avoids failures due to memory pressure or long task durations.

Scenario 32. How do you schedule a Databricks job to run automatically every day?

If I'm asked this question in an interview, I would say that scheduling a Databricks job to run every day is very simple and can be done directly from the Databricks UI using the Jobs feature. Databricks also allows more advanced scheduling through REST APIs or orchestration tools like Azure Data Factory, but for a daily job, the built-in scheduler is usually enough.

Here's how I would explain it step-by-step using simple words:

Step 1: Go to the Jobs section in Databricks

- From the left-hand menu in the Databricks workspace, I click on "Jobs".
- Then I click on "Create Job".

Step 2: Give the job a name

• I provide a name like "Daily_ETL_Job" or "Run_Sales_Processing" so it's easy to identify later.

Step 3: Select the task to run

In the Tasks section:

- I click "Add Task".
- I choose the notebook that I want to run (or I can choose a JAR, Python script, or SQL file).
- I also choose the cluster:

Either an existing all-purpose cluster

Or a new job cluster that is created just for this job run (recommended for cost saving)

I can also pass parameters to the notebook here if needed.

Step 4: Set the schedule to run daily

Now I go to the "Schedule" section of the job:

- I turn on the "Schedule" toggle.
- Set the frequency to "Daily".
- Choose the time of day when I want it to run (like 1:00 AM UTC).
- If I want the job to run on weekdays only, I can adjust the cron expression.

For example:

- Daily at 2 AM UTC → 0 2 * * *
- Every Monday at 7 AM → 0 7 * * 1

Databricks has a simple cron editor built-in for this.

Step 5: Set alerts or email notifications (optional)

If I want to get an email when the job fails or succeeds:

- I go to the "Notifications" section
- Add my email under "On failure" or "On success"

Step 6: Save the job

• Once everything is set, I click "Create" to save and activate the job.

Databricks will now run this notebook automatically every day at the scheduled time.

Step 7 (optional): Use Azure Data Factory or Logic Apps for more control

If the company uses Azure Data Factory, I can also trigger Databricks jobs from ADF pipelines using a Web activity or Databricks Notebook activity. This gives more control, such as:

- Running after a specific data file arrives
- Triggering multiple jobs in sequence
- Handling retries or dependencies

In summary:

To schedule a Databricks job to run automatically every day:

- 1. I go to the Jobs section in Databricks
- 2. Create a job with a name and select a notebook to run
- 3. Set the cluster, input parameters, and retry policy
- 4. Turn on scheduling and set it to run daily using time or cron
- 5. Optionally add email alerts for success or failure
- 6. Save and let Databricks handle the automation

This way, my notebook runs daily without manual intervention, and I can monitor or debug it easily using the built-in job history.

Scenario 33. How do you set up a CI/CD pipeline for deploying Databricks notebooks?

If I'm asked this in an interview, I would explain that CI/CD for Databricks means automating the development, testing, and deployment of notebooks using tools like Git, Azure DevOps, GitHub Actions, or Jenkins. The main goal is to keep the code versioned, tested, and automatically deployed to different environments like dev, test, and prod.

Here's how I wouldset it up step by step in simple words:

Step 1: Store notebooks in Git (Version Control)

The first thing I do is connect my Databricks workspace to a Git repository like GitHub, Azure DevOps, or Bitbucket.

- In Databricks, go to User Settings > Git Integration.
- Add my Git provider and personal access token.
- Now in the notebook UI, I can "Link to Git" and commit changes directly from Databricks.

This keeps my notebooks in source control so I can track versions, roll back, and collaborate with others.

Step 2: Organize the repository

my-databricks-project/

In my Git repo, I keep the project clean like this:

This makes it easy to test, deploy, and manage everything as code.

Step 3: Set up unit tests (optional but good practice)

```
If I have reusable logic in .py files, I write unit tests using pytest.

Example test:
```

```
def test_transform_data():
    df = spark.createDataFrame([...])
    result = transform_data(df)
    assert result.count() == 5
```

I run tests in CI before deployment.

Step 4: Set up CI/CD tool (Azure DevOps or GitHub Actions)

Now I build a pipeline that runs every time someone pushes code to Git.

For GitHub Actions:

```
I create a file like .github/workflows/deploy.yml:
```

name: Deploy to Databricks

on:

push:

branches:

- main

jobs:

deploy:

runs-on: ubuntu-latest

steps:

- name: Checkout code

uses: actions/checkout@v3

- name: Install dependencies

run: pip install databricks-cli

- name: Deploy notebooks

env:

DATABRICKS_HOST: \${{ secrets.DATABRICKS_HOST }}

DATABRICKS_TOKEN: \${{ secrets.DATABRICKS_TOKEN }}

run: |

databricks workspace import_dir notebooks /Workspace/ETL --overwrite

- I store my Databricks host URL and token as GitHub secrets.
- The command uploads notebooks from my repo into the Databricks workspace.

Step 5: Use Databricks CLI to deploy notebooks

The databricks-cli is the main tool I use in the pipeline to move code.

Example command:

databricks workspace import_dir ./notebooks /Workspace/ETL --overwrite

This copies all notebooks to the specified folder in Databricks.

Step 6: Use job JSON configs to deploy jobs too (optional)

If I have jobs defined in JSON, I can deploy them using:

databricks jobs create -- json-file jobs/job-configs.json

Or update an existing job:

databricks jobs reset -- job-id 1234 -- json-file jobs/job-configs.json

This lets me manage job definitions like schedule, cluster, and notebook path as code.

Step 7: Use environments (dev, test, prod)

To keep things organized, I use separate workspaces or folders for each environment. In the pipeline, I control deployment like this:

- Push to dev branch → deploy to dev workspace
- Push to main branch → deploy to prod workspace

I use different tokens and workspace URLs for each environment in the CI/CD pipeline.

In summary:

To set up a CI/CD pipeline for Databricks notebooks:

- 1. I store notebooks in Git (GitHub or Azure DevOps)
- 2. Use Databricks Repos or link notebooks to Git
- 3. Write tests if needed (using pytest)
- 4. Use a CI/CD tool like GitHub Actions or Azure DevOps
- 5. Use databricks-cli in the pipeline to upload notebooks and configure jobs
- 6. Use secrets for authentication and deploy to correct environment
- 7. Automate deployment on every push or pull request

This helps me move faster, reduce errors, and make sure the same code runs consistently in all environments.

Scenario 34. How do you troubleshoot a Databricks job failing due to long-running shuffle operations?

If I'm asked this in an interview, I would say that shuffle operations in Spark happen when data needs to be moved across partitions, such as during joins, groupBy, orderBy, or distinct operations. When these shuffle operations take too long or cause failures, I follow a structured approach to troubleshoot and optimize them.

Here is how I would explain it step by step in simple language:

Step 1: Understand the type of shuffle causing the issue

First, I look at what transformation is triggering the shuffle. Common operations that cause shuffling are:

- groupBy()
- join()
- orderBy()
- distinct()
- repartition()

So I review the notebook or job code to find these operations and note where the issue might be.

Step 2: Check the Spark UI (Job Details)

Then, I go to the Spark UI from the job run in Databricks:

- I open the failed job and click on the "Spark UI" tab.
- Under the "Stages" tab, I look for:
 - Long-running stages
 - Stages with high task time
 - Shuffle read/write size
 - Number of failed/retried tasks

This helps me identify which stage is the bottleneck and whether the shuffle is causing skew or memory problems.

Step 3: Look for data skew

A common reason for long shuffle operations is data skew, which happens when one or a few partitions contain most of the data. I check the task time distribution in Spark UI:

- If one task takes 10 minutes and others take 10 seconds, it's skew.
- I also check if one key has a very high number of records.

To fix skew, I use techniques like:

- salting the key before the shuffle
- using broadcast join if one table is small
- using repartition() to balance data

Step 4: Use broadcast joins when applicable

If I'm doing a join between a large and a small dataset, and the small dataset fits in memory, I change the join to a broadcast join.

from pyspark.sql.functions import broadcast

df_large.join(broadcast(df_small), "key")

This avoids shuffling the large table and improves performance.

Step 5: Repartition data properly

If the data is badly partitioned (e.g., one big partition), I repartition it using:

df = df.repartition(200) # or use a column: df.repartition("customer_id")

But I avoid too many partitions, because too many small tasks can also slow down the job.

I also avoid unnecessary coalesce() or repartition() if not needed.

Step 6: Increase shuffle partitions if needed

By default, Spark may not create enough shuffle partitions. I can increase it by setting:

spark.conf.set("spark.sql.shuffle.partitions", "300")

This spreads shuffle data across more partitions and reduces pressure on each task.

Step 7: Tune the cluster size

Sometimes, the issue is not with the code but with the cluster size. So I:

- Increase the driver and worker memory
- Use more executors or cores
- Use auto-scaling if needed

Also, I enable adaptive query execution (AQE):

spark.conf.set("spark.sql.adaptive.enabled", "true")

This allows Spark to optimize shuffle partitions and join strategies during runtime.

Step 8: Cache intermediate results if reused

If I'm using the same DataFrame multiple times after a shuffle, I cache it:

df.cache()

This avoids repeating the shuffle again.

In summary:

When a Databricks job fails or slows down due to shuffle operations, I do the following:

- 1. Check which transformation (e.g., join, groupBy) is causing the shuffle.
- 2. Use the Spark UI to identify slow stages and data skew.
- 3. Fix skew by salting, broadcasting, or repartitioning data.
- 4. Use broadcast joins if one table is small.
- 5. Adjust spark.sql.shuffle.partitions if needed.
- 6. Optimize cluster size and enable AQE.
- 7. Cache intermediate DataFrames to avoid repeated shuffling.

By following these steps, I can usually find the root cause and improve job performance significantly.

Scenario 35. How do you use a bronze-silver-gold pipeline to clean and curate data?

If I'm asked this question in an interview, I would explain that the bronze-silver-gold pipeline is a design pattern used in Databricks Medallion Architecture. It helps organize and manage data in layers to make it easy to process, clean, and analyze. I've used this pattern to build scalable and structured pipelines in real projects.

Here's how I would describe it step-by-step in simple terms:

Step 1: Understand the purpose of each layer

Bronze Layer (Raw Data):

This is where I land the raw data from source systems (like Azure Data Lake, Kafka, or REST APIs). No cleaning is done here. It's mainly used for backup and traceability.

Silver Layer (Cleaned and Filtered):

In this layer, I clean the data, apply business rules, and transform it into usable formats. Nulls are handled, types are cast, and duplicates are removed.

Gold Layer (Aggregated and Curated):

This is where I prepare final datasets for reporting, dashboards, machine learning, or business use. Data is grouped, joined, and enriched for specific use cases.

Step 2: Example use case

Let's say I'm working with sales transaction data.

- Source: CSV files are dropped into Azure Data Lake daily from an external sales system.
- Goal: Clean it, remove duplicates, and prepare daily revenue reports.

Step 3: Bronze layer - Ingest raw data

In this step, I load raw data from Azure Data Lake and write it into the bronze Delta table.

```
raw_df = spark.read.option("header", True).csv("/mnt/raw/sales")
```

raw_df.write.format("delta").mode("append").saveAsTable("bronze.sales_raw")

No filtering or transformation here. This keeps a raw backup.

Step 4: Silver layer - Clean and enrich data

Now I read from the bronze layer, remove bad records, correct data types, and remove duplicates.

```
bronze_df = spark.read.table("bronze.sales_raw")
silver_df = bronze_df.filter("amount IS NOT NULL") \
    .withColumn("amount", bronze_df["amount"].cast("double")) \
    .dropDuplicates(["transaction_id"])
```

silver_df.write.format("delta").mode("overwrite").saveAsTable("silver.sales_cleaned")

Now the data is clean and ready for reporting or joining with other sources.

Step 5: Gold layer - Aggregate and curate for business use

```
Next, I build the gold layer by preparing the final table, such as daily revenue by region.

silver_df = spark.read.table("silver.sales_cleaned")

gold_df = silver_df.groupBy("region", "date").agg(
    sum("amount").alias("total_sales")

)

gold_df.write.format("delta").mode("overwrite").saveAsTable("gold.daily_sales_report")

This gold table is optimized for Power BI dashboards or finance reports.
```

Step 6: Schedule and automate the pipeline

To make this run every day, I use Databricks Jobs with three tasks:

- Task 1: Bronze ingestion
- Task 2: Silver cleaning (depends on Task 1)
- Task 3: Gold aggregation (depends on Task 2)

I also enable alerts and retries to handle any failures.

Step 7: Add logging and quality checks (optional)

To make it production-ready, I add:

- Logging (e.g., log row counts, errors)
- Data quality checks (e.g., check if any nulls in key columns)
- Monitoring dashboards (e.g., success/failure rates)

In summary:

To use a bronze-silver-gold pipeline in Databricks:

- 1. Bronze: Ingest raw data without changes.
- 2. Silver: Clean and transform data (remove nulls, fix types, deduplicate).
- 3. Gold: Aggregate and prepare data for reporting or ML.

This design helps keep raw data safe, cleaned data reusable, and final data optimized for consumption. It also makes the pipeline scalable and easy to debug.

Scenario 36. What are some challenges you have faced when working with large datasets in Databricks, and how did you overcome them?

If I'm asked this question in an interview, I would answer honestly and practically by talking about real challenges that are common when handling large data in Databricks. Then I'd walk through the steps I took to solve them, using very simple language.

Here's how I would answer:

One of the biggest challenges I faced while working with large datasets in Databricks was performance and memory issues during complex transformations, especially joins and aggregations.

For example, I was working with a dataset of over 500 million customer transactions stored in Azure Data Lake. I had to join this with customer profiles and then calculate total spending by category.

The job was failing with "out of memory" errors and taking more than an hour to complete.

Challenge 1: Data skew during joins

I found that one of the tables (the customer transactions) had data skew, meaning most of the records belonged to just a few customers. This caused one task to process too much data, while others finished quickly.

Solution:

- I added salting to break the skewed key into multiple keys.
- I also used broadcast joins when the lookup table (like customer info) was small enough to fit in memory:

from pyspark.sql.functions import broadcast

df_result = df_transactions.join(broadcast(df_customers), "customer_id")

This reduced the shuffle and made the job run much faster.

Challenge 2: Cluster was too small

Another issue was that the cluster didn't have enough memory to handle such large data.

Solution:

- I increased the worker node size and used autoscaling clusters in Databricks.
- I also enabled photon runtime which helped improve performance significantly.
- In addition, I made sure to use job clusters instead of interactive ones to avoid memory leaks from previous runs.

Challenge 3: Too many small files in the data lake

The source data had thousands of small files. This slowed down the read performance and increased the number of tasks.

Solution:

- I used auto-merge and optimized writes in Delta Lake to compact the files.
- I also ran OPTIMIZE command after loading the bronze layer:

OPTIMIZE bronze.transactions

This helped reduce I/O and improved processing speed.

Challenge 4: Long shuffle operations

Some transformations like groupBy and orderBy were triggering large shuffle operations.

Solution:

• I tuned the configuration:

spark.conf.set("spark.sql.shuffle.partitions", 400)

spark.conf.set("spark.sql.adaptive.enabled", "true")

- Adaptive Query Execution helped Spark adjust the plan during runtime.
- I also used repartitioning where necessary to balance the data:

df = df.repartition("region")

Challenge 5: Slow data exploration in notebooks

When using very large datasets, sometimes the notebook became slow or unresponsive.

Solution:

I always worked with sampled data for testing:

 $df_sample = df.limit(1000)$

• And for heavy jobs, I used Databricks Jobs instead of running it manually in notebooks.

In summary:

Some of the challenges I faced with large datasets in Databricks were:

- 1. Data skew during joins
- 2. Memory and cluster resource limits
- 3. Too many small files in the data lake
- 4. Shuffle operations causing slow stages
- 5. Slow performance in notebooks

I overcame them by using broadcast joins, cluster tuning, adaptive query execution, optimizing files, and using job clusters. These techniques helped me make the pipelines more stable, faster, and cost-effective.

Scenario 37. What are some strategies for cost optimization in Databricks?

If I'm asked this question in an interview, I would answer by explaining that Databricks gives powerful tools, but it can get expensive if we don't manage resources properly. So I always try to make my notebooks, jobs, and clusters as efficient as possible.

Here are the cost optimization strategies I've used, explained step by step in simple words:

1. Use job clusters instead of all-purpose clusters

- All-purpose clusters stay running even when not used, which means you get charged continuously.
- Job clusters are created for each job and automatically terminate when the job finishes.
- I use job clusters for scheduled pipelines and batch jobs to reduce idle costs.

2. Enable auto termination

- I always set the auto termination feature on clusters.
- This shuts down the cluster automatically after a set period (e.g., 10 or 15 minutes) of inactivity.
- It avoids unnecessary charges when someone forgets to shut it down.

3. Use autoscaling

- I enable autoscaling on clusters so the number of worker nodes adjusts automatically based on the workload.
- During heavy processing, more nodes are added. When it's idle or light, fewer nodes are used.
- This avoids over-provisioning and reduces costs.

Example from the cluster config:

Min workers: 2

Max workers: 8

4. Choose the right cluster size and instance type

- I avoid using very large or expensive VM types unless required.
- For most ETL jobs, a medium-size worker node (like Standard_D4_v2) is enough.
- For memory-heavy jobs, I use memory-optimized VMs, but only when necessary.

5. Use Photon runtime

- When Photon is enabled on a cluster, Spark queries run faster using vectorized execution.
- This reduces compute time, which directly reduces cost.
- I've seen up to 30-50% performance gain with Photon for certain workloads.

I enable it while creating the cluster by selecting the Photon runtime option.

6. Use Delta Lake and Z-ordering

- Delta Lake improves performance by handling data updates efficiently and avoiding full table scans.
- I also use ZORDER BY to optimize read performance for specific columns.

OPTIMIZE sales_data ZORDER BY (customer_id)

This helps speed up queries and reduce the cost of long reads.

7. Compact small files using OPTIMIZE

- Reading from many small files slows down performance and increases costs.
- I use the OPTIMIZE command regularly to compact data in Delta tables:

OPTIMIZE bronze.transactions

• This reduces shuffle and I/O cost.

8. Use caching carefully

- I cache data only when the same DataFrame is reused multiple times.
- Caching unnecessary data can fill up memory and increase cluster cost.

df.cache()

9. Monitor usage and costs

- I use Databricks cost reports, cluster metrics, and Azure Cost Management to monitor how much is being spent.
- I also analyze logs to find jobs that run too frequently or take too long, and then optimize them.

10. Avoid frequent small job runs

- Instead of running small jobs every few minutes, I combine them into batch runs every hour or daily.
- This reduces overhead from cluster spin-up time and minimizes billing cycles.

In summary:

To keep Databricks costs low, I use strategies like:

- 1. Job clusters with auto termination
- 2. Autoscaling and correct VM sizes
- 3. Using Photon for faster jobs
- 4. Delta Lake with ZORDER and OPTIMIZE
- 5. Monitoring and adjusting job frequency
- 6. Caching smartly and combining small tasks

By following these, I make sure that my Databricks projects are efficient and budget-friendly while still delivering good performance.

Scenario 38. What is your strategy for migrating an on-premises Hadoop workload to Azure Databricks?

If I get this question in an interview, I would answer by explaining that I follow a step-by-step migration plan. I always start by understanding the existing system and then carefully modernize it using Azure services and Databricks. Here's how I usually approach such a migration using very simple and practical steps:

Step 1: Understand the existing Hadoop setup

Before starting the migration, I make sure to gather all the information about the current Hadoop system:

- What kind of data is being processed?
- What tools are used (Hive, MapReduce, Sqoop, Pig, Oozie)?
- What are the current workflows and schedules?
- Where is the data stored (HDFS, local storage)?
- Are there any custom scripts or third-party tools?

This helps me plan the migration properly without missing anything.

Step 2: Identify what can be replaced by Databricks and Azure services

Then, I map the existing Hadoop components to modern Azure tools:

On-prem Hadoop Azure Databricks Replacement

HDFS Azure Data Lake Storage (ADLS Gen2)

Hive or Impala Databricks with Spark SQL and Delta Lake

MapReduce PySpark in Databricks

Sqoop for ingestion Azure Data Factory, or Databricks JDBC reads

Oozie for scheduling Azure Data Factory pipelines, or Databricks Jobs

Pig scripts PySpark equivalent scripts in Databricks

This helps simplify the overall architecture and makes it cloud-native.

Step 3: Set up the Azure environment

I prepare the environment where the workloads will be moved:

- Create an Azure Data Lake Gen2 account for storage.
- Set up Azure Databricks workspace.
- Configure networking, RBAC, and secure access to external sources.
- If needed, enable Git integration, Unity Catalog, and Key Vault for secrets.

Step 4: Migrate the data

I move all the HDFS data into Azure Data Lake using tools like:

- AzCopy (for static files),
- DistCp (for large volumes),
- or Azure Data Factory (for pipeline-based transfers).

While copying data, I also ensure:

- Folder structure is maintained,
- File formats are preserved or upgraded to Parquet or Delta for better performance.

Step 5: Convert Hive or MapReduce jobs to PySpark in Databricks

- I go through each Hive or MapReduce job and convert the logic into PySpark or Spark SQL notebooks in Databricks.
- For example, a Hive query like:

SELECT customer_id, SUM(amount) FROM transactions GROUP BY customer_id;

Can be converted to:

df = spark.read.format("parquet").load("/mnt/datalake/transactions")

df.groupBy("customer_id").agg(sum("amount")).show()

Step 6: Rebuild workflows and scheduling

If the original system used Oozie or cron jobs, I move the scheduling logic to:

- Databricks Jobs (for notebook workflows), or
- Azure Data Factory (for orchestration across multiple services).

I recreate the same dependencies and order using task dependencies in Databricks Jobs.

Step 7: Testing and validation

I test the new Databricks workflows by comparing:

- Row counts,
- Aggregated results,
- Data quality reports,

...against the old system. I do parallel runs for a few days or weeks until I'm confident that the results match.

Step 8: Performance tuning and optimization

Once everything is working, I improve the pipeline by:

- Using Delta Lake for fast writes and updates,
- Partitioning data properly,
- Using Photon runtime and autoscaling clusters,
- Adding ZORDER for better file pruning.

Step 9: Monitoring and logging

I add monitoring for:

- Job failures,
- Performance metrics,
- Cost usage,

...using Azure Monitor, Databricks Job logs, and custom logging in notebooks.

Step 10: Decommission on-prem resources

After testing is complete and the Databricks version is stable, I shut down the old Hadoop jobs and document the new process fully.

In summary:

To migrate an on-prem Hadoop workload to Azure Databricks, I follow this strategy:

- 1. Analyze the current system (jobs, tools, and data).
- 2. Map each Hadoop component to an Azure equivalent.
- 3. Migrate data to Azure Data Lake.
- 4. Rewrite MapReduce/Hive logic using PySpark in Databricks.
- 5. Rebuild and schedule workflows using Databricks Jobs or ADF.
- 6. Test and compare results.
- 7. Optimize performance using Delta Lake and Photon.
- 8. Monitor jobs and track costs.
- 9. Fully decommission the old system.

This step-by-step method ensures the migration is smooth, efficient, and cloud-optimized.

Scenario 39. What metrics are important for monitoring Databricks jobs?

If I'm asked this in an interview, I would explain that monitoring is very important for keeping Databricks jobs reliable and efficient. I always keep an eye on both job-level and cluster-level metrics to make sure everything runs smoothly, quickly, and without wasting resources.

Here's how I usually monitor Databricks jobs, explained in simple and detailed steps:

1. Job Status and Duration

- I always check whether the job succeeded or failed.
- I also look at how long the job took compared to the past runs.
- This helps me catch performance issues or unexpected changes.

Where I check it:

From the Jobs UI in Databricks → It shows the run history, success/failure, and time taken.

2. Task Execution Time

- I monitor how long each notebook or task in the job is taking.
- If one task suddenly takes more time, it could mean data skew or inefficient code.

Example:

If Task 2 usually takes 2 minutes but now takes 10 minutes, I investigate that step.

3. Cluster CPU and Memory Usage

From the Spark UI or Ganglia metrics, I monitor:

- CPU utilization: High usage might mean the code is compute-heavy.
- Memory usage: If memory usage is near 100%, I may hit out-of-memory errors.

I use this to:

- Decide whether to scale up the cluster,
- Or optimize the code (like avoiding unnecessary caching or joins).

4. Shuffle Read and Write

This is one of the most critical metrics in Spark jobs.

- Shuffle operations happen when data is moved between nodes (like in joins or groupBy).
- If shuffle size is very large, the job may become slow or even fail.

In the Spark UI, I check:

- Shuffle Read Size
- Shuffle Write Size
- Number of spilled records (means it couldn't fit into memory)

If I see very large shuffle sizes, I may:

- · Repartition the data,
- Broadcast small tables,
- Or reduce wide transformations.

5. Stage and Task Failures

I monitor:

- Number of failed tasks
- Number of retries per task

If tasks are failing frequently, it could be due to:

- Bad data,
- Memory errors,
- Or cluster configuration issues.

Databricks retries tasks by default, but too many retries usually signal deeper problems.

6. Data Skew

In the Spark UI, I check:

- Task duration across partitions
- Size of records processed per task

If one task takes much longer or processes more data than others, it's a sign of data skew.

To fix skew, I may:

- Use salting or skew hints,
- Avoid large joins on skewed columns.

7. Input and Output Data Volume

I monitor how much data is being:

- Read into the job,
- Written out to storage (like ADLS or Delta Lake).

Sudden spikes may mean schema changes or bad upstream data.

I also check the number of files written, since writing too many small files can cause performance problems later.

8. Auto-scaling behavior (if enabled)

If autoscaling is on, I track:

- Whether the cluster scaled up or down correctly.
- How many workers were used at peak.

This helps me balance performance and cost. If the cluster is not scaling up when needed, I check if:

- The max worker limit is too low,
- Or if there's a slow task blocking the scale-up.

9. Delta Lake-specific metrics (if using Delta)

When using Delta Lake, I check:

- Number of files written or rewritten
- Compaction statistics after OPTIMIZE
- ZORDER stats if indexing is used

This helps me tune storage performance.

10. Logs and Alerts

I always check:

- Driver logs and executor logs for error messages or stack traces.
- Set up alerts using Azure Monitor or custom scripts if:
 - A job fails,
 - Job duration exceeds expected time,
 - Or costs spike beyond budget.

In summary:

The most important metrics I monitor in Databricks jobs are:

- 1. Job status and duration
- 2. Task execution times
- 3. CPU and memory usage
- 4. Shuffle read/write and spilled data
- 5. Task failures and retries
- 6. Data skew across partitions
- 7. Data volume read/written
- 8. Auto-scaling behavior
- 9. Delta Lake optimization stats
- 10. Logs and alerts for failures

Monitoring these helps me make sure jobs are reliable, fast, and cost-efficient. I use the built-in Jobs UI, Spark UI, Ganglia, and Azure Monitor to keep track of all of these.

Scenario 40. What steps are involved in implementing a data governance strategy in Azure Databricks?

If I'm asked this in an interview, I would explain that implementing a data governance strategy in Azure Databricks is all about making sure data is secure, well-managed, and used responsibly. I always break this into a few simple steps, covering access control, data cataloging, audit logging, and compliance.

Here's how I would explain it in detail, step-by-step:

Step 1: Set up Unity Catalog

The first step is to enable Unity Catalog, which is Databricks' centralized data governance layer. It helps manage:

- Data access control (down to table, column, and row levels),
- Data discovery (metadata, schema, lineage),
- And auditing.

Unity Catalog helps manage permissions across workspaces, users, and groups centrally.

Step 2: Organize data into catalogs, schemas, and tables

In Unity Catalog, I structure the data in a logical way:

- Catalog = top-level container (like a project or department)
- Schema = organizes tables/views (like a folder)
- Table = the actual data

For example:

finance db.transactions 2024

marketing_db.campaign_metrics

This structure helps with both security and discoverability.

Step 3: Set fine-grained access control

Using Unity Catalog and Azure AD groups, I assign role-based access:

- Data scientists get read access to curated tables.
- Analysts may get access to only specific columns (like non-PII).
- Engineers get write access to staging areas only.

Access can be managed using SQL:

GRANT SELECT ON TABLE finance_db.transactions TO `analyst_group`;

Or using the Databricks UI under "Data" tab > Permissions.

Step 4: Mask or encrypt sensitive data (PII, financial, etc.)

To protect sensitive data, I use:

- Column-level masking: Like showing only last 4 digits of a credit card.
- Tokenization or hashing: To anonymize data.
- Row-level security: To restrict certain rows based on user roles.

For example, with dynamic views:

CREATE OR REPLACE VIEW masked_view AS

SELECT

name,

CASE

WHEN is_member('finance_group') THEN salary

ELSE NULL

END AS salary

FROM employees;

Step 5: Track data lineage

Unity Catalog provides automatic lineage tracking, so I can see:

- · Which notebooks or jobs created or modified a dataset,
- · What the upstream sources are,
- And how data flows through bronze, silver, and gold layers.

This helps with auditing, debugging, and impact analysis before making changes.

Step 6: Enable auditing and logging

To keep track of data usage and changes:

- I enable audit logs in Databricks (via diagnostic settings).
- These logs are stored in Azure Log Analytics or Azure Storage.
- I can monitor for:
 - Unauthorized access,
 - Dropped tables,
 - Role changes.

This helps meet compliance requirements like GDPR and HIPAA.

Step 7: Set data quality rules and monitoring

I implement data quality checks using tools like:

- · Great Expectations inside Databricks notebooks,
- Or custom PySpark validation scripts to check for:
 - Null values,
 - Duplicates,
 - Outliers.

Example:

df.filter(col("email").isNull()).count()

I also schedule these checks using Databricks Jobs and alert if issues are found.

Step 8: Use catalogs and tags for discoverability

I tag datasets with metadata to make discovery easier:

- Owner
- Description
- Classification (PII, confidential, etc.)
- Source system

Databricks UI shows these tags, which help users understand what they are working with.

Step 9: Define and enforce naming conventions and lifecycle rules

I maintain proper naming conventions like:

- project_schema.table_name
- bronze_, silver_, gold_ prefixes

And I use table versioning and retention policies to manage old data using:

VACUUM delta_table RETAIN 7 HOURS;

Step 10: Educate users and document policies

Finally, I make sure all team members are trained on:

- Access policies,
- How to request permissions,
- What's allowed and what's restricted.

I usually maintain a wiki or Confluence page explaining:

- How Unity Catalog is structured,
- How to request new data access,
- Security do's and don'ts.

In summary:

To implement a strong data governance strategy in Azure Databricks, I follow these steps:

- 1. Enable Unity Catalog
- 2. Organize data with catalogs, schemas, tables
- 3. Set role-based access using Azure AD
- 4. Mask or encrypt sensitive data
- 5. Use data lineage to track data flow
- 6. Enable audit logging for monitoring
- 7. Add data quality checks and alerts
- 8. Tag data assets for easier discovery
- 9. Define naming/lifecycle conventions
- 10. Document everything and educate users

This helps ensure that data is secure, trustworthy, and easy to manage at scale.

Scenario 41. What steps can be taken to avoid data loss in Databricks?

If I'm asked this question in an interview, I would explain that preventing data loss in Databricks is very important, especially when handling critical business data. In my experience, I follow a set of practical steps that help ensure data is safely stored, backed up, and not accidentally deleted or overwritten.

Here's how I would answer it step by step, in simple words:

Step 1: Always Write Data to Durable Storage like ADLS or S3

Instead of writing data to the local disk on a Databricks cluster (which gets deleted when the cluster shuts down), I always write and read data from durable cloud storage like:

- Azure Data Lake Storage Gen2 (ADLS)
- Or Azure Blob Storage

Example:

df.write.format("delta").save("abfss://datalake@storageaccount.dfs.core.windows.net/bronze/customers")

This way, data is stored safely even if the cluster is deleted or restarted.

Step 2: Use Delta Lake for ACID Transactions and Versioning

Delta Lake helps avoid data loss during write or update operations. It supports:

- ACID transactions which means even if a job fails midway, it won't leave data in an inconsistent or corrupted state.
- Time travel lets me go back to a previous version of the table if something goes wrong.

Example to restore old data:

SELECT * FROM my_table VERSION AS OF 5

Or using timestamp:

SELECT * FROM my_table TIMESTAMP AS OF '2023-08-01T00:00:00Z'

Step 3: Enable Auto-Backup or Create Checkpoints

For long-running streaming or batch processes, I use checkpointing to store progress safely.

In streaming jobs:

df.writeStream \

- .format("delta") \
- .option("checkpointLocation", "/mnt/checkpoints/stream1") \
- .start("/mnt/output")

This helps avoid re-processing data from the beginning in case of failure.

For backup, I sometimes keep a copy of the final Delta tables (especially critical ones) in a backup folder.

Step 4: Use Save Modes Carefully

When writing data, using the wrong save mode like overwrite can delete existing data.

I make sure to choose the correct save mode:

- append to add new data without deleting old
- overwrite only if I really intend to replace data
- errorlfExists safest if I want to avoid accidental overwrite

Example:

df.write.mode("append").format("delta").save("/mnt/delta/sales")

Step 5: Use Unity Catalog for Access Control and Protection

With Unity Catalog, I can:

- Prevent unauthorized users from modifying or dropping tables
- Use role-based access control (RBAC)
- · Apply column-level permissions and row filtering

This helps avoid accidental or unauthorized deletion of data.

Step 6: Use Table Constraints (Optional)

Delta Lake supports basic constraints like NOT NULL and CHECK.

They help in maintaining data integrity and prevent invalid or incorrect data from being inserted, which can reduce the chance of needing to roll back or recover corrupted data.

Example:

CREATE TABLE sales (

id INT,

amount DOUBLE NOT NULL

) USING DELTA;

Step 7: Schedule Regular Backups (for critical data)

For very important datasets, I take a snapshot of the table periodically and store it in another safe location (maybe another container or region).

This is usually done via a scheduled Databricks Job:

Copy data to backup folder

df.write.mode("overwrite").format("delta").save("/mnt/backup/daily/sales_snapshot")

Step 8: Enable Audit Logs to Track Changes

By enabling audit logging using Azure Monitor or Azure Log Analytics, I can track:

- Who deleted or updated data
- When the changes happened

This doesn't prevent data loss but helps in investigations and recovery.

Step 9: Test All Pipelines in Development First

Before running any data pipeline in production, I first test it thoroughly in dev or staging environments with sample data.

This helps catch logic bugs that could cause data loss, like:

- Wrong join keys,
- · Wrong filters,
- Accidental overwrites.

Step 10: Educate Team Members and Set Permissions

Many data loss issues happen due to human error. So, I:

- Restrict write/delete permissions to only those who need it,
- Educate the team on using save modes and overwrite carefully,
- Document best practices.

In summary:

To avoid data loss in Databricks, I follow these steps:

- 1. Store all data in cloud storage (like ADLS, not local disk)
- 2. Use Delta Lake for versioning and recovery
- 3. Set up checkpoints for streaming pipelines
- 4. Choose write/save modes carefully
- 5. Use Unity Catalog for fine-grained access control
- 6. Add table constraints where needed
- 7. Create regular backups for important datasets
- 8. Monitor changes using audit logs
- 9. Always test code in dev before deploying to prod
- 10. Educate the team and apply strict permissions

By combining all these steps, I make sure that the chances of accidental data loss are very low in any Databricks project I work on.

Scenario 42. What would you do if your Databricks job is running out of memory?

If I'm in an interview and asked what I would do when a Databricks job is running out of memory, I would explain that this is a very common issue, especially when dealing with large datasets. When this happens, I try to identify the root cause first and then apply memory optimization techniques step by step.

Here's how I would explain it in simple and detailed steps:

Step 1: Check the job and cluster logs for memory errors

The first thing I do is go to the "Spark UI" from the job run or notebook and look at:

- Stages tab to see if any stage is failing with memory-related errors like:
 - java.lang.OutOfMemoryError
 - GC overhead limit exceeded
- Executors tab to check if any executor has high memory usage or is crashing.

This helps me confirm that the problem is indeed memory-related and not something else.

Step 2: Increase the cluster size or memory

If the cluster has very little memory to start with, I try to increase the driver and executor memory by selecting a more powerful instance type in the cluster configuration.

For example, I might switch from:

Standard DS3 v2 (14 GB RAM) → to Standard DS5 v2 (56 GB RAM)

Or increase the number of worker nodes to distribute the load.

Step 3: Repartition or coalesce the data

One of the most common reasons for memory issues is skewed partitions or too few partitions.

To fix this, I might repartition the DataFrame to spread the data more evenly across executors:

df = df.repartition(200)

Or, if there are too many small partitions and causing shuffle overhead, I use:

df = df.coalesce(50)

The right number depends on the data size. My rule of thumb is to aim for partition sizes of around 100–200 MB.

Step 4: Avoid wide transformations before filtering

Sometimes, we perform joins or aggregations before filtering the data, which uses more memory than necessary.

So, I try to filter the data early, before heavy transformations:

Bad

df.join(large_df).filter(col("status") == "active")

Better

filtered_df = df.filter(col("status") == "active")

filtered_df.join(large_df)

This reduces the amount of data that needs to be kept in memory.

Step 5: Cache only if needed, and uncache when done

Caching is useful, but it also consumes memory. If I'm caching large DataFrames:

df.cache()

I make sure to uncache them when I no longer need them:

spark.catalog.uncacheTable("df")

This frees up memory and avoids unnecessary pressure on the cluster.

Step 6: Use broadcast joins for small tables

When joining a small table with a large table, I use broadcast joins to avoid large shuffles:

from pyspark.sql.functions import broadcast

df.join(broadcast(small_df), "id")

This sends the small table to all executors and prevents Spark from performing a costly shuffle operation.

Step 7: Optimize data formats (use Parquet or Delta)

If I'm reading raw CSV or JSON files, they use more memory. Instead, I convert them to Parquet or Delta format, which is optimized and compressed.

Example:

df.write.format("delta").save("/mnt/delta/optimized_data")

Delta format also supports schema enforcement and indexing, which makes reads faster and lighter on memory.

Step 8: Optimize joins using bucketing or skew hints (for large joins)

If memory issues happen during joins, I sometimes use:

- Bucketing: partition and sort data before join
- Skew hints: to tell Spark that a column has skewed data

Example using skew hint:

```
df1.hint("skew").join(df2, "id")
```

This tells Spark to handle skewed keys more efficiently, avoiding memory spikes on a few nodes.

Step 9: Break the job into smaller stages

If the dataset is too big to process at once, I divide the job into smaller chunks using filters or date partitions.

Example:

```
for year in range(2018, 2023):
df\_year = df.filter(col("year") == year)
```

process(df_year)

This allows Spark to process only a portion of the data in memory at a time.

Step 10: Use Delta caching (if running interactive queries)

If I'm using interactive clusters with Delta tables, I turn on Delta caching, which speeds up reads without putting everything into driver memory.

This is done automatically if I'm using a supported cluster with Photon or optimized autoscaling.

In summary:

If a Databricks job is running out of memory, I follow these steps:

- 1. Check Spark UI logs to confirm memory errors
- 2. Scale up the cluster size or memory
- 3. Repartition or coalesce data based on size
- 4. Filter data early, before joins or aggregations
- 5. Cache only when necessary and uncache after use
- 6. Use broadcast joins for small lookup tables
- 7. Convert raw data to optimized formats like Delta
- 8. Handle skewed joins with hints or bucketing
- 9. Process data in smaller chunks if possible
- 10. Use Delta caching for faster reads (if available)

Scenario 43. What would you do if your Databricks job is slow due to shuffle operations?

If I'm asked this question in an interview, I would explain that shuffle operations are often the main reason for slow Spark jobs in Databricks. Shuffles happen when data is moved between executors across the cluster, usually during operations like joins, groupBy, distinct, and orderBy. They're expensive in terms of time and resources.

Here's how I would approach this problem in a step-by-step way, in simple words:

Step 1: Identify if the shuffle is really the issue

First, I go to the Spark UI of the job run and check the following:

- Look at the Stages tab to see which stages have the longest time.
- If I see operations like shuffle read, shuffle write, or exchange taking a long time or involving large data transfers, then I confirm it's a shuffle issue.
- I also check if the task time is uneven, which could be a sign of data skew.

Step 2: Repartition the data properly

If there are very few partitions or the partitions are unbalanced, Spark can't distribute the shuffle evenly. So I:

• Use repartition() to increase the number of partitions before the shuffle-heavy operation:

df = df.repartition(200) # or based on number of executors

• If I want fewer, larger partitions (after shuffle), I use coalesce():

df = df.coalesce(50)

This helps balance the workload across the cluster.

Step 3: Use broadcast join for small lookup tables

If the shuffle is happening because of a **join**, and one of the tables is small, I use broadcast join to avoid moving large data between executors.

from pyspark.sql.functions import broadcast

large_df.join(broadcast(small_df), "id")

This way, the small table is copied to every executor, and no shuffle is needed.

Step 4: Use bucketing or Z-Ordering for repeated joins

If I need to do repeated joins on the same column (like user_id), I might use bucketing or Z-Ordering (with Delta Lake).

Bucketing helps avoid full shuffles because the data is pre-partitioned.

Example (in Delta):

CREATE TABLE sales_bucketed

USING DELTA

CLUSTERED BY (user_id) INTO 50 BUCKETS

AS SELECT * FROM sales;

Or I use Z-Ordering if I want better data skipping:

OPTIMIZE sales ZORDER BY (user_id);

Step 5: Apply filtering before the shuffle

Sometimes a join or groupBy happens on a very large dataset that could be filtered earlier. So I move filters before the join or groupBy:

Instead of this

df1.join(df2, "user_id").filter(df1["country"] == "US")

I do this

filtered_df1 = df1.filter(df1["country"] == "US")

filtered_df1.join(df2, "user_id")

This reduces the amount of data shuffled.

Step 6: Handle skewed data

If one key (like a specific user ID or customer ID) has a huge amount of data compared to others, that causes data skew and slows down the job.

In that case, I use the skew join hint:

df1.hint("skew").join(df2, "id")

Or, if needed, I use a technique called salting, where I add a random number to the key to distribute it more evenly.

Step 7: Avoid wide transformations if not needed

Operations like groupBy, distinct, orderBy all cause shuffle. I always ask myself: Do I really need this operation?

If I can avoid it, or replace it with an alternative (like reduceByKey or approx_count_distinct), I do that.

Step 8: Tune Spark config settings

If needed, I adjust Spark configurations in the notebook or job cluster:

spark.conf.set("spark.sql.shuffle.partitions", "200")

This controls how many partitions Spark creates after a shuffle. The default is 200, but I tune it based on the data size.

Step 9: Use Photon or optimized compute engines

If the workspace supports it, I use Photon-enabled clusters, which are optimized for better performance and can handle shuffle operations more efficiently.

Step 10: Break down the job into smaller stages

If nothing else works, I split the job into smaller parts. For example, I might process one month or one customer group at a time, and then merge the results.

This reduces the amount of data that needs to be shuffled in one go.

In summary:

If my Databricks job is slow because of shuffle operations, I take these steps:

- 1. Confirm it's a shuffle issue by checking Spark UI logs
- 2. Repartition or coalesce data to balance partitions
- 3. Use broadcast joins for small tables
- 4. Apply filters early to reduce data size
- 5. Use bucketing or Z-Ordering for frequent joins
- 6. Handle skewed keys with hints or salting
- 7. Avoid unnecessary wide transformations
- 8. Tune Spark shuffle configurations
- 9. Use Photon for better performance
- 10. Break the job into smaller, manageable pieces

These actions have helped me fix slow jobs due to shuffle in my past Databricks projects.

Scenario 44. What would your strategy be to migrate an existing on-prem job to Databricks?

If I'm asked this in an interview, I would explain that migrating an on-premises job to Databricks is not just about copying code from one system to another. It requires a well-planned approach to make sure the job runs efficiently in the cloud and takes advantage of Databricks' features.

Here's how I would explain my strategy in very simple and step-by-step terms:

Step 1: Understand the existing on-prem job

Before I migrate anything, I first try to fully understand how the current on-prem job works:

- What is the source of the data? (e.g., Oracle, SQL Server, flat files)
- What kind of processing is being done? (e.g., ETL, aggregations, joins, ML)
- What are the dependencies? (e.g., scripts, cron jobs, third-party tools)
- How often does the job run? (e.g., batch, near real-time)
- What technologies are being used? (e.g., Hadoop, Spark, Python scripts)

This helps me estimate the complexity of the migration and plan accordingly.

Step 2: Identify the equivalent services in Azure

Once I understand the current setup, I map the on-prem components to Azure services. For example:

- On-prem HDFS → Azure Data Lake Storage Gen2
- On-prem Spark/Hadoop → Azure Databricks
- On-prem scheduler → Azure Data Factory or Databricks Jobs
- On-prem relational DB → Azure SQL Database or Synapse

This mapping helps me design the new architecture.

Step 3: Set up the Databricks workspace and cluster

Then I set up the Databricks environment in Azure:

- Create a Databricks workspace
- Set up a compute cluster (with autoscaling and autosuspend enabled for cost control)
- Configure access to Azure Data Lake or any other data sources

I make sure the cluster has the required libraries and runtime (e.g., Spark 3.x with Delta Lake support).

Step 4: Migrate data storage first

If the source data is on-premise, I move it first to Azure Data Lake or Blob Storage.

- For one-time transfer, I can use tools like AzCopy or Azure Data Factory copy activity
- For ongoing sync, I can use Data Factory Integration Runtimes

This allows the data to be accessible from Databricks notebooks.

Step 5: Refactor the existing job code

If the on-prem job is written in Spark (Scala or PySpark), then the code can often be reused in Databricks with minor changes.

But if the job is written in something else (e.g., Hive scripts, shell scripts, SQL Server stored procedures), I refactor it into PySpark or Spark SQL code that works well in Databricks notebooks or jobs.

I also break the code into logical steps for better readability, for example:

- 1. Read source data
- 2. Clean/transform
- 3. Apply business logic
- 4. Write to destination (like Delta or Azure SQL)

Step 6: Use Delta Lake for output tables

If the job writes output files, I prefer to write them in Delta format instead of plain CSV or Parquet.

df.write.format("delta").mode("overwrite").save("/mnt/silver/sales_data")

Delta gives me benefits like ACID transactions, time travel, and schema enforcement, which are not available in normal file formats.

Step 7: Schedule the job in production

After testing, I schedule the job using one of the following:

- Databricks Jobs for simple scheduling (daily, hourly)
- Azure Data Factory for more complex orchestration (like multiple steps, dependencies, monitoring)

In Databricks Jobs, I can set retries, cluster configuration, email alerts, and parameters.

Step 8: Set up monitoring and alerting

I make sure the migrated job has proper monitoring in place:

- Set email alerts for failures in Databricks Job settings
- Log important information to a storage location or monitoring system
- Use tools like Azure Monitor or Log Analytics for better visibility

Step 9: Optimize the job for cloud

Once the job is running, I optimize it to reduce cost and improve performance:

- Use autoscaling clusters
- Use caching and partitioning wisely
- Use Photon-enabled clusters for faster performance
- Tune shuffle partitions or memory settings if needed

Step 10: Decommission the on-prem job

Finally, once the cloud job has been tested and validated, I work with the team to safely turn off the on-prem job, update documentation, and inform all stakeholders.

In summary:

My strategy to migrate an on-prem job to Databricks would be:

- 1. Understand the current job and its components
- 2. Map each component to its Azure equivalent
- 3. Set up Databricks and cloud storage
- 4. Move data from on-prem to Azure
- 5. Refactor and test the job code in Databricks
- 6. Write outputs using Delta Lake
- 7. Schedule the job using Databricks Jobs or ADF
- 8. Add monitoring and alerts
- 9. Optimize for performance and cost
- 10. Retire the old on-prem job after successful testing

This step-by-step approach helps ensure the migration is smooth, cost-effective, and improves over the original on-prem setup.

Scenario 45. What's your approach to processing petabytes of data daily in Databricks?

If I'm asked this question in an interview, I would answer by explaining that processing petabyte-scale data requires careful design to make sure the pipeline is scalable, efficient, and fault-tolerant. I would focus on distributed processing, storage format, partitioning, and optimizing performance at every level.

Here's how I would explain it step by step in very simple words:

Step 1: Use Delta Lake format for storage

When dealing with very large volumes of data, I always store the data in Delta Lake format instead of plain Parquet or CSV.

Delta Lake gives me several benefits:

- Faster reads and writes
- ACID transactions (safe for parallel writes)
- Schema evolution
- Time travel
- Easy partitioning and compaction

Example:

df.write.format("delta").mode("append").save("/mnt/bronze/events")

Step 2: Implement the Medallion Architecture

To organize the data pipeline clearly and avoid reprocessing everything, I use the bronze-silvergold model:

- Bronze layer: raw data ingested as-is from source
- Silver layer: cleaned and enriched data
- Gold layer: aggregated data for reporting or machine learning

This structure helps me process in stages, which improves performance and makes troubleshooting easier.

Step 3: Partition the data properly

When data is this huge, partitioning is very important.

I usually partition by a high-cardinality field like date, region, or event_type. This helps Spark skip data when reading.

Example:

```
df.write.partitionBy("event_date").format("delta").save("/mnt/silver/events")

If needed, I also use Z-Ordering to cluster data for faster querying:
```

OPTIMIZE events ZORDER BY (user_id)

Step 4: Use auto-scaling clusters with high concurrency

To handle petabytes of data, I use cluster autoscaling so that more worker nodes are added automatically based on workload.

- I use Job clusters for ETL tasks
- For high-throughput streaming, I use high-concurrency clusters
- I enable Photon (if available) for faster processing with native engine

This way, the system can grow and shrink based on the data volume.

Step 5: Process incrementally using streaming or incremental batch

I never load all data every day. Instead, I process only new or changed data daily using:

- Structured Streaming for real-time data
- Incremental batch using Delta Lake's MERGE, UPDATE, and CHANGE DATA FEED features

For example, I can write streaming logic like:

```
spark.readStream.format("cloudFiles") \
.option("cloudFiles.format", "json") \
.load("/mnt/raw/") \
.writeStream \
.format("delta") \
.option("checkpointLocation", "/mnt/checkpoints/") \
.start("/mnt/bronze/")
```

This helps reduce cost and speeds up processing.

Step 6: Use Data Skipping and Caching

When reading petabyte-scale data, performance can suffer. So I use:

- Data skipping through Delta partitioning and Z-order
- Caching for reused DataFrames (temporarily, only if needed)

Example:

df.cache()

df.count() # trigger the cache

But I use cache carefully to avoid memory issues.

Step 7: Monitor and tune Spark jobs

At this scale, performance tuning is critical. I regularly:

- Check Spark UI to find bottlenecks
- Tune spark.sql.shuffle.partitions based on data size
- Avoid wide transformations like large joins or groupBy
- Use broadcast joins for small lookup tables

Example:

spark.conf.set("spark.sql.shuffle.partitions", "1000")

Step 8: Handle failures and retries

For large pipelines, I make sure the jobs are:

- Idempotent: so they can be re-run safely
- Configured with retries (in Databricks Jobs)
- Using checkpoints for streaming jobs
- Logging errors to storage or monitoring tools

Step 9: Use ADF or Workflows for orchestration

I orchestrate the full pipeline using either:

- Azure Data Factory for scheduling, dependencies, and alerts
- Or Databricks Workflows to manage jobs in sequence with retries and notifications

Step 10: Monitor performance and cost

I regularly monitor:

- Cluster CPU/memory usage
- Shuffle size and skew
- Job run times
- Storage costs (optimize with VACUUM and OPTIMIZE)

This helps me keep the solution cost-effective while handling massive volumes.

In summary:

To process petabytes of data daily in Databricks, I follow these steps:

- 1. Use Delta Lake for scalable storage
- 2. Design the pipeline using Medallion (bronze-silver-gold)
- 3. Partition data properly and use Z-ordering
- 4. Use autoscaling and Photon-enabled clusters
- 5. Process data incrementally (streaming or batch)
- 6. Optimize with caching and data skipping
- 7. Tune Spark settings and avoid wide operations
- 8. Add retries, checkpoints, and error logging
- 9. Orchestrate using ADF or Databricks Workflows
- 10. Monitor performance, cost, and optimize continuously

This approach has worked for me in big data environments where data volume is massive and needs to be processed quickly and reliably.

Scenario 46. What is your approach to implementing row-level security in Databricks?

If I'm asked this question in an interview, I would explain that row-level security in Databricks means controlling which rows of data a user is allowed to see, based on their role, department, or identity. Since Databricks is built on Apache Spark, and Delta Lake doesn't have native row-level security like some data warehouses, I usually build custom logic to implement it using views, filters, or Unity Catalog.

Here is my approach, explained step by step in simple words:

Step 1: Understand the business rule for row-level security

First, I ask questions like:

- Should each department see only its own data?
- Are there user roles like Admin, Manager, Analyst?
- Is the user information (like role or region) stored somewhere?

This helps me know how to define the filtering logic for example, region-based access, department-based access, or role-based access.

Step 2: Store user-role mapping

I create or use an existing mapping of users and their allowed access. This can be stored in a table like:

username, department, region

alice@company.com, HR, US

bob@company.com, Finance, EU

I save this as a Delta table, say user_access_control.

Step 3: Use current_user() function (Unity Catalog) or manually pass user info

If Unity Catalog is enabled in Databricks, I can get the current user using:

SELECT current_user()

Otherwise, in older setups, I may have to pass the username into the notebook or workflow as a parameter.

Step 4: Join user info with the main data table

Now, I use a SQL view or a dynamic filter in Spark to return only the rows the user is allowed to see.

For example, let's say we have a sales_data table with columns: region, sales_amount, product I create a view like this:

CREATE OR REPLACE VIEW secure_sales_data AS

SELECT s.*

FROM sales_data s

JOIN user_access_control u

ON s.region = u.region

WHERE u.username = current_user()

Now, when a user runs a query on secure_sales_data, they will only see data for their own region.

If Unity Catalog is not available, I pass the username as a parameter and filter it in PySpark: user = dbutils.notebook.entry_point.getDbutils().notebook().getContext().userName().get() user_df = spark.table("user_access_control").filter(F.col("username") == user) region = user_df.select("region").first()["region"] secure_df = spark.table("sales_data").filter(F.col("region") == region)

Step 5: Hide the base table and expose only the secure view

To make sure users cannot access full raw data, I do the following:

- Restrict access to the base table (sales_data)
- Grant SELECT permission only on the secure view (secure_sales_data)
- Use Unity Catalog or workspace permissions to enforce this

Step 6: Optional Add logic for Admin roles

Sometimes admins need access to all data. So I enhance the logic like this:

CREATE OR REPLACE VIEW secure_sales_data AS

SELECTs.*

FROM sales_data s

JOIN user_access_control u

ON (s.region = u.region OR u.role = 'Admin')

WHERE u.username = current_user()

So admins can see all rows, while regular users see only filtered data.

Step 7: Test with different users

Before going live, I test the view by switching user contexts (if possible) or simulating different usernames. I check if each user can only see the data they are supposed to.

Step 8: Keep the access control table up to date

The user_access_control table must be maintained regularly. I usually automate this by syncing it with Active Directory or Azure AD using APIs or scheduled updates.

Step 9: Use Unity Catalog (if available) for built-in row-level security

If my Databricks workspace uses **Unity Catalog**, I can now define **row-level access policies** using SQL commands.

For example:

CREATE ROW FILTER region_filter AS (region STRING) -> region = current_user_region()

Then I apply it:

ALTER TABLE sales_data SET ROW FILTER region_filter ON (region)

This feature makes row-level security easier and more manageable without needing custom views.

In summary:

To implement row-level security in Databricks, I follow this approach:

- 1. Understand how rows need to be restricted based on user identity
- 2. Create or use a user-role mapping table (user access control)
- 3. Use current_user() (with Unity Catalog) or manually detect the user
- 4. Filter the main data based on user access using views or dynamic filters
- 5. Restrict access to raw tables and expose only secure views
- 6. Add admin override logic if needed
- 7. Test the logic thoroughly
- 8. Keep the access table updated regularly
- 9. Use Unity Catalog's native row-level security if available

This strategy helps me ensure that users only see the data they are allowed to, keeping things secure and compliant with data privacy policies.

Scenario 47. You need to collaborate on notebooks with version control how would you set that up?

If I'm asked this question in an interview, I would explain that Databricks now supports Gitbased version control using Databricks Repos, which allows real-time collaboration while keeping notebooks synced with Git providers like GitHub, Azure DevOps, or GitLab.

Here's how I would set this up step by step, explained in simple words:

Step 1: Choose a Git provider

First, I confirm which Git system my team is using for example:

- GitHub
- Azure DevOps
- GitLab
- Bitbucket

This will be the central place where all code changes are stored and tracked.

Step 2: Set up Git credentials in Databricks

Before I can connect to the Git repo from Databricks, I need to link my Git account to my Databricks user.

To do this:

- 1. Go to User Settings in Databricks.
- 2. Click on the Git Integration tab.
- 3. Choose the Git provider (e.g., GitHub).
- 4. Paste a personal access token (PAT) from my Git provider.

This allows Databricks to push and pull notebooks from Git.

Step 3: Create or clone a repo using Databricks Repos

To start using Git version control in notebooks, I go to the Repos section in the left sidebar and click "Add Repo".

Now I have two options:

- 1. Clone an existing Git repo:
 - I paste the Git URL of the existing project.
 - Example: https://github.com/my-team/databricks-etl-pipeline.git

2. Create a new repo:

• I can create notebooks first in Databricks and later push them to a new Git repository from there.

Databricks automatically pulls the latest code from Git and keeps track of changes.

Step 4: Work on the notebook and commit changes

Now, when I open a notebook inside the repo, I can:

- Make changes to code or markdown
- Click "Git" in the top-right menu to:
 - View changes
 - Commit the code
 - Push to the remote Git branch

This is just like working with Git in VS Code or other editors, but inside Databricks.

Step 5: Collaborate using Git branches

To work with teammates without overwriting each other's work, we use branches.

- I create a feature branch: feature/add-ingestion-logic
- My teammate creates another: feature/add-validation-step

We each commit and push code to our branches, and later create pull requests in GitHub or Azure DevOps to merge the changes into main.

This helps keep work organized and prevents conflicts.

Step 6: Use Git best practices

To keep everything clean and collaborative, I follow good practices:

- Always pull the latest changes before editing
- Use meaningful commit messages
- Create pull requests with descriptions
- Review others' PRs before merging
- Use .gitignore to exclude unnecessary files (e.g., checkpoints)

Step 7: Use %run or modularization for reusability

To make notebooks more reusable and team-friendly:

- I break large workflows into multiple smaller notebooks (modular design)
- Use %run ./utils/cleaning to reuse code
- Or package logic into Python files and import functions

This helps everyone work on different parts independently.

Step 8: Automate testing and deployment (optional)

Once Git is set up, I can integrate it with:

- CI/CD tools like GitHub Actions or Azure Pipelines
- Automated tests using pytest or notebooks with test scripts
- Deployment pipelines to run or deploy notebooks after PRs

This makes the collaboration full end-to-end from development to deployment.

In summary:

To collaborate on notebooks with version control in Databricks, I do the following:

- 1. Set up Git credentials in Databricks user settings
- 2. Clone or create a Git repo using Databricks Repos
- 3. Make changes and commit them through the Databricks Git UI
- 4. Use branches to work safely with teammates
- 5. Create pull requests for review and merge
- 6. Follow Git best practices
- 7. Modularize notebooks and use %run for clean code
- 8. (Optional) Integrate with CI/CD for automated testing and deployment

This workflow helps me collaborate with my team just like we would in any software engineering project, while keeping the work versioned and safe.

Scenario 48. You need to implement a data governance strategy in Azure Databricks what steps would you take?

If I get this question in an interview, I would explain that data governance is about protecting, managing, and monitoring data so that it's used in the right way by the right people. In Databricks, I usually use tools like Unity Catalog, access controls, data lineage, and auditing to set up a full governance strategy.

Here's how I would do it step by step, using simple language:

Step 1: Understand data governance requirements

Before doing anything technical, I would talk to data owners, security teams, and compliance teams to understand:

- What data is sensitive or regulated (like PII or financial data)?
- Who should have access to what data?
- What audit or logging is needed?
- What compliance rules we need to follow (e.g., GDPR, HIPAA)?

This helps define the scope of the governance plan.

Step 2: Organize data using Unity Catalog

I would use Unity Catalog to create a centralized way to manage data access and structure. It gives me a way to define:

- Metastores (one per region or environment)
- Catalogs (logical groupings, like finance, hr, marketing)
- Schemas (inside catalogs, like folders e.g., finance.raw, finance.cleaned)
- Tables and views (actual data)

For example:

main_catalog

finance

--- raw

--- silver

gold

This structure keeps everything clean and easier to secure.

Step 3: Apply fine-grained access control

Once data is organized, I define **access rules** using Unity Catalog. This helps make sure only the right people can access certain data.

I do this by:

- Granting permissions like SELECT, INSERT, MODIFY on specific tables
- Assigning access based on groups (e.g., finance_team, hr_team)
- Using row-level and column-level security where needed

Example SQL command to restrict access:

GRANT SELECT ON TABLE finance.gold.revenue TO `finance_team`;

Or to block access to a PII column:

CREATE MASKING POLICY mask_ssn AS (val STRING) -> return 'XXX-XX-XXXX';

ALTER TABLE hr.gold.employees ALTER COLUMN ssn SET MASKING POLICY mask_ssn;

Step 4: Tag and classify sensitive data

To make data easier to govern, I add tags to columns that contain sensitive information like:

- email_address
- credit_card_number
- social_security_number

Databricks allows me to use column tags to label data. This is helpful for:

- Auditing
- Access control
- Data classification

Step 5: Enable auditing and monitoring

To keep track of who accessed or modified data, I enable logging through:

- Unity Catalog audit logs (sent to Azure Monitor or Log Analytics)
- Cluster and job logs
- Access history for tables and views

This helps in compliance reporting and in investigating issues.

Step 6: Implement data lineage

I use Unity Catalog's built-in data lineage feature to automatically track:

- Where the data came from (input source)
- What transformations were applied
- Where the data went (outputs, dashboards, reports)

This is very helpful for debugging, auditing, and understanding data flow.

Example: I can see that a gold table was built from a silver table, which was built from raw logs.

Step 7: Secure notebooks and jobs

To ensure governance at the code level:

- I use Git version control (via Databricks Repos)
- I manage access to notebooks using workspace permissions
- I restrict who can run or edit jobs and workflows
- I use parameters instead of hardcoded values (especially for secrets)

Step 8: Protect secrets and credentials

For any passwords, API keys, or connection strings, I never store them in notebooks. Instead, I use:

- Databricks Secrets
- Store them in a secret scope (backed by Azure Key Vault)
- Access them securely inside code like this:

dbutils.secrets.get(scope="prod_scope", key="sql_password")

Step 9: Set up data retention and lifecycle policies

I define policies for:

- How long we keep raw, intermediate, or output data
- Automatic cleanup of old files or logs
- Archiving and backup strategies

This helps with storage cost control and compliance.

Step 10: Train users and enforce policies

Finally, governance only works if everyone follows it. So I:

- Educate users about policies and best practices
- Monitor access patterns
- Set up alerts for suspicious or unauthorized access
- Review permissions regularly

In summary:

To implement a data governance strategy in Databricks, I would:

- 1. Understand the business and compliance requirements
- 2. Use Unity Catalog to structure and organize data
- 3. Set fine-grained access controls (row, column, table-level)
- 4. Tag sensitive data for classification
- 5. Enable logging and audit trails
- 6. Use data lineage to track data movement
- 7. Secure notebooks, jobs, and workflows
- 8. Use secret scopes for managing credentials securely
- 9. Define data retention and lifecycle rules
- 10. Train users and monitor usage regularly

This full approach helps me make sure that data in Databricks is secure, trusted, and used in a responsible way.

Scenario 49. You need to optimize a Databricks job that processes petabytes of data what would you do?

If I am asked this in an interview, I would explain that working with petabytes of data requires special attention to performance, memory, and cost. I would take a structured approach to optimize the job in both Spark logic and Databricks environment. Here's how I would explain it step by step in simple terms:

Step 1: Understand the job and data flow

Before changing anything, I first try to understand:

- What is the job doing? (e.g., ingestion, transformation, aggregation)
- What are the input and output data formats?
- Where is the data coming from (e.g., ADLS, S3)?
- Is the job CPU-bound, memory-bound, or I/O-bound?

I review the job run history and check how long each stage is taking, especially shuffle-heavy ones.

Step 2: Use efficient file formats and compression

I make sure that the job is reading and writing in optimized formats like:

- Delta Lake (best for performance and ACID support)
- Or Parquet if Delta is not used

I also make sure compression is enabled (like Snappy), which reduces I/O and speeds up reading and writing.

Example:

df.write.format("delta").mode("overwrite").save("/mnt/datalake/output/")

Step 3: Partition and Z-Order the data

To speed up reads and writes, I use partitioning:

- For example, partitioning by date or region if queries filter on those
- I make sure not to over-partition (avoid tiny files problem)

Also, if I'm using Delta Lake, I apply Z-Ordering on frequently filtered columns:

OPTIMIZE my_table ZORDER BY (customer_id, product_id);

This makes scanning large datasets much faster.

Step 4: Reduce shuffle and expensive operations

Shuffle is often the main reason why jobs are slow or crash at petabyte scale. I do these:

- Avoid groupBy if reduceByKey or mapGroups can be used
- Avoid wide joins unless needed
- Use broadcast joins when joining with small tables:

from pyspark.sql.functions import broadcast

df.join(broadcast(small_df), "id")

• Coalesce or repartition only when needed and strategically

Step 5: Use cluster wisely

I check if the cluster size is appropriate:

- Use a cluster with enough memory and CPUs to handle large partitions
- Turn on autoscaling so it adjusts based on the workload
- Choose Photon runtime if possible (faster and optimized for big data)

Also, I use job clusters for each run to avoid leftover states from previous jobs.

Step 6: Use caching carefully

If the same intermediate data is used multiple times, I cache it:

df.cache()

But I also monitor memory usage caching huge datasets can cause memory errors, so I use it only when it saves time and fits in memory.

Step 7: Use adaptive query execution (AQE)

Databricks automatically enables AQE in newer runtimes, but I always make sure it's on:

- It automatically optimizes joins and partition sizes during runtime
- Helps when data skew is present

No code needed it works under the hood in Spark 3+ runtimes.

Step 8: Monitor performance using Spark UI

To find bottlenecks, I check the Spark UI:

- Look at Stage DAG to see if any stages are slow
- Check task skew if one task takes much longer, I try to fix the skewed data
- Look at Shuffle Read/Write and executor memory

This helps me target the exact spot where optimization is needed.

Step 9: Tune configurations if needed

Sometimes I tune the Spark job with config settings like:

spark.conf.set("spark.sql.shuffle.partitions", "500")

spark.conf.set("spark.sql.adaptive.enabled", "true")

If I notice too many partitions or too much shuffle, I adjust the values accordingly.

Step 10: Break large jobs into stages (optional)

If the job is too large and does multiple things like ingestion, cleaning, joining, aggregation, and writing I try to split it into multiple smaller jobs.

For example:

- First job: ingest and clean raw data
- Second job: perform joins and enrich the data
- Third job: run aggregations and write to output

This modular approach makes the pipeline easier to debug and scale.

Summary

To optimize a Databricks job that processes petabytes of data, I would:

- 1. Understand the job's logic and where time is spent
- 2. Use efficient file formats like Delta and compression
- 3. Partition and Z-Order the data for faster access
- 4. Reduce shuffle operations and use broadcast joins
- 5. Use autoscaling clusters with Photon runtime
- 6. Cache only when useful and safe
- 7. Use adaptive query execution (AQE)
- 8. Monitor the Spark UI to find bottlenecks
- 9. Tune Spark configurations if needed
- 10. Split big jobs into multiple modular stages

This way, I ensure the job runs faster, more reliably, and uses resources efficiently even at petabyte scale.

53. You need to migrate a Hadoop job to Databricks how do you plan it?

If I'm asked this in an interview, I would say that migrating a Hadoop job to Databricks is not just a "lift and shift" process. I would break it down into a few clear steps: understanding the existing job, re-engineering it for Spark, testing it, and finally deploying and optimizing it in Databricks. I'll try to explain it in very simple words below.

Step 1: Understand the existing Hadoop job

First, I take time to understand what the current Hadoop job is doing:

- What is the input and output data format (e.g., CSV, JSON, ORC)?
- Is it written in MapReduce, Hive, Pig, or Spark?
- What are the transformations, joins, or aggregations being performed?
- What are the dependencies (scripts, JARs, external tools)?
- How often is the job scheduled?

This helps me know what needs to be rebuilt or migrated.

Step 2: Set up the Databricks environment

I create the target environment in Databricks:

- · Set up the Databricks workspace
- Mount Azure Data Lake (ADLS) or any storage where data will reside
- Create clusters with the right compute and autoscaling settings
- Set up access to external systems (like Azure SQL, Synapse, etc.)
- Configure Unity Catalog if needed for managing data and permissions

Step 3: Choose the right compute engine (PySpark or SQL)

If the old job was in MapReduce or Hive, I usually migrate it to PySpark or Spark SQL inside Databricks notebooks.

For example:

- Hive → Spark SQL
- MapReduce → PySpark
- Pig → PySpark or SQL

If the original job is already in Spark (like running on a Hadoop YARN cluster), I reuse the PySpark logic and adapt it to Databricks style.

Step 4: Migrate the code

This is the main part. I start rewriting the job as a notebook or a job script in Databricks.

- If it's a Hive job: I copy the SQL queries and run them using Spark SQL
- If it's a MapReduce job: I rewrite it in PySpark using DataFrame API

Example: a MapReduce word count job becomes:

```
text_df = spark.read.text("/mnt/input_data/book.txt")
words_df = text_df.selectExpr("explode(split(value, ' ')) as word")
```

word_counts = words_df.groupBy("word").count()

word_counts.write.mode("overwrite").parquet("/mnt/output_data/wordcount")

This is much simpler and more efficient than MapReduce.

Step 5: Migrate data to cloud storage

If the Hadoop job used HDFS, I move that data to Azure Data Lake Storage (ADLS) or Blob Storage, because Databricks works with cloud storage.

- I use tools like DistCp or AzCopy to move large files
- For small migrations, I can download from HDFS and upload manually to ADLS

Step 6: Replace any custom libraries or scripts

Hadoop jobs sometimes use external scripts or JARs. In Databricks:

- I install Python libraries using %pip install or attach JARs to the cluster
- I rewrite shell scripts or bash jobs as Databricks workflows

Step 7: Test the migrated job on a small dataset

Before running the full job on terabytes or petabytes of data, I test it on a small sample to make sure:

- Output matches the original job
- Performance is acceptable
- No errors during runtime

I use sample filters like limit 1000 or where date = '2023-01-01' to reduce volume during testing.

Step 8: Automate with jobs and workflows

Once tested, I schedule the job using Databricks Jobs:

- Set the cluster, notebook, parameters, retry logic, and email alerts
- Use task workflows if there are multiple steps (e.g., extract → transform → load)

Step 9: Monitor and optimize

After running in production, I monitor performance using:

- Spark UI
- Ganglia metrics
- Databricks job run logs

If needed, I optimize with:

- Caching
- Partitioning
- Delta Lake
- Better cluster sizing

Step 10: Retire old Hadoop infra

Once the Databricks job is live and stable, I decommission the old Hadoop job and release the infrastructure to save cost.

Summary

To migrate a Hadoop job to Databricks, I follow this step-by-step plan:

- 1. Understand what the Hadoop job does
- 2. Set up the Databricks environment
- 3. Choose the right engine (PySpark or SQL)
- 4. Rewrite the job logic using Spark in Databricks
- 5. Migrate the data from HDFS to ADLS
- 6. Replace any custom scripts or JARs
- 7. Test the new job on sample data
- 8. Schedule and automate the job
- 9. Monitor and optimize performance
- 10. Retire the old Hadoop job after validation

This way, I make sure the migration is clean, efficient, and future-ready.

Scenario 51. You need to manage secrets in Databricks securely how would you approach it?

If I'm asked this in an interview, I would say that securely managing secrets in Databricks is very important when working with databases, APIs, or storage systems. Secrets include things like passwords, tokens, and connection strings. I always make sure that these values are not hardcoded in notebooks and are securely stored and accessed.

Here's how I would handle it step by step, in very simple words.

Step 1: Decide where to store secrets

Databricks provides two main options for secret management:

- 1. Databricks-backed secret scope Secrets are stored within Databricks
- 2. Azure Key Vault-backed scope Secrets are stored in Azure Key Vault and Databricks just accesses them

I usually prefer Azure Key Vault if we want centralized secret management across teams and services, especially in enterprise environments.

Step 2: Create a secret scope

If using Databricks-backed secret scope:

1. I use the Databricks CLI to create the scope:

databricks secrets create-scope --scope my-secrets

2. Then I add secrets to it:

databricks secrets put --scope my-secrets --key db-password

If using Azure Key Vault, I link it like this:

- Go to Azure portal and create a Key Vault (if not already)
- Add secrets in it (e.g., connection string, API key)
- Then in Databricks, create a Key Vault-backed scope:

databricks secrets create-scope --scope akv-secrets --scope-backend-type AZURE_KEYVAULT --resource-id <vault-resource-id> --dns-name <vault-dns>

Now Databricks can read from the Key Vault.

Step 3: Access the secret inside notebooks

```
Once the secrets are added to a scope, I use them like this in the notebook:

db_password = dbutils.secrets.get(scope="my-secrets", key="db-password")

Now I can use db_password safely in my JDBC connection or API call.

Example:

jdbc_url = "jdbc:sqlserver://mydb.database.windows.net;databaseName=mydb"

connection_properties = {

"user": "myuser",

"password": db_password

}

df = spark.read.jdbc(url=jdbc_url, table="sales", properties=connection_properties)
```

Step 4: Set permissions on secret scopes

For extra safety, I control who can access which scope.

This way, the password is never shown in plain text.

- In Databricks workspace, go to "Admin Console" → "Secrets"
- Set permissions like:
 - Only certain users or groups can read secrets
 - Some users can write but not read

This prevents unauthorized access.

Step 5: Never log or print secrets

A very important practice I follow is to never print secrets in notebooks or logs. I make sure I don't do this:

```
print(db_password) # X Bad practice
```

Instead, I use the secret only where needed (like connections), and avoid exposing it.

Step 6: Rotate secrets regularly

In production environments, I also recommend rotating secrets like database passwords or tokens periodically (e.g., every 90 days), and updating the values in Key Vault or secret scope.

Databricks will automatically use the updated value if it reads from Key Vault.

Step 7: Use Unity Catalog (optional)

If using Unity Catalog, I can manage access to external resources securely using External Locations, Data Access Connectors, and Credential Passthrough. In some cases, I can avoid using raw secrets at all by using Azure identities.

Summary

To manage secrets securely in Databricks, I follow these steps:

- 1. Choose between Databricks-backed or Azure Key Vault-backed secret scopes
- 2. Create secret scopes and add secrets (like passwords or tokens)
- 3. Access secrets securely using dbutils.secrets.get()
- 4. Set permissions on secret scopes to control access
- 5. Never print secrets in logs or notebooks
- 6. Rotate secrets regularly for security
- 7. Use identity-based access (like Unity Catalog or managed identity) if possible

This approach helps me keep sensitive credentials safe and avoid security risks in Databricks projects.

Scenario 52. You are given a large dataset with evolving schema how do you handle it in Delta Lake?

If this question comes up in an interview, I would explain that schema evolution is something that happens often in real-world data pipelines. For example, the source system might start sending new columns, change column order, or change data types. Delta Lake handles this much better than traditional formats because it supports schema evolution and schema enforcement.

Here's how I would handle evolving schema in a large dataset using Delta Lake, step by step, in simple language.

Step 1: Understand what kind of schema changes are happening

Before doing anything, I try to understand:

- Are new columns being added over time?
- · Are any columns being removed or renamed?
- Are data types changing (like string to int)?
- Are fields nested (like JSON or struct types)?

This helps me plan how to handle it.

Step 2: Use Delta Lake for writing data

I always use Delta Lake instead of CSV or Parquet because Delta supports schema evolution directly.

So first, I write the data like this:

df.write.format("delta").mode("append").save("/mnt/datalake/raw/events")

But if the schema is changing (like new columns are coming), this will fail unless I explicitly enable schema evolution.

Step 3: Enable schema evolution during write

If new columns are expected in the future, I use the mergeSchema option.

```
df.write.format("delta") \
    .option("mergeSchema", "true") \
    .mode("append") \
    .save("/mnt/datalake/raw/events")
```

This tells Delta Lake to merge the incoming schema with the existing one, without failing.

Delta automatically updates the schema in the Delta log.

Step 4: Handle schema changes in merge (upsert) operations

If I'm doing merge operations (like updating existing records), and the schema might change, I also need to set spark.databricks.delta.schema.autoMerge.enabled to true.

spark.conf.set("spark.databricks.delta.schema.autoMerge.enabled", "true")

from delta.tables import DeltaTable

target = DeltaTable.forPath(spark, "/mnt/datalake/raw/events")

target.alias("target").merge(

df.alias("source"),

"target.id = source.id"

).whenMatchedUpdateAll().whenNotMatchedInsertAll().execute()

Without this, Delta will throw a schema mismatch error.

Step 5: Use schema inference carefully for semi-structured data

If I'm reading JSON or CSV data where the schema may change, I let Spark infer schema but save it afterward.

df = spark.read.option("inferSchema", True).json("/mnt/raw_data/events/")

df.write.format("delta").mode("append").option("mergeSchema", "true").save("/mnt/datalake/raw/events")

To avoid surprises in production, I may also use a defined schema and update it when needed.

Step 6: Track and audit schema history

Delta Lake keeps a history of schema changes. I can check this using:

spark.read.format("delta").load("/mnt/datalake/raw/events").history()

This helps me troubleshoot when something breaks due to unexpected schema drift.

Step 7: Use table constraints for critical fields

If I want to make sure some important fields always have valid values (like id or date), I add constraints to the Delta table.

ALTER TABLE events ADD CONSTRAINT valid_id CHECK (id IS NOT NULL);

This avoids bad data even when schema evolves.

Step 8: Communicate schema changes with downstream teams

Since schema changes affect everyone, I make sure to:

- Document the schema changes
- Notify downstream teams (BI, ML) when new columns are added
- Keep versioned schemas using a schema registry (if needed)

Summary

To handle evolving schema in Delta Lake, I follow these steps:

- 1. First understand what kind of schema changes are happening
- 2. Use Delta format to store the data
- 3. Enable mergeSchema during write to accept new columns
- 4. Enable schema.autoMerge for merge operations (upserts)
- 5. Use schema inference carefully with JSON or CSV
- 6. Use DESCRIBE HISTORY to track schema changes over time
- 7. Apply constraints to enforce important fields
- 8. Communicate changes with other teams and document them

This way, my pipeline becomes flexible and can handle schema changes without breaking or requiring manual fixes.

Scenario 53. You want to perform ETL on data from ADLS to Azure SQL how would you build it in Databricks?

If I'm asked this in an interview, I would say that moving data from Azure Data Lake Storage (ADLS) to Azure SQL Database using Databricks is a very common ETL use case. I would explain it step-by-step in simple terms, from connecting to the source, transforming the data, and finally loading it into Azure SQL.

Step 1: Connect Databricks to ADLS (source)

To read data from ADLS, I first mount or connect to the ADLS container. If I'm using Azure Data Lake Gen2, I usually configure it like this:

- 1. Get the Storage Account Name and Access Key or use Azure Service Principal
- 2. Set the credentials using Spark config (if not mounting):

```
spark.conf.set(
   "fs.azure.account.key.<storage_account>.dfs.core.windows.net",
   "<your_access_key>"
)
```

3. Then I read the data (for example, from a CSV or Parquet file):

df = spark.read.format("parquet").load("abfss://mycontainer@storageaccount.dfs.core.windows.net/sales/")

Step 2: Transform the data in Databricks

Once the data is loaded into a DataFrame (df), I perform the required transformations. This can include:

- Renaming columns
- Filtering records
- Aggregating values
- Handling nulls or duplicates
- Changing data types

Example:

Step 3: Prepare for writing to Azure SQL Database (target)

To load into Azure SQL, I need:

- JDBC URL
- SQL server name, database name, username, password
- The target table (created in advance or created from Databricks)

I also use a secret scope or key vault to keep the password safe (not hardcoded).

jdbc_url = "jdbc:sqlserver://<server-name>.database.windows.net:1433;database=<databasename>"

```
connection_properties = {
   "user": "my_user",
   "password": dbutils.secrets.get(scope="sql-secrets", key="sql-password"),
   "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
}
```

Step 4: Write the transformed data to Azure SQL

Now I write the transformed_df to Azure SQL using the .write.jdbc() method.

```
transformed_df.write \
```

.mode("overwrite") \ # or "append" depending on the use case

.jdbc(url=jdbc_url, table="dbo.sales_cleaned", properties=connection_properties)

If the table does not exist, I can let Databricks create it automatically during the first write.

Step 5: Schedule the ETL job

Once everything works, I create a Databricks job to automate this process:

- 1. Go to the Jobs tab in Databricks
- 2. Create a job and select the notebook
- 3. Set the schedule (e.g., daily at midnight)
- 4. Add retries and alerts if needed

Step 6: Monitor and log the pipeline

To make the ETL job production-ready, I also include:

- Logging (how many records were read/written)
- Error handling using try-except
- Job run history and alerts in case of failures

Example:

print("Records written:", transformed_df.count())

Optional: Use Delta Tables for intermediate stages

If the transformation is complex or has multiple steps, I might store the intermediate results in a Delta table:

transformed_df.write.format("delta").mode("overwrite").save("/mnt/datalake/temp/sales_clea ned")

Then I can further process or validate before final load.

Summary

To build an ETL from ADLS to Azure SQL using Databricks, my steps are:

- 1. Connect to ADLS and read raw data
- 2. Clean and transform the data using PySpark
- 3. Prepare the JDBC connection to Azure SQL
- 4. Write the transformed data into the target SQL table
- 5. Schedule the ETL job using Databricks Jobs
- 6. Add logging, error handling, and monitoring
- 7. (Optional) Use Delta Lake for intermediate steps

This makes the pipeline reliable, scalable, and production-ready.

Scenario 54. You need to troubleshoot a shuffle-heavy job that's failing how do you approach it?

If I get this question in an interview, I would say that shuffle operations are one of the most expensive parts of a Spark job. They involve moving data across different nodes in the cluster, which can easily cause issues like out of memory, slow performance, or job failure. So when I troubleshoot such a job in Databricks, I follow a step-by-step process to identify and fix the problem.

Step 1: Understand what causes shuffles

First, I remind myself of the common actions that cause shuffle:

- Joins (especially between large tables)
- GroupBy or Aggregations
- Distinct
- OrderBy or Sort
- Repartition

So I check if any of these are being used in the code and whether they're acting on large DataFrames.

Step 2: Review the Spark UI in Databricks

I go to the Spark UI under the failed job's details. This is where I find useful information like:

- Which stage took the longest
- Which tasks failed or retried
- How much data was shuffled
- Skewed partition sizes
- Executor memory usage

I look for stages with high shuffle read/write and any signs of task skew (like one task taking much longer than the others).

Step 3: Identify data skew

Data skew happens when one partition has a lot more data than others. This causes only one executor to do most of the work while others stay idle.

To check for skew, I run:

df.groupBy("join_column").count().orderBy("count", ascending=False).show(10)

If a single key appears millions of times, I know the data is skewed.

Step 4: Optimize joins (if applicable)

If the issue is due to a join, I check:

Can I use a broadcast join? (Only if one table is small)

from pyspark.sql.functions import broadcast

result = big_df.join(broadcast(small_df), "id")

• Can I filter or reduce data before the join?

Example: Instead of joining full tables, I do this:

```
filtered_df = big_df.filter("status = 'active'")
```

result = filtered_df.join(broadcast(small_df), "id")

Step 5: Repartition the data

If the partitions are not balanced, I try to repartition the data by a more even column.

```
df = df.repartition(100, "customer_id")
```

I avoid df.repartition() without a column, as it can trigger a full shuffle. If I just want fewer partitions, I use coalesce() instead:

df = df.coalesce(10)

Step 6: Avoid wide transformations where possible

I try to reduce the number of wide transformations like:

- groupBy
- orderBy
- distinct

For example, if ordering is not essential, I skip orderBy.

Step 7: Increase shuffle partitions if needed

If the data volume is huge, Spark may create too few shuffle partitions, which leads to each task handling too much data. I increase this setting:

spark.conf.set("spark.sql.shuffle.partitions", "800")

Default is usually 200, but for large jobs, 800 or even 1000 can help.

Step 8: Check executor memory and cluster size

If executors are running out of memory during shuffles, I do one of the following:

- Increase executor memory in the cluster settings
- Use a larger cluster with autoscaling
- Enable adaptive query execution if not already enabled

spark.conf.set("spark.sql.adaptive.enabled", "true")

Adaptive query execution helps Spark optimize the number of shuffle partitions dynamically at runtime.

Step 9: Cache intermediate results (if reused)

If an intermediate DataFrame is used multiple times, I cache it to avoid recomputation and repeated shuffles.

df.cache()

Step 10: Rerun and monitor

After applying changes, I rerun the job and monitor:

- Shuffle size
- Task distribution
- Execution time

I repeat this until the shuffle-heavy stage is balanced and the job completes successfully.

Summary

To troubleshoot a shuffle-heavy job that's failing, I follow these steps:

- 1. Use the Spark UI to identify slow/failing stages
- 2. Look for signs of data skew
- 3. Try using broadcast joins for small tables
- 4. Filter data early before joining
- 5. Repartition data to avoid skew
- 6. Avoid wide transformations like orderBy and distinct
- 7. Increase spark.sql.shuffle.partitions if needed
- 8. Check executor memory and enable AQE
- 9. Use caching when reusing data
- 10. Rerun and validate performance

This approach helps me fix shuffle-related issues and improve both the reliability and speed of the Spark job.

Scenario 55. Your team needs to collaborate on a shared notebook with version control how do you enable it?

If this question is asked in an interview, I would explain that in Azure Databricks, the best way for a team to collaborate on notebooks with proper version control is by using Databricks Repos, which integrates with Git providers like Azure DevOps, GitHub, and Bitbucket. This helps us manage code changes, collaborate without overwriting each other's work, and maintain a history of changes.

Here's how I would approach it step by step:

Step 1: Set up a Git repository

First, I make sure that the team has access to a Git repository. It could be in:

- Azure DevOps Repos
- GitHub
- GitLab
- Bitbucket

The repo will hold our notebooks and other code files.

Step 2: Configure Git integration in Databricks

- 1. Go to the User Settings in Databricks (top-right corner)
- 2. Click on Git Integration
- 3. Choose the Git provider (e.g., GitHub, Azure DevOps)
- 4. Generate a personal access token (PAT) from the Git provider
- 5. Paste the token in Databricks to connect your Git account

This is a one-time setup per user.

Step 3: Create or clone a repo inside Databricks

- 1. Go to the Repos tab on the left sidebar in Databricks
- 2. Click Add Repo
- 3. If the repo already exists, select Clone Repo and paste the Git URL Example for GitHub:

https://github.com/myorg/data-pipeline-notebooks.git

4. If we want to create a new repo, we start with a new notebook and then push it to Git from within Databricks.

Step 4: Collaborate using Git branches

- Every developer creates or checks out their own branch inside the repo folder.
- They make changes to notebooks or code files in their branch.
- Once changes are tested and reviewed, they commit and push the updates back to Git using the built-in UI or CLI.

Inside the notebook, there's a Git panel where I can:

- View diffs (changes made)
- Commit messages
- Push changes
- Pull latest updates
- Create or switch branches

This helps avoid conflicts and makes the workflow smoother.

Step 5: Use Pull Requests for merging code

Once a team member completes their work:

- 1. They push the changes to Git (from Databricks or a local Git tool)
- 2. They create a Pull Request (PR) in the Git platform
- 3. The team reviews and approves the PR
- 4. Once approved, the code is merged into the main branch

This process ensures code quality and team collaboration.

Step 6: Schedule jobs from the versioned notebooks (optional)

If we want to schedule a job using a notebook in a Git repo:

- 1. We make sure the notebook is synced with Git
- 2. In the Jobs UI, we select the notebook from the repo path
- 3. This way, the job always runs from the version-controlled code

Step 7: Set permissions and access control

- Use Workspace Access Control to manage who can edit or view notebooks
- Set branch protection rules in Git (e.g., require PR approval)
- Use GitHub Actions or Azure Pipelines to automate deployment if needed

Summary

To enable collaboration and version control on Databricks notebooks, I would:

- 1. Set up a Git repo for the team
- 2. Configure Git integration in each user's Databricks account
- 3. Clone the repo into Databricks Repos
- 4. Use branches to work on features individually
- 5. Commit and push changes through Databricks Git UI
- 6. Use pull requests to review and merge code
- 7. Set proper permissions and optionally automate deployments

This approach makes teamwork on notebooks smooth, traceable, and reliable, just like working on any software project.

Scenario 56. Your Spark job is failing due to skew how would you resolve it in Databricks?

If this question is asked in an interview, I would explain that data skew happens when one or more partitions have a lot more data than the others. This causes a few tasks to take a very long time, while others finish quickly. In worst cases, it can even cause executor memory issues and job failures.

Here's exactly how I would handle skew in Databricks, step by step:

Step 1: Confirm that skew is the issue

I first check the Spark UI from the failed job run:

- I go to the "Stages" tab and look for a stage that is taking too long.
- Then I look at the task duration and input size of each task.
- If one task is taking way more time or has much more input than the rest, then it's clear that data skew is causing the problem.

Step 2: Find the skewed key or column

Most of the time, skew happens during **joins** or **groupBy** operations. So I check which column is causing the issue.

To do this, I run a simple aggregation to count values in the key column:

df.groupBy("join_column").count().orderBy("count", ascending=False).show(10)

This shows me if any value is appearing much more than the rest. For example, if one customer_id appears 5 million times and the rest only appear a few times, that's a skewed key.

Step 3: Use broadcast join if one side is small

If I'm doing a join, and one of the DataFrames is small (less than a few MBs), I can use a broadcast join. This avoids shuffle altogether.

from pyspark.sql.functions import broadcast

result = large_df.join(broadcast(small_df), "customer_id")

This works well if small_df fits in memory.

Step 4: Filter or pre-aggregate before join

Sometimes, we're joining a full table when we don't need to. So I filter the data before joining:

filtered_df = large_df.filter("status = 'active'")

result = filtered_df.join(small_df, "id")

Or I group or aggregate early so that less data is shuffled.

Step 5: Salting technique for severe skew

If skew is really bad, I use a trick called salting. It means we artificially spread out the skewed key so that shuffle is balanced.

Step-by-step salting:

1. Add a random number to one side of the join:

from pyspark.sql.functions import rand, floor left_df = skewed_df.withColumn("salt", floor(rand() * 10))

2. Duplicate the right side data across all salt values:

from pyspark.sql.functions import explode, array, lit right_df = small_df.withColumn("salt", explode(array([lit(i) for i in range(10)])))

3. Perform join on both key and salt:

result = left_df.join(right_df, ["key", "salt"])

This way, we distribute the heavy key across 10 partitions instead of 1.

Step 6: Use Adaptive Query Execution (AQE)

Databricks supports Adaptive Query Execution, which can automatically optimize skewed joins.

I enable AQE in the notebook like this:

spark.conf.set("spark.sql.adaptive.enabled", "true")

This can sometimes fix skew issues automatically by splitting skewed partitions at runtime.

Step 7: Repartition the data

If skew happens because of bad partitioning, I repartition the DataFrame on a better column.

df = df.repartition(100, "customer_id")

This helps to balance partitions before doing any shuffle.

Step 8: Coalesce for smaller datasets

If the issue is not skew but too many small tasks, I use coalesce() to reduce partitions:

df = df.coalesce(10)

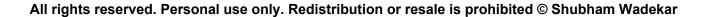
But I never use coalesce() before a shuffle-heavy operation it's more for writing output.

Summary

To fix skew in a Spark job on Databricks, I would:

- 1. Use the Spark UI to confirm task skew
- 2. Identify the skewed key using groupBy counts
- 3. Use broadcast joins when one side is small
- 4. Filter or aggregate early to reduce shuffle
- 5. Use salting to distribute skewed keys evenly
- 6. Enable Adaptive Query Execution
- 7. Repartition the DataFrame on a better column
- 8. Use coalesce wisely when writing output

By combining these methods, I can resolve skew-related failures and improve job performance.



Scenario 57. How would you analyze website traffic using Databricks to find popular pages?

If this question is asked in an interview, I would explain that the goal is to process website logs (like clickstream data), analyze which web pages are getting the most traffic, and maybe identify trends or patterns. Azure Databricks is ideal for this because it can handle large-scale log data using Spark and Delta Lake.

Here's how I would approach this step by step, using very simple logic:

Step 1: Ingest website traffic logs into Databricks

Usually, website logs are stored in formats like JSON, CSV, or Apache log format, and they might be coming from Azure Data Lake, blob storage, or an event hub.

Let's assume logs are stored in Azure Data Lake in JSON format. I would load the data like this:

df = spark.read.json("abfss://logs@mydatalake.dfs.core.windows.net/website_logs/")

Step 2: Understand the structure of the data

I'd first check what fields are available in the logs. Normally, website logs will have:

- timestamp
- user_id
- session_id
- url
- referrer
- user_agent
- response_code

I check it using:

df.printSchema()

df.show(5)

Let's say we have a field called url that tells us which page the user visited.

Step 3: Clean and filter the data

Next, I remove records that are not relevant (like failed requests or bots), so we analyze only real traffic.

from pyspark.sql.functions import col

df_cleaned = df.filter((col("response_code") == 200) & ~col("user_agent").contains("bot"))

Step 4: Count page visits to find popular pages

Now I group by url and count how many times each page was visited.

from pyspark.sql.functions import desc

page_counts = df_cleaned.groupBy("url").count().orderBy(desc("count"))

This gives me the most visited pages at the top.

Step 5: (Optional) Time-based analysis

If I want to see how traffic changes over time (e.g., hourly or daily), I use timestamp transformation:

from pyspark.sql.functions import to_date

df_with_date = df_cleaned.withColumn("date", to_date("timestamp"))

daily_traffic = df_with_date.groupBy("date", "url").count().orderBy("date", desc("count"))

This helps understand which pages are trending on which days.

Step 6: Store results into Delta table or visualize

To store the results:

page_counts.write.format("delta").mode("overwrite").save("/mnt/delta/popular_pages")

Or if I want to visualize it in Databricks notebook:

display(page_counts.limit(10))

This gives a bar chart of top 10 pages.

Step 7: Schedule as a daily job (optional)

If this is a recurring task, I'd use Databricks Jobs to run this notebook daily.

Summary

To find popular pages from website traffic in Databricks, I would:

- 1. Load the log data from storage
- 2. Understand the data structure
- 3. Filter out bots and failed requests
- 4. Group by url and count visits
- 5. Optionally analyze trends by date
- 6. Store results in Delta or visualize them
- 7. Schedule the notebook if needed

This approach is scalable, and we can easily extend it for more advanced analytics like user journeys or conversion tracking later.

Scenario 58. How would you implement a real-time fraud detection pipeline in Databricks?

If this question is asked in an interview, I would explain that fraud detection in real-time usually means processing streaming data (like credit card transactions), applying some logic or machine learning to detect suspicious activity, and then taking quick action like sending alerts.

I would break this down step by step, using a simple explanation:

Step 1: Set up a real-time data stream

Let's say we are getting live credit card transactions through Azure Event Hubs or Kafka. In Databricks, I use Structured Streaming to connect to that source.

Example with Event Hubs:

from pyspark.sql.types import StructType, StringType, DoubleType, TimestampType from pyspark.sql.functions import from_json, col

```
# Define the schema of incoming transactions
schema = StructType() \
    .add("transaction_id", StringType()) \
    .add("user_id", StringType()) \
    .add("amount", DoubleType()) \
    .add("location", StringType()) \
    .add("timestamp", TimestampType())

# Read stream from Event Hub
raw_stream = (spark.readStream
    .format("eventhubs")
    .options(**event_hub_config)
    .load())

# Parse the body of the message as JSON
transactions = raw_stream.select(from_json(col("body").cast("string"), schema).alias("data")).select("data.*")
```

Step 2: Define fraud detection logic

I can start with simple rules like:

- If a transaction amount is over a threshold (say \$10,000)
- Or if there are multiple transactions from different locations within a short time for the same user

Example:

```
from pyspark.sql.functions import when
fraud_candidates = transactions.withColumn(
   "is_fraud",
   when(col("amount") > 10000, True).otherwise(False)
)
```

Later, I can replace this logic with a trained machine learning model.

Step 3: Apply a machine learning model (optional)

If I have trained a model for fraud detection (like using logistic regression or XGBoost), I can use MLflow to load it and apply it in real time.

Example:

import mlflow.pyfunc

model = mlflow.pyfunc.load_model("models:/fraud_model/Production")

```
def predict_fraud(batch_df, batch_id):
    predictions = model.predict(batch_df.toPandas())
```

convert back to Spark DataFrame if needed
spark_df = spark.createDataFrame(predictions)

spark_df.write.format("delta").mode("append").save("/mnt/delta/fraud_predictions")

fraud_candidates.writeStream.foreachBatch(predict_fraud).start()

Step 4: Write suspicious transactions to a sink (alerts or database)

I can store flagged fraud records in a Delta table, send them to Azure SQL, or trigger alerts through email or webhook.

Example to write to Delta:

```
(fraud_candidates.filter("is_fraud == true")
```

- .writeStream
- .format("delta")
- .outputMode("append")
- .option("checkpointLocation", "/mnt/checkpoints/fraud_detection")
- .start("/mnt/delta/fraud_alerts"))

Step 5: Monitor and scale

I monitor the streaming job using the Databricks UI (Streaming tab) and configure auto-scaling clusters so it can handle increasing traffic.

Step 6: Optional - Use Delta Live Tables for better orchestration

If I want a cleaner and more reliable pipeline, I can define this as a Delta Live Table (DLT) pipeline. That helps with error handling, monitoring, and managing dependencies between steps.

Summary

To implement a real-time fraud detection pipeline in Databricks, I would:

- 1. Ingest live transactions using Structured Streaming from Event Hubs or Kafka
- 2. Parse and clean the data
- 3. Apply rule-based or ML-based fraud logic
- 4. Filter and store or alert on suspicious transactions
- 5. Use Delta tables or SQL sinks for storage
- 6. Monitor and scale using the Databricks UI and auto-scaling
- 7. Optionally, use Delta Live Tables to simplify and manage the full pipeline

This pipeline gives fast detection and is scalable to handle millions of records per day.

Scenario 59. How would you configure auto-scaling in Databricks clusters to reduce cost?

If this question is asked in an interview, I would explain that auto-scaling helps manage compute resources efficiently by automatically increasing or decreasing the number of workers in a cluster based on the workload. This avoids over-provisioning and helps reduce costs without compromising performance.

Here's how I would explain it step by step, in simple terms:

Step 1: Understand what auto-scaling means in Databricks

In Databricks, a cluster can scale up (add more worker nodes) when the job requires more resources and scale down (remove idle nodes) when the job is less active. This helps in saving cost, especially in batch jobs, notebooks, and streaming workloads.

Step 2: When creating the cluster, enable auto-scaling

When I create a new cluster or set up a job cluster, I would configure the minimum and maximum number of workers.

For example:

- Minimum workers = 1
- Maximum workers = 8

That means the cluster can start small with 1 worker and automatically grow up to 8 if needed.

In the Databricks UI:

- 1. Go to Clusters → Create Cluster
- 2. Under Worker Type, choose your node size (e.g., Standard_DS3_v2)
- 3. Under Autoscaling, check the box Enable autoscaling
- 4. Set:

"max_workers": 8

"spark_version": "13.3.x-scala2.12",

"node_type_id": "Standard_DS3_v2"

},

}

```
Min workers = 1

Max workers = 8 (or based on workload size)

In JSON format (for jobs), it looks like:

{

"num_workers": 1,

"autoscale": {

"min_workers": 1,
```

Step 3: Use job clusters instead of all-purpose clusters (when possible)

If I'm running a scheduled job, I prefer job clusters instead of keeping all-purpose clusters always on. Job clusters spin up when the job starts and shut down when it finishes, which saves cost.

While scheduling a job:

- I choose "New Job Cluster"
- Set autoscaling in the cluster config

Step 4: Reduce cluster idle time

To save more cost, I set the auto termination time. For example, I set it to 10 minutes, which means if the cluster is idle for 10 minutes, it automatically shuts down.

In the cluster settings:

• Set Terminate after = 10 minutes of inactivity

Step 5: Use smaller worker nodes if the workload allows

Sometimes instead of using large nodes, using more smaller nodes with auto-scaling is more cost-effective.

I choose a smaller worker type like Standard_DS3_v2 or Standard_D4_v2 to reduce per-hour costs.

Step 6: Monitor cluster usage and tune settings

Databricks provides cluster metrics like CPU usage, memory usage, and number of active tasks. Based on that, I may adjust:

- Minimum workers (if under-utilized)
- Maximum workers (if hitting limits too often)
- Worker type (switch to spot instances if cost is a concern)

Step 7: Use Pools (optional)

If I want to reduce startup time and cost, I can use cluster pools. These keep idle VMs ready, so I don't pay for full usage but get faster startups with cost benefits.

Summary

To configure auto-scaling in Databricks to reduce cost, I would:

- 1. Enable autoscaling by setting min and max workers
- 2. Use job clusters instead of all-purpose clusters
- 3. Set idle timeout for auto-termination
- 4. Choose appropriate (smaller) worker node types
- 5. Monitor cluster metrics and adjust settings
- 6. Optionally use cluster pools to improve cost and speed

This way, I make sure the cluster uses only what it needs and shuts down when idle, which helps reduce the overall Databricks cost.

Scenario 60. How do you ensure compliance when handling regulated data in Databricks?

If this question is asked in an interview, I would explain that when working with regulated data like healthcare, financial, or personal information (PII), I need to follow strict rules to keep that data secure, private, and auditable. In Databricks, I take several steps to make sure I meet compliance requirements like GDPR, HIPAA, or SOC2.

Here's how I would explain it in simple and detailed steps:

Step 1: Classify the data

First, I try to understand and label what kind of regulated data I'm working with. This could include:

- PII (like name, email, SSN)
- Financial data (like credit card numbers)
- Health data (like medical records)

Based on that, I tag the datasets and apply the right security rules.

Step 2: Use Unity Catalog for access control and governance

Databricks Unity Catalog helps me manage access to sensitive data easily:

- I assign permissions at table, schema, and column levels
- I use column-level access control to restrict PII fields only to authorized users
- I can mask or hide sensitive columns using views

Example:

CREATE OR REPLACE VIEW secure_view AS

SELECT

name,

email,

CASE

WHEN current_user() IN ('compliance_user') THEN ssn

ELSE '***MASKED***'

END AS ssn

FROM raw_data;

This way, unauthorized users don't see sensitive data even if they can query the table.

Step 3: Secure data storage using encryption

All data stored in Databricks (like in Delta Lake or ADLS) is encrypted at rest and in transit.

- For ADLS Gen2, I enable Microsoft-managed or customer-managed keys (CMK)
- Databricks automatically encrypts data while transferring between notebooks and storage

Step 4: Use secret scopes for credentials

I never hardcode passwords, API keys, or connection strings in my notebooks.

Instead, I store them securely using Databricks secrets:

username = dbutils.secrets.get(scope="db_scope", key="sql_user")

password = dbutils.secrets.get(scope="db_scope", key="sql_pass")

This prevents sensitive credentials from being leaked or logged.

Step 5: Set up audit logging

To be compliant, I make sure all activities (reads, writes, queries, access) are tracked.

In Databricks:

- I enable audit logging to log events like job runs, notebook access, user logins
- Logs can be sent to Azure Log Analytics, Storage Account, or SIEM tools for monitoring

This helps in proving compliance during audits.

Step 6: Enable data retention and purging policies

For regulations like GDPR, I must ensure users can request deletion of their data.

In Delta Lake, I can use time travel to access historical data, but I also use the VACUUM command to remove deleted data permanently:

DELETE FROM users WHERE user_id = '1234';

VACUUM users RETAIN 0 HOURS;

This ensures the data is really gone.

Step 7: Network and workspace security

- I use private link or VNet injection to isolate my Databricks workspace from the public internet
- I enable IP access lists to restrict who can access Databricks
- I configure SAML or Azure AD for secure user authentication and role-based access

Step 8: Documentation and training

Finally, I document all my data flows, access policies, and compliance steps, and make sure my team is trained on handling regulated data securely.

Summary

To ensure compliance when handling regulated data in Databricks, I follow these steps:

- 1. Classify sensitive data like PII or financial records
- 2. Use Unity Catalog for fine-grained access control
- 3. Encrypt data at rest and in transit
- 4. Manage secrets using secret scopes
- 5. Enable audit logging for traceability
- 6. Handle data deletion and retention properly
- 7. Secure the network and workspace
- 8. Maintain documentation and team awareness

These practices help me meet regulatory requirements and build secure, compliant pipelines in Databricks.

Scenario 61. How would you automate ETL workflows using Databricks jobs and triggers?

If this is asked in an interview, I would explain that automating ETL workflows in Databricks means setting up a repeatable, scheduled process to extract, transform, and load data using Databricks jobs. I can use triggers like time-based schedules or external events to run the jobs automatically without manual effort.

Here's how I would approach this in a simple and step-by-step way:

Step 1: Build your ETL logic inside a notebook or script

First, I create a Databricks notebook (or a Python/Scala script) that contains all my ETL steps. For example:

- Extract data from Azure Data Lake or SQL database
- Clean, filter, and transform the data using PySpark or SQL
- Load the transformed data into Delta Lake, Azure SQL, or another destination

Example ETL code in a notebook:

```
# Extract
```

df = spark.read.format("csv").option("header",
"true").load("abfss://raw@datalake.dfs.core.windows.net/sales.csv")

Transform

df_cleaned = df.filter(df["amount"] > 0).withColumn("month", month("date"))

Load

df_cleaned.write.format("delta").mode("overwrite").save("/mnt/silver/sales")

Step 2: Create a Databricks Job to automate the notebook

Now that my ETL notebook is ready, I create a Databricks job to run it automatically.

- 1. Go to the Jobs tab in the Databricks workspace.
- 2. Click Create Job
- 3. Enter a job name like daily_sales_etl
- 4. Under Task, choose:

Type: Notebook (or Python script)

Notebook path: The path to the ETL notebook

- 5. Choose or create a job cluster to run the job
- 6. Optionally add parameters if the notebook accepts input

Step 3: Set up a schedule trigger

To automate the job, I add a trigger:

- Choose Schedule as the trigger type
- Set a CRON expression or simple frequency like:
 - Every day at 1 AM
 - Every hour
- Example CRON: 0 1 * * * (means run at 1:00 AM daily)

Databricks will now run this ETL job every day at the set time.

Step 4: Add job dependencies (optional for multi-step workflows)

If my ETL has multiple stages, like Bronze \rightarrow Silver \rightarrow Gold, I can create multiple tasks in the same job and set dependencies:

- Task 1: Load raw data into Bronze
- Task 2: Clean and write to Silver (depends on Task 1)
- Task 3: Aggregate and write to Gold (depends on Task 2)

This creates a job DAG (Directed Acyclic Graph) where each task runs in order.

Step 5: Set retry policies and alerts

To make the ETL job reliable, I:

- Enable retry on failure (e.g., 2 retries with 10-minute delay)
- Configure email or webhook alerts if the job fails or succeeds

This helps me monitor the automation and avoid manual checks.

Step 6: Trigger the job from external systems (optional)

Sometimes, I need to run the ETL job from another tool like Azure Data Factory, Power Automate, or a REST API.

Databricks provides a Jobs REST API that can be used to trigger jobs:

curl -X POST https://<databricks-instance>/api/2.1/jobs/run-now \

-H "Authorization: Bearer <token>" \

-d '{"job_id": 123}'

This lets me start a job programmatically after a file lands in the lake or when another pipeline finishes.

Step 7: Monitor job runs and logs

Databricks automatically keeps **job run history** with logs, output, and status. I use this to:

- Check which runs succeeded or failed
- See how long each task took
- Debug any errors from logs

This helps me ensure the ETL is stable and trustworthy.

Summary

To automate ETL workflows using Databricks jobs and triggers, I follow these steps:

- 1. Build the ETL logic inside a notebook or script
- 2. Create a Databricks job to run that notebook
- 3. Add a schedule trigger (like daily or hourly)
- 4. Create multiple tasks with dependencies if needed
- 5. Set retry logic and failure alerts
- 6. Optionally trigger jobs from external systems using the API
- 7. Monitor job history and logs for visibility

This setup helps me run ETL pipelines automatically with minimal manual work and better reliability.

Scenario 62. How do you test notebooks and pipelines in Databricks before production deployment?

If this is asked in an interview, I would explain that before moving my notebooks or ETL pipelines to production in Databricks, I always make sure they are fully tested to avoid failures, data loss, or incorrect results. Testing in Databricks involves validating both the code logic and the data flow in a safe environment.

Here's how I approach it step by step in a very simple and detailed way:

Step 1: Use a development or staging environment

I don't run my notebooks directly in production. I always test them in a separate workspace, or at least with test data, so real production data is not affected.

For example:

- I use a test storage location like /mnt/dev/ instead of /mnt/prod/
- I create a dev database like dev_sales_db with sample or anonymized data

This way, I test the full logic safely.

Step 2: Test with sample data

I create small datasets that are easy to inspect and use them to test the logic.

Example:

from pyspark.sql import Row

sample_data = [Row(customer="Alice", amount=100), Row(customer="Bob", amount=-50)]

df = spark.createDataFrame(sample_data)

df.filter("amount > 0").show()

This helps me quickly check if filters, joins, and calculations are working as expected.

Step 3: Use assertions to validate logic

Just like writing unit tests in Python, I use **assertions** in notebooks to validate intermediate results.

Example:

assert df.count() == 2, "Row count is not matching"

assert df.filter("amount < 0").count() == 1, "Negative amount not handled"

If these assertions fail, I know my logic is incorrect before deploying.

Step 4: Use Databricks Repos for version control and testing branches

If I'm working in a team, I use Databricks Repos (connected to GitHub or Azure DevOps). I:

- Create a new branch for testing
- Make and test changes in that branch
- Once it's tested and approved, merge to main for deployment

This gives a clean separation between development and production.

Step 5: Run notebook as a Databricks job in test mode

I create a Databricks job for the notebook and run it with test parameters.

For example, if the notebook accepts a date or environment as a parameter:

```
dbutils.widgets.text("env", "dev")
env = dbutils.widgets.get("env")
if env == "prod":
    data_path = "/mnt/prod/sales"
else:
    data_path = "/mnt/dev/sales"
```

This way, I can test how the notebook behaves with different inputs.

Step 6: Validate pipeline DAGs

If I'm working with multiple notebooks in a pipeline (job with multiple tasks), I test the full sequence.

- I make sure task dependencies work
- I check that intermediate outputs are correct
- I simulate failures by inserting wrong inputs or forcing errors

This helps identify weak points before running the job in production.

Step 7: Use Unit Testing Libraries (optional but good practice)

For more advanced testing, especially for Python code used in Databricks, I use libraries like pytest.

I move common logic to Python modules and write unit tests like:

```
def test_cleaning_function():
  input_data = [{"amount": 100}, {"amount": -10}]
  cleaned = clean_data(input_data)
  assert len(cleaned) == 1
```

I run these tests locally or using %run in Databricks.

Step 8: Peer review before deployment

Before deployment, I always ask a team member to review the code. This helps catch:

- Hardcoded values
- Forgotten test paths
- Performance issues

In Databricks Repos, this is done through pull requests.

Summary

To test notebooks and pipelines in Databricks before production deployment, I follow these steps:

- 1. Use a development or staging environment
- 2. Test with small, sample datasets
- 3. Use assertions to validate logic inside notebooks
- 4. Use Databricks Repos to test in a Git branch
- 5. Run the notebook as a test job with parameters
- 6. Test full pipelines with task dependencies
- 7. Use pytest for unit testing Python functions if applicable
- 8. Do a peer review before merging or deploying

This process helps me catch bugs early, maintain quality, and ensure a smooth deployment to production.

Scenario 63: How would you structure code and notebooks for a large team project in Databricks?

Answer:

For a large team project in Databricks, I like to follow a clean and consistent folder and notebook structure so that the whole team can work together smoothly without confusion. Here's how I generally approach it:

1. Project Folder Structure:

I start by creating a root folder in the Databricks Workspace for the project, and inside that, I organize the work into clear subfolders like this:

- 01_Ingestion: Contains notebooks or scripts that connect to data sources (like Azure Data Lake, SQL DB, APIs, etc.) and load the raw data.
- **02_Processing:** Here we clean and structure the data (like parsing, schema corrections).
- **03_Transformation:** This is where business logic is applied joins, aggregations, feature engineering, etc.
- **04_ML_Models:** If the project has any machine learning component.
- **05_Validation:** Data quality checks or test cases to ensure data integrity.
- **06_Visualization:** Dashboards or notebooks for exploratory data analysis and business insights.
- 07_Utilities: Helper functions, common reusable code, UDFs, or libraries.
- **08_Configs:** Contains parameter files, secrets scope references, and environment-specific settings.
- 99_Archive: Deprecated or backup notebooks (helps avoid deleting things directly).

2. Notebook Structure:

Each notebook usually follows a similar pattern:

- 1. **Header** with project name, author, purpose of the notebook, and last updated date (as markdown).
- 2. **Imports and Configs**: All required libraries and configs (like DB secrets, file paths) are defined at the top.
- 3. **Functions Section**: I define helper functions here, even if they are used only once helps in readability.
- 4. Main Logic: The core ETL logic goes here in step-by-step cells.
- 5. **Output or Display**: For final display, write, or register table/view.

3. Code Reusability and Collaboration:

For large teams, we usually follow these best practices:

- **Use Git Integration:** I connect the Databricks workspace to a Git repo (like Azure DevOps or GitHub) so multiple team members can work with version control.
- Create Shared Libraries (if needed): For large projects, we sometimes move reusable code (like transformations or validations) into .py files and install them as shared libraries in Databricks.
- **Parameterization:** I use dbutils.widgets to make notebooks dynamic especially for scheduled jobs or notebooks reused in multiple pipelines.

4. Cluster and Job Management:

- We tag notebooks clearly if they are meant to be run as a Job (ETL pipelines).
- We separate development clusters and production job clusters to avoid conflict.
- We use Job clusters with pre-defined libraries to ensure consistency across environments.

5. Documentation and Naming:

- I make sure notebook names are prefixed with numbers to indicate sequence (like 01_Load_Customer_Data).
- We write markdown documentation inside the notebooks to explain logic helpful when someone new joins the project.
- For each folder, we also add a README notebook explaining the purpose of that folder.

So overall, the goal is to keep things modular, well-documented, version-controlled, and easy to understand so that even if a new developer joins the project tomorrow, they can easily figure out where to look and what each notebook is doing.

Scenario 64: How would you structure code and notebooks for a large team project in Databricks?

How do you handle large joins and optimize them in Databricks?

When I'm working with large joins in Databricks, especially with big datasets, I follow a few important steps to make sure the join runs efficiently and doesn't fail due to memory or performance issues.

Let me explain how I handle and optimize large joins, step by step:

1. Use Broadcast Join When One Table is Small

If one of the tables is small enough to fit into the memory of each executor, I use broadcast join.

In Databricks, Spark automatically decides to broadcast small tables, but sometimes I use:

from pyspark.sql.functions import broadcast

df_joined = large_df.join(broadcast(small_df), "id", "inner")

This is very fast because the small table is copied to all executors so there's no shuffling.

2. Filter Early to Reduce Data Size

Before joining, I always try to filter both tables as much as possible:

```
large_df_filtered = large_df.filter("country = 'US'")
```

small_df_filtered = small_df.filter("status = 'active'")

This reduces the number of rows Spark has to process during the join. It saves memory and speeds up performance.

3. Select Only Required Columns Before Join

I never do joins on full tables I only select the columns I really need:

```
df1 = df1.select("id", "name")
```

df2 = df2.select("id", "order_total")

Less data = less shuffle = better performance.

4. Join on Partitioned or Bucketed Columns (if available)

If I know the data is already **partitioned or bucketed** by the join column (like user_id), I try to leverage that.

For example, in Delta tables, I may use:

OPTIMIZE my_table ZORDER BY (user_id)

This helps Spark read only the necessary files during the join which improves I/O and join time.

5. Use Appropriate Join Type

I choose the right join type depending on the use case:

- Inner Join Fastest and safest when I only want matching records.
- Left Join / Outer Join Only when needed, since they are slower and bring more nulls.
- Semi/Anti Join Tuse them if I just need to check existence:

df1.join(df2, "id", "left_semi")

This is faster than a full join.

6. Check Spark UI for Skew and Shuffle Issues

After running the join, I check the Spark UI to see:

- If any tasks are taking too long
- If there is data skew (like one key has too many rows)
- If shuffle size is very high

If there's data skew, I sometimes use salting by adding a random number to the join key to distribute the load more evenly.

7. Repartition Before Join (When Needed)

If both tables are huge and being shuffled, I do a controlled repartition before the join:

df1 = df1.repartition(100, "user_id")

df2 = df2.repartition(100, "user_id")

This makes sure the data is evenly distributed across workers, and prevents data skew and outof-memory issues.

8. Use Delta Tables and Caching Wisely

If I use Delta tables, I make sure the tables are OPTIMIZED and VACUUMED regularly.

If I'm reusing the same dataset multiple times in the notebook, I use:

df.cache()

This avoids recomputation during join or transformations.

So to summarize:

- I filter early, select only needed columns
- Use broadcast joins when possible
- Repartition for large shuffles
- Check for skew in the Spark UI
- Use semi/anti joins for existence checks
- And always keep the data optimized and clean, especially in Delta format

This way, even when working with very large datasets like in terabytes the joins are fast, stable, and cost-effective.

Scenario 64: Design an end-to-end ETL pipeline from an on-premise database to Databricks

Overview: The ETL pipeline will extract data from an on-premise database, transform it in Databricks, and load it into a target storage for analytics.

Extract

- Source Connection: Use JDBC drivers to connect to the on-premise database (e.g., SQL Server, Oracle, MySQL). Configure connection strings in Databricks secrets to securely store credentials.
- **Data Ingestion**: Use Apache Spark's spark.read.jdbc to extract data incrementally (e.g., based on a timestamp or ID column). For example:

df = spark.read.jdbc(url="jdbc:mysql://on-prem-host:3306/db", table="table_name",
properties={"user": "user", "password": dbutils.secrets.get("scope", "password")})

- **Incremental Loads**: Implement a watermark or high-water mark (e.g., last_updated column) to fetch only new or updated records.
- Data Transfer: Stage data in a cloud storage layer (e.g., ADLS, S3) using Databricks' dbutils.fs or cloud SDKs to transfer files securely.

Transform

- **Data Processing**: Use Databricks notebooks with Spark SQL or PySpark for transformations (e.g., filtering, joins, aggregations). Leverage Delta Lake for ACID transactions and reliability.
- **Optimization**: Apply Spark optimizations like caching, partitioning (e.g., by date or region), and broadcast joins for small tables.
- **Schema Management**: Use Delta Lake's schema enforcement or schema evolution to handle changes in source data structure.

Load

- **Target Storage**: Write transformed data to Delta Lake tables in a cloud storage layer (e.g., ADLS Gen2, S3) using df.write.format("delta").save("/path/to/delta").
- **Table Creation**: Register the Delta table in the Databricks metastore for querying:

spark.sql("CREATE TABLE target_table USING DELTA LOCATION '/path/to/delta'")

• **Orchestration**: Use Databricks Workflows (or Apache Airflow) to schedule and orchestrate the pipeline, ensuring dependencies and retries.

Monitoring

- Use Databricks' job monitoring and logging to track pipeline performance and failures.
- Set up alerts for job failures or SLA breaches using Databricks notifications or integration with tools like PagerDuty.

Tools

- Databricks Runtime for Spark processing.
- Delta Lake for reliable storage.
- Cloud storage (ADLS/S3) for intermediate and final data.
- Databricks Workflows for orchestration.

Scenario 65. How do you handle schema evolution in a Parquet file?

Schema evolution in Parquet files (used in Delta Lake) allows for changes in the data schema (e.g., adding, dropping, or modifying columns) while maintaining compatibility.

Challenges with Parquet:

- Parquet is schema-on-read but not tolerant to changes unless handled properly.
- Issues arise when columns are added, reordered, or changed types.

Approach

- **Use Delta Lake**: Parquet alone is immutable, but Delta Lake, built on Parquet, supports schema evolution:
- **Schema-on-Write**: By default, Delta enforces schema. To allow evolution, set mergeSchema to true when writing:

df.write.format("delta").option("mergeSchema", "true").save("/path/to/delta")

- This merges new columns into the existing schema.
- **Schema-on-Read**: Use spark.read.option("mergeSchema", "true") to handle varying schemas during reads.

Handling Specific Changes

- Adding Columns: New columns are automatically added with mergeSchema.
- **Dropping Columns**: Delta doesn't physically delete columns but marks them as dropped in metadata. Use ALTER TABLE DROP COLUMN if needed.
- **Type Changes**: Handle cautiously, as type mismatches can break pipelines. Use explicit casting in transformations or update schemas manually with ALTER TABLE.

Best Practices

- Validate schema changes in a staging environment before applying to production.
- Use Delta's time travel to recover from problematic schema updates (RESTORE TABLE).
- Document schema changes in a data catalog (e.g., Unity Catalog).
- Create **versioned folders** per schema change.
- Validate schemas using StructType before ingestion.
- Leverage cloudFiles.schemaEvolutionMode = addNewColumns in Auto Loader.

Append new data with additional columns new_df.write.format("delta").option("mergeSchema", "true").mode("append").save("/delta/table")

Scenario 66. What is watermarking in streaming data processing?

Watermarking in streaming data processing (e.g., Spark Structured Streaming) defines a threshold for how late data can arrive before it's ignored, ensuring bounded state and timely processing.

 Manages late-arriving events in event-time processing by setting a watermark based on the event timestamp.

df = df.withWatermark("event_timestamp", "10 minutes")

• This means Spark will ignore events arriving more than 10 minutes after the latest event timestamp seen.

Example

from pyspark.sql.functions import window streaming_df = df.groupBy(window("event_timestamp", "5 minutes")).count()

Effect:

- Ensures efficient state cleanup.
- Protects against unbounded state growth.
- Balances completeness vs performance.

Behavior

- Spark tracks the maximum event timestamp and discards events older than max_timestamp watermark_interval.
- Ensures stateful operations (e.g., aggregations) don't grow indefinitely.

Best Practices

- Set watermark duration based on expected data latency.
- Combine with outputMode("update") or append for efficient writes to Delta tables.
- Monitor watermark delays using Spark UI to tune the threshold.

Scenario 67. How do you handle late-arriving data in a batch pipeline?

Options:

- Delta Lake MERGE
- Partition repair: Re-run or merge new records into the existing partitions (e.g., hourly/daily).
- Change Data Feed (CDF) in Delta: Capture and reprocess late records.
- Flagging and tagging: Add is_late boolean and actual_event_ts vs ingest_ts.

Use Delta Lake's MERGE operation to handle late-arriving data by updating existing records or inserting new ones:

```
spark.sql("""

MERGE INTO target_table t

USING source_table s

ON t.id = s.id AND t.event_date = s.event_date

WHEN MATCHED THEN UPDATE SET *

WHEN NOT MATCHED THEN INSERT *

""")
```

Partitioning: Partition data by a time-based column (e.g., event_date) to limit the scope of updates:

df.write.format("delta").partitionBy("event_date").save("/delta/target")

Incremental Processing:

- Track the last processed timestamp or ID in a metadata table.
- Use a sliding window (e.g., last 7 days) to reprocess late data:

df = spark.read.table("source").filter(f"event_date >= '{last_processed_date}'")

• **Reprocessing**: For significant delays, trigger a full or partial reprocessing job for affected partitions using Databricks Workflows.

Best Practices:

- Use Delta Lake for efficient upserts and versioning.
- Log late-arriving data for auditing and monitoring.
- Test reprocessing logic in a staging environment.

Scenario 68. How do you implement data quality validation and unit testing inside Databricks notebooks?

Validation Strategy:

Use PyDeequ, Great Expectations, or custom code for:

- Null checks, schema conformance, referential integrity.
- Distribution checks, uniqueness, value thresholds.

```
from great_expectations.dataset import SparkDFDataset ge_df = SparkDFDataset(df)
ge_df.expect_column_values_to_not_be_null("id")
ge_df.expect_column_values_to_be_unique("id")
validation_results = ge_df.validate()
```

Custom Rules: Write custom Spark SQL or PySpark checks:

```
from pyspark.sql.functions import col
null_check = df.filter(col("id").isNull()).count()
if null_check > null_check > 0:
    raise ValueError("Null IDs detected")
```

Delta Live Tables (DLT): Use DLT for built-in data quality checks:

```
@dlt.table
def validated_table():
    df = dlt.read("source_table")
    return df.where(col("id").isNotNull()).expect("no_null_ids", "id IS NOT NULL")
```

Unit Testing

• Use pytest with Databricks Connect or Databricks notebooks for unit tests:

```
def test_data_quality():
    df = spark.read.table("test_table")
    assert df.filter(col("id").isNull()).count() == 0, "Null IDs found"
```

Create small, representative datasets for testing transformations:

```
test_data = [(1, "A"), (2, "B")]
test_df = spark.createDataFrame(test_data, ["id", "name"])
```

• **Notebook Structure**: Organize notebooks with modular functions for transformations, making them easier to test:

```
def transform_data(df):
    return df.withColumn("new_col", col("id") * 2)
```

Best Practices: Store validation results in a Delta table for auditing. Integrate tests into CI/CD pipelines using Databricks Workflows or GitHub Actions. Use Unity Catalog for governance and metadata tracking.

Scenario 69. Cleaning and Transforming Messy Sales Data for Reporting Tables.

Cleaning

• **Remove Duplicates**: Use dropDuplicates based on a unique key (e.g., order_id):

df = df.dropDuplicates(["order_id"])

Handle Missing Values:

• Fill numeric columns with defaults (e.g., 0 for quantities):

df = df.fillna({"quantity": 0, "price": 0.0})

• Drop rows with critical missing fields (e.g., order_id):

df = df.dropna(subset=["order_id"])

Standardize Formats:

Convert dates to a consistent format:

from pyspark.sql.functions import to_date
df = df.withColumn("order_date", to_date("order_date", "yyyy-MM-dd"))

• Normalize strings (e.g., trim, lowercase):

df = df.withColumn("customer_name", lower(trim(col("customer_name"))))

Transformation

• Aggregations: Compute metrics like total sales by region or product:

sales_agg = df.groupBy("region", "product").agg(sum("price").alias("total_sales"))

• Enrich Data: Join with dimension tables (e.g., product catalog, customer data):

df = df.join(product_df, "product_id", "left")

• Feature Engineering: Add derived columns like profit margin:

df = df.withColumn("profit_margin", (col("price") - col("cost")) / col("price"))

Designing Reporting Tables:

- **Fact Table**: Store transactional data (e.g., sales_fact with columns: order_id, customer_id, product_id, order_date, quantity, price).
- **Dimension Tables**: Create tables for customers, products, and dates for efficient joins.
- **Star Schema** sales_fact(Fact table), dim_customer, dim_product, dim_date (Dimension tables).
- Partitioning: Partition fact tables by order_date for performance:

df.write.format("delta").partitionBy("order_date").saveAsTable("sales_fact")

• Materialized Views: Use Delta Lake to create aggregated views for common queries (e.g., monthly sales by region).

Tools:

- Delta Lake for storage and versioning.
- Unity Catalog for governance.
- Databricks SQL for dashboarding and reporting.

Scenario 70. How would you design a data pipeline to handle daily logs from multiple sources?

Ingestion

• **Sources**: Use Databricks Auto Loader to ingest logs from cloud storage (e.g., S3, ADLS) incrementally:

```
df = spark.readStream.format("cloudFiles") \
    .option("cloudFiles.format", "json") \
    .load("s3://logs/")
```

- Schema Inference: Auto Loader infers schemas for JSON/CSV or allows manual schema definition to handle schema drift.
- File Discovery: Use cloudFiles to automatically detect new files.

Processing

• **Unify Schema**: Normalize logs from different sources into a common schema using PySpark transformations:

```
from pyspark.sql.functions import col
unified_df = df.select(
  col("source1_field").alias("common_field"),
  col("timestamp").cast("timestamp")
)
```

• Handle Errors: Use try-catch blocks or filter invalid records to a quarantine table:

```
invalid_df = df.filter(~col("timestamp").cast("timestamp").isNotNull())
invalid_df.write.format("delta").save("/quarantine")
```

• Enrich Data: Add metadata like ingestion date or source identifier:

```
from pyspark.sql.functions import current_date
df = df.withColumn("ingestion_date", current_date())
```

Storage

Write to Delta Lake tables, partitioned by date and source:

```
df.writeStream.format("delta") \
    .partitionBy("ingestion_date", "source") \
    .outputMode("append") \
    .save("/delta/logs")
```

• Register as a table for querying:

spark.sql("CREATE TABLE logs USING DELTA LOCATION '/delta/logs'")

Orchestration

- Use Databricks Workflows to schedule daily jobs.
- Monitor pipeline health with Databricks' streaming metrics and alerts.

Best Practices:

- Use Auto Loader for scalability and schema evolution.
- Implement data quality checks (e.g., Great Expectations).
- Optimize for cost by using serverless compute or spot instances.

Scenario 71. Design an optimal schema to store event logs (like clicks/swipes) for high-velocity web traffic.

Schema Design

Table

event_logs (Delta Lake table).

Columns

- event_id: UUID or auto-incrementing ID (primary key).
- user_id: String or integer (foreign key to user table).
- event_type: String (e.g., "click", "swipe").
- event_timestamp: Timestamp (event time).
- session_id: String (to group events in a session).
- page_url: String (URL or page identifier).
- device_id: String (device tracking).
- metadata: Struct or JSON (flexible field for additional attributes).
- ingestion_timestamp: Timestamp (when the event was ingested).

Optimizations:

- Enable CDF for downstream processing
- Use Auto Loader with schema evolution for flexibility

Partitioning

 Partition by event_date (derived from event_timestamp) and optionally event_type for query performance:

df.write.format("delta").partitionBy("event_date", "event_type").saveAsTable("event_logs")

• **Indexing**: Use Z-Order indexing on frequently queried columns (e.g., user_id, session_id):

spark.sql("OPTIMIZE event_logs ZORDER BY (user_id, session_id)")

Considerations:

- Scalability: Delta Lake handles high write throughput with ACID transactions.
- Schema Evolution: Use mergeSchema for adding new metadata fields.
- Retention: Implement time-based retention policies with VACUUM to manage storage costs:

spark.sql("VACUUM event_logs RETAIN 30 DAYS")

Storage

• Use cloud storage (ADLS/S3) with Delta Lake for durability and performance.

Scenario 72. Explain how you'd model a "swipe payment" API in a relational schema.

Table

swipe_payments

Columns

- payment_id: BIGINT (Primary Key, auto-incrementing).
- user_id: BIGINT (Foreign Key to users table).
- merchant_id: BIGINT (Foreign Key to merchants table).
- card_id: BIGINT (Foreign Key to cards table).
- amount: DECIMAL(10,2) (Payment amount).
- currency: STRING (e.g., "USD").
- payment_timestamp: TIMESTAMP (When the swipe occurred).
- status: STRING (e.g., "success", "failed", "pending").
- transaction_ref: STRING (Unique reference from payment gateway).
- device_id: STRING (Device used for swipe).
- created_at: TIMESTAMP (Record creation time).

Related Tables

- users: (user_id, name, email, etc.)
- merchants: (merchant_id, name, location, etc.)
- cards: (card_id, user_id, last_four, card_type, etc.)

Relations

- Use FKs between transactions → cards → users.
- Normalize dimensions for easy joins.
- Use fact_transaction table for analytics.

API Data Flow

Input

```
JSON payload from the swipe payment API:
```

```
{
  "user_id": 123,
  "merchant_id": 456,
  "card_id": 789,
  "amount": 50.00,
  "currency": "USD",
  "transaction_ref": "TXN_001",
  "device_id": "POS_001"
}
```

Processing

- Validate input (e.g., check user_id, merchant_id existence).
- Insert into swipe_payments table using Delta Lake:

```
df = spark.read.json("path/to/api_data")
df.write.format("delta").mode("append").saveAsTable("swipe_payments")
```

Output

Return payment_id and status to the API caller.

Best Practices

- Use foreign key constraints (if supported) or enforce via application logic.
- Partition swipe_payments by payment_date (derived from payment_timestamp).
- Use Delta Lake for auditability and time travel.

Scenario 73. How would you optimize a slow-performing JOIN query joining large tables? What indexing or partitioning strategies would you use?

Optimization Strategies

Analyze Query Plan

• Use EXPLAIN to inspect the query plan:

spark.sql("EXPLAIN SELECT * FROM orders o JOIN customers c ON o.customer_id =
c.customer_id").show()

• Identify issues like shuffle joins or full table scans.

Partitioning:

 Partition tables by frequently filtered columns (e.g., order_date for orders, region for customers):

orders_df.write.format("delta").partitionBy("order_date").saveAsTable("orders")

• Ensure JOIN keys are co-partitioned to avoid shuffles:

customers_df.write.format("delta").partitionBy("customer_id").saveAsTable("customers")

Indexing

Apply Z-Order indexing on JOIN keys:

spark.sql("OPTIMIZE orders ZORDER BY (customer_id)")
spark.sql("OPTIMIZE customers ZORDER BY (customer_id)")

• This reduces data scanned during joins.

Broadcast Joins

• If one table is small (<10 GB), use broadcast join to avoid shuffling:

from pyspark.sql.functions import broadcast
df = orders_df.join(broadcast(customers_df), "customer_id")

Caching

• Cache frequently accessed tables or intermediate results:

spark.sql("CACHE TABLE customers")

Skew Handling

- Identify data skew in JOIN keys (e.g., a few customer_id values with many rows).
- Use salting to distribute data evenly:

```
orders_df = orders_df.withColumn("salt", (col("customer_id") % 100))
customers_df = customers_df.withColumn("salt", (col("customer_id") % 100))
df = orders_df.join(customers_df, ["customer_id", "salt"])
```

Query Rewriting

Filter data early to reduce join input size:

filtered_orders = orders_df.filter("order_date >= '2025-01-01'") df = filtered_orders.join(customers_df, "customer_id")

• Use semi-joins for existence checks:

df = orders_df.join(customers_df.select("customer_id"), "customer_id", "left_semi")

File Size Tuning:

- Use OPTIMIZE to compact small files.
- Set target file size (~1GB).

Delta Features:

Use Delta Caching and Data Skipping Stats

Best Practices:

- Monitor query performance using Databricks' Spark UI.
- Use Delta Lake's OPTIMIZE and ZORDER for ongoing maintenance.
- Test optimizations in a staging environment.

Scenario 74. How does Databricks manage the cluster lifecycle, and what is the difference between interactive and job clusters?

Databricks manages the cluster lifecycle through a combination of user configuration, automation (autoscaling, auto-termination), and integration with cloud provider infrastructure. The lifecycle includes:

- Creation: Users define cluster configurations (e.g., instance types, number of workers, runtime version) via the Databricks UI, CLI, or API. Databricks provisions the underlying compute resources in the cloud.
- **Startup**: The cluster is initialized with the specified Databricks Runtime, libraries, and configurations. This includes setting up Spark, installing dependencies, and configuring networking.
- Running: The cluster executes workloads (notebooks, jobs, or queries). Databricks
 provides features like auto-scaling to dynamically adjust the number of workers based
 on workload demand.
- **Termination**: Clusters can be terminated manually or automatically. Interactive clusters remain running until manually stopped or configured to terminate after a period of inactivity (e.g., 2 hours). Job clusters are ephemeral and terminate after job completion.
- **Maintenance**: Databricks handles patching and updating the Databricks Runtime, ensuring security and performance improvements without user intervention.

Interactive vs. Job Clusters

Interactive Clusters

- **Purpose**: Designed for ad-hoc, exploratory data analysis, development, and collaboration in notebooks or dashboards.
- **Lifecycle**: Long-running, manually created, and persist until terminated (manually or via auto-termination).
- **Use Case**: Data scientists and engineers use interactive clusters for iterative development, running notebooks, and visualizing results.
- **Features**: Support multiple concurrent users, notebook execution, and real-time collaboration. Auto-scaling is optional.
- **Cost**: Can be more expensive due to longer runtime unless auto-termination is configured.

Job Clusters

- **Purpose**: Designed for automated, production-grade ETL pipelines or scheduled tasks.
- **Lifecycle**: Ephemeral, created automatically when a job starts and terminated when the job completes.
- **Use Case**: Running scheduled Databricks Jobs (e.g., ETL workflows) with minimal user intervention.
- **Features**: Optimized for single-purpose execution, reducing costs by terminating after completion. No manual interaction required.

• **Cost**: More cost-efficient for scheduled tasks, as resources are only used during job execution.

Key Difference: Interactive clusters are long-running and suited for interactive work, while job clusters are ephemeral, cost-efficient, and optimized for automated, scheduled jobs.

Scenario 75. How would you design an end-to-end ETL pipeline using Databricks notebooks and orchestrate it using Databricks Jobs?

Here's a step-by-step design using Databricks notebooks and orchestration with Databricks Jobs:

Extract

- **Source**: Data can come from cloud storage (e.g., S3, ADLS, GCS), databases (via JDBC), or streaming sources (e.g., Kafka).
- Notebook: Create a notebook to read data using Spark APIs. For example, to read CSV files from S3:

df = spark.read.csv("s3://bucket/source/data.csv", header=True, inferSchema=True)

Use Auto Loader for incremental ingestion (covered later).

Transform

• **Notebook**: Create one or more notebooks for transformations (e.g., filtering, joining, aggregating). Use PySpark or Spark SQL for scalable processing. For example:

from pyspark.sql.functions import col transformed_df = df.filter(col("age") > 18).groupBy("region").agg({"sales": "sum"})

- Leverage Delta Lake for reliable storage and transformation (e.g., upserts, schema evolution).
- Modularize transformations into separate notebooks for reusability and maintainability.

Load/Aggregate

• **Notebook**: Write the transformed data to a target (e.g., Delta Lake table, database, or data warehouse). For example, writing to a Delta table:

transformed_df.write.format("delta").mode("overwrite").save("dbfs:/mnt/delta/target_table")

• Register the Delta table in the metastore for querying:

spark.sql("CREATE TABLE target_table USING DELTA LOCATION 'dbfs:/mnt/delta/target_table'")

Orchestration with Databricks Jobs

- **Job Creation**: In the Databricks UI, create a job and add tasks corresponding to each notebook (e.g., Extract, Transform, Load).
- **Task Dependencies**: Define dependencies between tasks (e.g., Transform depends on Extract, Load depends on Transform). This ensures sequential execution.
- **Cluster Configuration**: Assign a job cluster to the job for cost efficiency or use an existing interactive cluster for development.
- **Scheduling**: Set a schedule (e.g., daily at 2 AM) using the Databricks Jobs UI or API. For example: (Trigger: Cron expression like 0 0 2 * * ? for daily at 2 AM.)
- Parameters: Pass runtime parameters to notebooks (e.g., date ranges) using Databricks widgets or job parameters.
- **Monitoring**: Use the Jobs UI to monitor runs, view logs, and set up alerts for failures (e.g., via email or webhooks).

Best Practices

- Use Delta Lake for reliable storage and ACID transactions.
- Enable auto-scaling on job clusters to optimize resource usage.
- Implement error handling in notebooks (e.g., try-catch blocks) and retry policies in jobs.
- Use Databricks secrets for secure access to credentials (e.g., database passwords).

Scenario 76. How do you handle large-scale file ingestion using Databricks Auto Loader, and what are its benefits over standard Spark file streaming?

Databricks Auto Loader provides a scalable and efficient way to ingest large volumes of files incrementally from cloud storage. Here's how to implement it:

Setup

- Configure Auto Loader to monitor a cloud storage directory (e.g., S3, ADLS, GCS).
 Uses cloudFiles with incremental processing, checkpointing, and schema inference.
- Specify the file format (e.g., CSV, JSON, Parquet) and schema options.
- Example in PySpark:

from pyspark.sql.types import StructType, StructField, StringType, IntegerType

```
schema = StructType([
    StructField("id", IntegerType(), True),
    StructField("name", StringType(), True)
])

df = (spark.readStream
    .format("cloudFiles")
    .option("cloudFiles.format", "csv")
    .option("cloudFiles.schemaLocation", "dbfs:/mnt/checkpoints/schema")
    .schema(schema)
    .load("s3://bucket/source/"))
```

Processing

- Use Auto Loader's streaming capabilities to process new files as they arrive.
- Write the stream to a Delta table for downstream processing:

```
query = (df.writeStream
    .format("delta")
    .option("checkpointLocation", "dbfs:/mnt/checkpoints/stream")
    .table("target_table"))
```

Scaling

- Auto Loader automatically detects new files and processes them incrementally, leveraging Spark's streaming engine.
- Use auto-scaling clusters to handle varying ingestion rates.
- Configure cloudFiles.maxFilesPerTrigger to control the number of files processed per micro-batch.

Benefits of Auto Loader over Standard Spark File Streaming

- Schema Inference and Evolution: Auto Loader automatically infers schemas and handles schema changes (e.g., new columns) without manual intervention, unlike standard Spark streaming, which requires predefined schemas.
- **File Discovery**: Uses cloud-native notifications (e.g., S3 Event Notifications, Azure Event Grid) for efficient file discovery, reducing latency and cost compared to Spark's directory listing.
- **Resilience**: Automatically handles file backlogs and failures, with robust checkpointing to avoid reprocessing.
- **Simplicity**: Simplifies configuration with options like cloudFiles.schemaLocation for schema management, reducing boilerplate code.
- **Performance**: Optimized for cloud storage, with better handling of large directories and incremental processing.

Scenario 77. How do you schedule and monitor Databricks Jobs with multiple tasks and dependencies?

Job Creation

- In the Databricks UI, navigate to the "Jobs" section and create a new job.
- Add tasks (e.g., notebooks, JARs, or Python scripts) and assign them to a cluster (job or interactive).
- Define timeout per task, retry limits, parameter passing, and cluster settings.
- Integrate email/slack/webhook alerts.

Task Dependencies

- Define dependencies in the Jobs UI by specifying which tasks depend on others (e.g., Transform depends on Extract).
- This creates a Directed Acyclic Graph (DAG) for task execution.

Scheduling

- Set a schedule using the UI or API with a cron expression (e.g., 0 0 2 * * ? for daily at 2 AM).
- Configure retries (e.g., 3 attempts) and timeouts for robustness.
- Example: Schedule a job to run an ETL pipeline every 6 hours.

Parameters

Pass parameters to tasks (e.g., date ranges) via the Jobs UI or API:

```
{
    "parameters": {"start_date": "2025-07-01", "end_date": "2025-07-10"}
}
```

Access parameters in notebooks using dbutils.widgets.get("start_date").

Monitoring Databricks Jobs

- Job Runs dashboard
- Run history
- Gantt charts
- Task logs
- Output artifacts

Programmatic Monitoring:

• Use the Databricks Jobs API to programmatically check job status:

curl -X GET https://<databricks-instance>/api/2.1/jobs/runs/get?run_id=<run_id> \ -H "Authorization: Bearer < token>"

• Integrate with monitoring tools like Azure Monitor or AWS CloudWatch.

Scenario 78. How does Delta Lake ensure ACID compliance and support concurrent read/write operations in a distributed environment?

Delta Lake ensures ACID (Atomicity, Consistency, Isolation, Durability) compliance on top of Spark by using a transaction log stored in a distributed file system (e.g., S3, ADLS). Here's how it works:

- Atomicity: Delta Lake uses a transaction log (stored as JSON files in _delta_log) to record all operations (e.g., inserts, updates, deletes). Transactions are committed atomically; if a write fails, no partial changes are applied. Writes are atomic via writeahead logs + file renames.
- **Consistency**: The transaction log enforces a consistent view of data. Readers see a snapshot of the data at a specific version, ensuring consistency even during concurrent writes. Uses MVCC (Multi-Version Concurrency Control).
- **Isolation**: Delta Lake implements optimistic concurrency control. Writers append new transaction log entries, and conflicts are resolved by checking the log version. If a conflict occurs (e.g., concurrent writes), Delta retries or fails the transaction. Every operation is serializable with snapshot isolation.
- **Durability**: Once a transaction is committed to the transaction log and data files are written to cloud storage, changes are durable, even in the event of failures.

Concurrent Read/Write Support

- Snapshot Isolation: Delta Lake uses versioning to provide snapshot isolation. Each
 read operation sees a consistent snapshot of the data based on the latest committed
 transaction log version.
- Optimistic Concurrency: Writers append to the transaction log without locking, and
 Delta resolves conflicts by ensuring no overlapping changes occur. For example, if two
 writers attempt to update the same partition, Delta detects the conflict and retries or
 fails one transaction.
- **Scalability**: The transaction log is stored in a distributed file system, allowing concurrent access by multiple Spark tasks. Delta Lake's log compaction (e.g., via OPTIMIZE) reduces log size for efficiency.
- **File Management**: Delta Lake manages data files (Parquet) and ensures that readers only access committed files, avoiding partial or uncommitted data.

Scenario 79. How do you perform upserts (merge operations) in Delta Lake using PySpark in Databricks?

Delta Lake's MERGE operation allows updating existing records or inserting new ones based on a condition. Here's how to perform an upsert using PySpark:

Setup

- Assume a Delta table (target_table) and a DataFrame (updates_df) with new/updated data.
- Example schema: id, name, updated_at.

Merge Operation

```
from delta.tables import DeltaTable
from pyspark.sql.functions import col

# Load the Delta table
delta_table = DeltaTable.forPath(spark, "dbfs:/mnt/delta/target_table")

# Perform the merge
(delta_table.alias("target")
.merge(
    updates_df.alias("source"),
    "target.id = source.id"
)
.whenMatchedUpdate(set={"name": "source.name", "updated_at": "source.updated_at"})
.whenNotMatchedInsert(values={"id": "source.id", "name": "source.name", "updated_at":
"source.updated_at"})
.execute())
```

Explanation

- **Condition**: target.id = source.id matches rows between the target Delta table and the source DataFrame.
- When Matched: Updates the name and updated_at columns for matching rows.
- When Not Matched: Inserts new rows from the source DataFrame.
- The operation is ACID-compliant, ensuring no data corruption during concurrent writes.

Best Practices

- Use partitioning on the merge key (e.g., id) to optimize performance.
- Enable OPTIMIZE and ZORDER on the Delta table to improve merge performance:

spark.sql("OPTIMIZE target_table ZORDER BY (id)")

 Handle schema evolution if the source DataFrame has new columns (see schema evolution question).

Scenario 80. How would you implement Slowly Changing Dimensions (SCD Type 1 and 2) using Delta Lake in Databricks?

SCD Type 1 (Overwrite):

Updates existing records with new values, overwriting historical data.

Implementation:

- Use a MERGE operation to update or insert records.
- Example for a customer table with customer_id, name, and email:

from delta.tables import DeltaTable

```
# Load target Delta table and source DataFrame
delta_table = DeltaTable.forPath(spark, "dbfs:/mnt/delta/customers")
source_df = spark.read.csv("s3://bucket/source/customers.csv")

# Perform SCD Type 1 merge
(delta_table.alias("target")
.merge(
    source_df.alias("source"),
    "target.customer_id = source.customer_id"
)
.whenMatchedUpdate(set={"name": "source.name", "email": "source.email"})
.whenNotMatchedInsert(values={"customer_id": "source.customer_id", "name": "source.name", "email": "source.name", "email": "source.email"})
.execute())
```

SCD Type 2 (Maintain History)

Maintains historical records by adding new rows with validity timestamps or flags.

Implementation:

- Add columns to the Delta table for tracking history (e.g., start_date, end_date, is_active).
- Example:

from delta.tables import DeltaTable from pyspark.sql.functions import current_timestamp, lit

```
# Load target Delta table and source DataFrame
delta_table = DeltaTable.forPath(spark, "dbfs:/mnt/delta/customers_scd2")
source_df = spark.read.csv("s3://bucket/source/customers.csv")
```

```
# Perform SCD Type 2 merge
(delta_table.alias("target")
.merge(
  source_df.alias("source"),
  "target.customer_id = source.customer_id AND target.is_active = true"
)
.whenMatchedUpdate(set={"is_active": lit(False), "end_date": current_timestamp()})
.whenNotMatchedInsert(values={
  "customer id": "source.customer id",
  "name": "source.name",
  "email": "source.email",
  "start_date": current_timestamp(),
  "end_date": lit(None),
  "is_active": lit(True)
})
.execute())
# Insert new records for updated customers
source_df_with_dates = source_df.withColumn("start_date", current_timestamp()) \
              .withColumn("end_date", lit(None)) \
              .withColumn("is_active", lit(True))
source_df_with_dates.write.format("delta").mode("append").save("dbfs:/mnt/delta/customers
_scd2")
```

• **Result**: Updates mark the old record as inactive (is_active = False, end_date set) and insert a new active record for the updated data.

Best Practices

- Partition the Delta table by a high-cardinality column (e.g., customer_id) to optimize merges.
- Use OPTIMIZE and ZORDER to improve query performance on SCD Type 2 tables.
- Vacuum old records periodically to manage storage (e.g., VACUUM customers_scd2 RETAIN 168 HOURS).

Scenario 81. How would you implement Slowly Changing Dimensions (SCD Type 2 vs Type 3)? Explain the differences and when to use each.

SCD Type 2

Maintains full history by creating a new row for each change, with validity timestamps or flags (e.g., start_date, end_date, is_active).

Implementation

• Using Delta Lake's MERGE to update existing records and append new ones.

When to Use:

- When full historical tracking is required (e.g., auditing, regulatory compliance).
- Example: Tracking customer address changes over time in a data warehouse.

Pros:

Preserves all historical data, supports time-based queries.

Cons

• Increases storage and complexity due to multiple rows per entity.

SCD Type 3

 Maintains limited history by adding a column to store the previous value of a specific attribute (e.g., previous_email).

Implementation

- Add a column to the Delta table for the previous value (e.g., previous_email).
- Example for a customer table:

```
from delta.tables import DeltaTable from pyspark.sql.functions import col
```

```
delta_table = DeltaTable.forPath(spark, "dbfs:/mnt/delta/customers_scd3")
source_df = spark.read.csv("s3://bucket/source/customers.csv")
```

```
# Perform SCD Type 3 merge
(delta_table.alias("target")
.merge(
    source_df.alias("source"),
    "target.customer_id = source.customer_id"
)
.whenMatchedUpdate(set={
    "previous_email": "target.email",
    "email": "source.email",
    "name": "source.name"
})
.whenNotMatchedInsert(values={
    "customer_id": "source.customer_id",
    "name": "source.name",
    "email": "source.email",
    "previous_email": lit(None)
```

```
})
.execute())
```

When to Use

- When only the current and previous value of a specific attribute need to be tracked.
- Example: Tracking the current and previous job title of an employee.

Pros

Simpler and less storage-intensive than SCD Type 2.

Cons

• Limited to tracking only one previous value, not full history.

Key Differences

- **History Tracking**: SCD Type 2 maintains full history with multiple rows, while SCD Type 3 tracks only the current and previous value in a single row.
- **Storage**: SCD Type 2 requires more storage due to multiple rows per entity, while SCD Type 3 is more compact.
- **Complexity**: SCD Type 2 is more complex to query and manage, while SCD Type 3 is simpler but less flexible.
- **Use Case**: Use SCD Type 2 for full historical analysis (e.g., auditing); use SCD Type 3 for lightweight tracking of a single attribute's history.

Scenario 82. What are the key advantages of using Delta Lake in Databricks over raw Parquet or CSV files?

ACID Transactions

 Delta Lake ensures ACID compliance via its transaction log, allowing reliable updates, deletes, and merges. Raw Parquet/CSV files lack transactional guarantees, leading to potential data corruption during concurrent writes.

Schema Enforcement and Evolution

- Delta Lake enforces schemas, preventing data quality issues from mismatched schemas. It also supports schema evolution (e.g., adding columns) during writes.
- Raw Parquet/CSV files require manual schema management, which can lead to errors.

Time Travel

- Delta Lake's versioning allows querying historical data or rolling back to previous versions (e.g., SELECT * FROM table VERSION AS OF 1).
- Raw Parquet/CSV files do not support versioning natively.

Upserts and Deletes

- Delta Lake supports MERGE for upserts and efficient deletes, enabling complex ETL workflows.
- Raw Parquet/CSV files require full table rewrites for updates or deletes, which is inefficient.

Performance Optimization

- Delta Lake supports OPTIMIZE for file compaction and ZORDER for data skipping, improving query performance.
- Raw Parquet/CSV files lack these optimizations, leading to slower queries on large datasets.

Concurrent Read/Write

- Delta Lake handles concurrent operations with snapshot isolation and optimistic concurrency control.
- Raw Parquet/CSV files can result in conflicts or partial reads during concurrent access.

Data Reliability

- Delta Lake's transaction log ensures data integrity, even in the case of failures.
- Raw Parquet/CSV files are prone to partial writes or corruption without transactional support.

Unified Batch and Streaming

- Delta Lake supports both batch and streaming workloads with the same table, simplifying pipelines.
- Raw Parquet/CSV files require separate handling for batch and streaming.