# Databricks Data Engineer Associate(By Vijay Bhaskar Reddy)
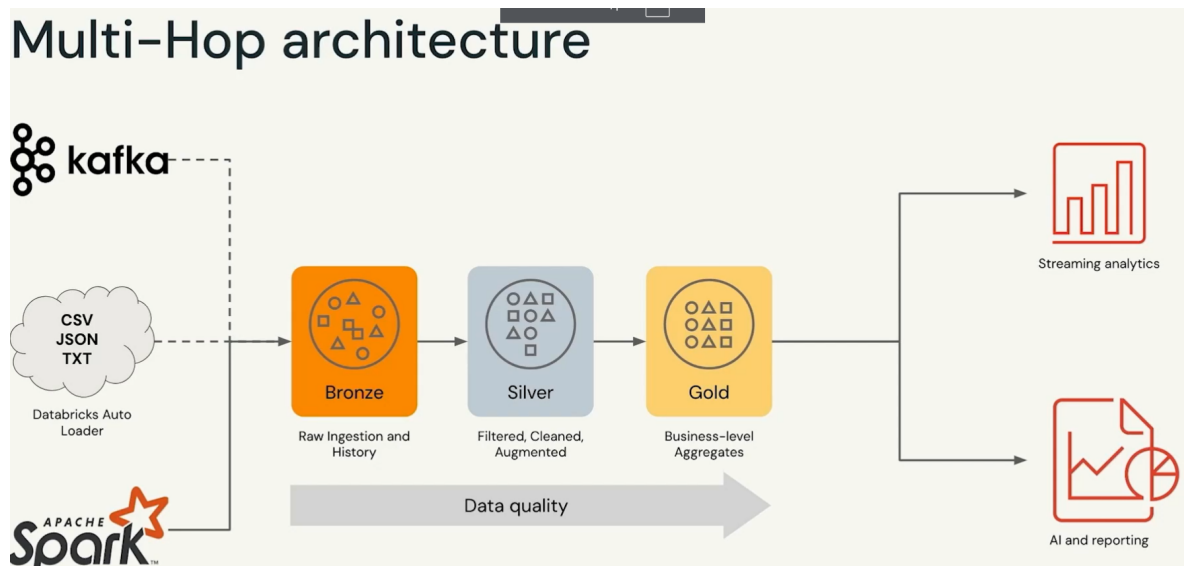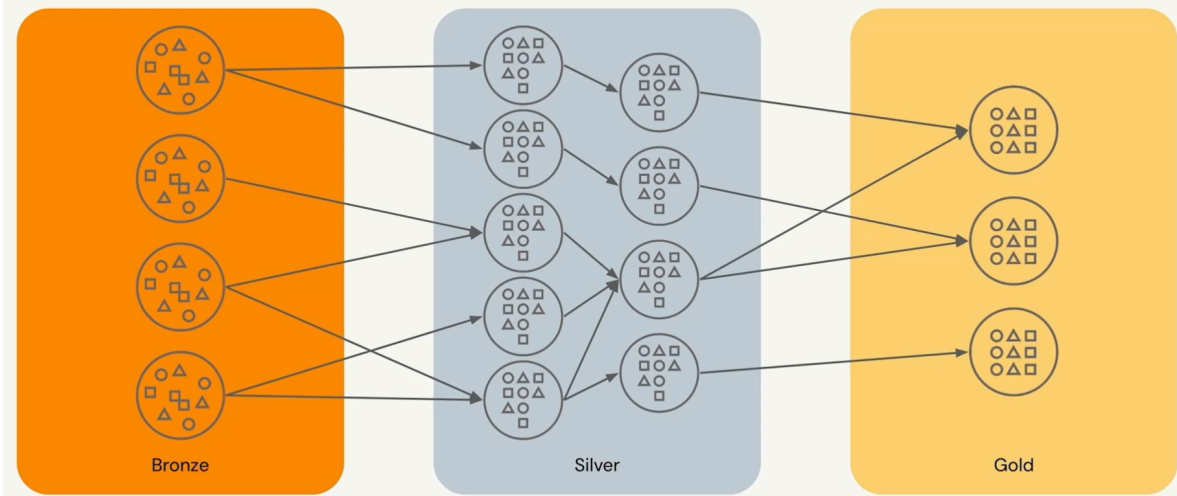
The content is from databricks exam syllabus:

https://www.databricks.com/sites/default/files/2024-05/databricks-certified-data-engineer-associate-exam-guide.pdf

1. Databricks Lakehouse Platform – 24%

## Section 1: Databricks Lakehouse Platform (24%)

- ▶ 1.Describe the relationship between the data lakehouse and the data warehouse.
- ▶ 2.Identify the improvement in data quality in the data lakehouse over the data lake.
- ▶ 3.Compare and contrast silver and gold tables, which workloads will use a bronze table as a source, which workloads will use a gold table as a source.
- ▶ 4.Identify elements of the Databricks Platform Architecture, such as what is located in the data plane versus the control plane and what resides in the customer's cloud account
- ▶ 5.Differentiate between all-purpose clusters and jobs clusters.
- ▶ 6.Identify how cluster software is versioned using the Databricks Runtime.
- ▶ 7.Identify how clusters can be filtered to view those that are accessible by the user
- ▶ 8.Describe how clusters are terminated and the impact of terminating a cluster.
- ▶ 9.Identify a scenario in which restarting the cluster will be useful.
- ▶ 10.Describe how to use multiple languages within the same notebook.
- ▶ 11.Identify how to run one notebook from within another notebook.
- ▶ 12.Identify how notebooks can be shared with others.
- ▶ 13.Describe how Databricks Repos enables CI/CD workflows in Databrick
- ▶ 14.Identify Git operations available via Databricks Repos.
- ▶ 15.Identify limitations in Databricks Notebooks version control functionality relative to Repos

2. ELT With Spark SQL and Python – 29%

## Section 2: ELT with Apache Spark(29%)

- ▶ 16.Extract data from a single file and from a directory of files
- ▶ 17.Identify the prefix included after the FROM keyword as the data type
- ▶ 18.Create a view, a temporary view, and a CTE as a reference to a file
- ▶ 19.Identify that tables from external sources are not Delta Lake tables.
- ▶ 20.Create a table from a JDBC connection and from an external CSV file
- ▶ 21.Identify how the count_if function and the count where x is null can be used
- ▶ 22.Identify how the count(row) skips NULL values.
- ▶ 23.Deduplicate rows from an existing Delta Lake table.
- ▶ 24.Create a new table from an existing table while removing duplicate rows.
- ▶ 25.Deduplicate a row based on specific column
- ▶ 26. Validate that the primary key is unique across all rows.
- ▶ 27.Validate that a field is associated with just one unique value in another field.
- ▶ 28.Validate that a value is not present in a specific field
- ▶ 29.Cast a column to a timestamp.
- ▶ 30.Extract calendar data from a timestamp.
- ▶ 31.Extract a specific pattern from an existing string column
- ▶ 32.Utilize the dot syntax to extract nested data field
- ▶ 33.Identify the benefits of using array function
- ▶ 34.Parse JSON strings into structs.
- ▶ 35.Identify which result will be returned based on a join query.
- ▶ 36.Identify a scenario to use the explode function versus the flatten function
- ▶ 37.Identify the PIVOT clause as a way to convert data from a long format to a wide format.
- ▶ 38.Define a SQL UDF
- ▶ 39.Identify the location of a function.
- ▶ 40.Describe the security model for sharing SQL UDFs.
- ▶ 41.Use CASE/WHEN in SQL code.
- ▶ 42.Leverage CASE/WHEN for custom control flow

3. Incremental Data Processing – 22%

# Section 3: Incremental Data Processing(22%)

▶ 43.Identify where Delta Lake provides ACID transactions

▶ 44.Identify the benefits of ACID transaction

▶ 45.Identify whether a transaction is ACID-compliant.

▶ 46.Compare and contrast data and metadata.

▶ 47.Compare and contrast managed and external tables.

▶ 48.Identify a scenario to use an external table.

▶ 49.Create a managed table.

▶ 50.Identify the location of a table.

▶ 51.Inspect the directory structure of Delta Lake file.

▶ 52.Identify who has written previous versions of a table. || 53.Review a history of table transactions.

▶ 54.Roll back a table to a previous version.

▶ 55.Identify that a table can be rolled back to a previous version.| 56.Query a specific version of a table

▶ 57.Identify why Z-ordering is beneficial to Delta Lake table

▶ 58.Identify how vacuum commits deletes.

▶ 59.Identify the kind of files Optimize compact.

▶ 60.Identify CTAS as a solution.

▶ 61.Create a generated column.

▶ 62.Add a table comment.

▶ 63.Use CREATE OR REPLACE TABLE and INSERT OVERWRITE

▶ 64.Compare and contrast CREATE OR REPLACE TABLE and INSERT OVERWRITE

▶ 65.Identify a scenario in which MERGE should be used.

▶ 66.Identify MERGE as a command to deduplicate data upon writing.

▶ 67.Describe the benefits of the MERGE command

▶ 68.Identify why a COPY INTO statement is not duplicating data in the target table.

▶ 69.Identify a scenario in which COPY INTO should be used.

▶ 70.Use COPY INTO to insert data.

> 71.Identify the components necessary to create a new DLT pipeline.
>
> 72.Identify the purpose of the target and of the notebook libraries in creating a pipeline.
>
> 73.Compare and contrast triggered and continuous pipelines in terms of cost and latency
>
> 74.Identify which source location is utilizing Auto Loader.
>
> 75.Identify a scenario in which Auto Loader is beneficial
>
> 76.Identify why Auto Loader has inferred all data to be STRING from a JSON source
>
> 77.Identify the default behavior of a constraint violation
>
> 78.Identify the impact of ON VIOLATION DROP ROW and ON VIOLATION FAIL UPDATE for a constraint violation
>
> 79.Explain change data capture and the behavior of APPLY CHANGES INTO
>
> 80.Query the events log to get metrics, perform audit log in, examine lineage.
>
> 81.Troubleshoot DLT syntax: Identify which notebook in a DLT pipeline produced an error, identify the need for LIVE in create statement, identify the need for STREAM in from clause.

4. Production Pipelines – 16%

# Section 4: Production Pipelines(16%)

> 82.Identify the benefits of using multiple tasks in Job
>
> 83.Set up a predecessor task in Jobs.
>
> 84.Identify a scenario in which a predecessor task should be set up.
>
> 85.Review a task's execution history.
>
> 86.Identify CRON as a scheduling opportunity.
>
> 87.Debug a failed task.
>
> 88.Set up a retry policy in case of failure.
>
> 89.Create an alert in the case of a failed task.
>
> 90.Identify that an alert can be sent via email.

5. Data Governance – 9%

## Section 5: Data Governance(9%)

- ▸ 91.Identify one of the four areas of data governance.
- ▸ 92.Compare and contrast meta stores and catalogs.
- ▸ 93.Identify Unity Catalog securable.
- ▸ 94.Define a service principal
- ▸ 95.Identify the cluster security modes compatible with Unity Catalog.
- ▸ 96.Create a UC-enabled all-purpose cluster.
- ▸ 97.Create a DBSQL warehouse.
- ▸ 98.Identify how to query a three-layer namespace.
- ▸ 99.Implement data object access control
- ▸ 100.Identify colocating meta stores with a workspace as best practice.
- ▸ 101.Identify using service principals for connections as best practice.
- ▸ 102.Identify the segregation of business units across catalog as best practice.

# Section 1: Databricks Lakehouse Platform (24%)

▼ **1.Describe the relationship between the data lakehouse and the data warehouse.**

- The relationship between a **data lakehouse** and a **data warehouse** is quite complementary, as they each serve distinct but interconnected roles in data management and analytics

### Data Warehouse

- ○ **Purpose**: Designed for high-performance analytics and real-time insights.
- ○ **Data Type**: Stores processed, structured data.
- ○ **Schema**: Uses a "schema-on-write" approach, meaning data is structured before being stored.
- ○ **Users**: Typically used by business analysts and professionals who need to generate reports and perform business intelligence tasks.
- ○ **Governance**: Strong data governance and rigorous structure make it easier to maintain.

### Data Lakehouse

- ○ **Purpose**: Combines the flexible storage of a data lake with the high-performance analytics of a data warehouse.
- ○ **Data Type**: Can store data in any format—structured, semi-structured, or unstructured.

- **Schema**: Uses a "schema-on-read" approach, allowing for more flexibility as data is structured when read.
- **Users**: Ideal for data scientists and engineers working on machine learning, AI, and data science projects.
- **Governance**: Requires robust governance to prevent data chaos, but offers the flexibility to handle diverse data types.

## Relationship

- **Integration**: Data lakehouses can integrate with data warehouses to provide a unified data management solution. Raw data can be stored in the lakehouse and then processed and moved to the warehouse for specific analytics tasks.
- **Flexibility and Performance**: While data warehouses are optimized for performance with structured data, data lakehouses offer the flexibility to handle various data types and support advanced analytics and machine learning

- data lakehouse can be seen as an evolution that bridges the gap between the raw data storage capabilities of a data lake and the structured, high-performance analytics of a data warehouse

▼ **2.Identify the improvement in data quality in the data lakehouse over the data lake.**

- The data lakehouse architecture offers several improvements in data quality over traditional data lakes:

  1. **Structured Data Management**: Data lakehouses **integrate** the structured data management capabilities of data warehouses, ensuring that **data is organized and easily accessible**. This contrasts with data lakes, which often store raw, unstructured data that can be difficult to manage and query.

  2. **Data Governance and Compliance**: Data lakehouses provide robust data governance frameworks, which help maintain data quality by enforcing policies and standards. This is **crucial** for compliance with regulations and for **ensuring data integrity**.

  3. **Schema Enforcement**: While data lakes use a "schema-on-read" approach, data lakehouses can enforce schemas at both write and read times. This ensures that data conforms to predefined structures, reducing errors and inconsistencies.

  4. **Data Lineage and Auditing**: Data lakehouses offer better tools for tracking data lineage and auditing changes. This transparency helps in identifying and correcting data quality issues more efficiently.

  5. **Integrated Data Quality Tools**: Many data lakehouse platforms come with built-in data quality tools that automate the detection and correction of data anomalies. These tools help maintain high data quality standards across the data lifecycle.

- Overall, the data lakehouse architecture enhances data quality by combining the flexibility of data lakes with the **structured, governed** approach of **data warehouses**. This hybrid

model ensures that data is not only **stored efficiently** but also remains **accurate, consistent, and reliable.**

▼ **3.Compare and contrast silver and gold tables, which workloads will use a bronze table as a source, which workloads will use a gold table as a source.**

**Comparison of Bronze, Silver, and Gold Tables**

**Bronze Tables**:

- **Purpose**: Store raw, unprocessed data directly from various sources.
- **Data Quality**: Data is in its original format, with minimal to no transformation.
- **Usage**: Initial data ingestion and archiving.
- **Users**: Data engineers and data scientists for preliminary data exploration and processing.

**Silver Tables**:

- **Purpose**: Store cleansed and conformed data.
- **Data Quality**: Data is validated, deduplicated, and transformed to a certain extent.
- **Usage**: Intermediate layer for further processing and analytics.
- **Users**: Data engineers, data analysts, and data scientists for more refined datasets.

**Gold Tables**:

- **Purpose**: Store enriched and aggregated data.
- **Data Quality**: Data is highly processed, often aggregated, and ready for high-performance analytics.
- **Usage**: Final layer for business intelligence, reporting, and advanced analytics.
- **Users**: Business analysts, BI developers, executives, and decision-makers for generating insights and making data-driven decisions.

## Workloads Using Bronze Tables as a Source

- **Data Ingestion**: Initial data loading from external sources like cloud storage, Kafka, or databases.
- **Data Archiving**: Maintaining historical records and raw data for compliance and audit purposes.
- **Preliminary Data Processing**: Basic transformations and validations before moving data to the silver layer.

## Workloads Using Gold Tables as a Source

- **Business Intelligence (BI)**: Generating dashboards, reports, and visualizations for business insights.

- **Advanced Analytics**: Performing complex queries, predictive analytics, and machine learning model training.
- **Operational Reporting**: Real-time and batch reporting for operational decision-making.

By understanding the roles and uses of these different table layers, organizations can effectively manage their data lifecycle, ensuring data quality and optimizing performance for various analytical workloads.

▼ **4.Identify elements of the Databricks Platform Architecture, such as what is located in the data plane versus the control plane and what resides in the customer's cloud account**

- The Databricks Platform Architecture is divided into two main components: the **control plane** and the **data plane**. Here's a detailed breakdown of each

## Control Plane

The control plane includes the backend services that Databricks manages within its own account. These services are responsible for:

- **Authentication and Authorization**: Managing user access and security.
- **Job Scheduling**: Orchestrating and scheduling jobs.
- **Web Application**: The Databricks workspace interface.
- **REST APIs**: Providing programmatic access to Databricks features.

## Data Plane

The data plane is where your data is processed and resides within your cloud account. It includes:

- **Compute Resources**: Virtual machines (VMs) or clusters that run your data processing tasks.
- **Storage**: Data stored in your cloud storage (e.g., AWS S3, Azure Blob Storage, Google Cloud Storage).

## Customer's Cloud Account

In the customer's cloud account, the following elements are present:

- **Compute Resources**: Managed within the customer's virtual network (VNet) or Virtual Private Cloud (VPC).
- **Storage**: Data stored in the customer's cloud storage buckets.
- **Networking**: Network configurations and security groups that control access to the compute resources.

Here's a simplified diagram to illustrate the architecture:

```
+--------------------+        +--------------------+
|   Control Plane    |        |     Data Plane     |
|--------------------|        |--------------------|
| - Authentication   |        | - Compute Resources|
| - Job Scheduling   |        | - Storage          |
| - Web Application  |        | - Networking       |
| - REST APIs        |        |                    |
+--------------------+        +--------------------+
          |                            |
          |                            |
          v                            v
+---------------------------------------------------+
|               Customer's Cloud Account            |
|---------------------------------------------------|
| - Compute Resources (VMs, Clusters)               |
| - Storage (S3, Blob Storage, GCS)                 |
| - Networking (VNet, VPC, Security Groups)         |
+---------------------------------------------------+
```

- while Databricks manages the control plane, the data plane remains within the customer's cloud account, providing better control over data and resources

▼ **5.Differentiate between all-purpose clusters and jobs clusters.**

Certainly! Here's a comparison between **all-purpose clusters** and **jobs clusters** in Databricks:

## All-Purpose Clusters

- **Use Case**: Ideal for interactive data analysis, development, and collaboration.
- **Lifespan**: These clusters can be manually started and stopped. They remain active until manually terminated or until they auto-terminate due to inactivity.
- **Access**: Multiple users can share these clusters, making them suitable for collaborative work.
- **Cost**: Typically more expensive as they are designed to be always available for interactive use.
- **Configuration**: Can be configured through the UI, CLI, or REST API.

## Jobs Clusters

- **Use Case**: Designed for running automated jobs and workflows.

- **Lifespan**: These clusters are ephemeral; they are created when a job starts and terminated automatically once the job completes.
- **Access**: Dedicated to a single job run, ensuring isolation and resource efficiency.
- **Cost**: Generally more cost-effective as they only run for the duration of the job.
- **Configuration**: Created and managed by the Databricks job scheduler.

▼ **6.Identify how cluster software is versioned using the Databricks Runtime.**

Databricks Runtime versions the cluster software to ensure compatibility, performance, and security. Here's how it works:

## Databricks Runtime Versions

- **Versioning**: Each Databricks Runtime version is identified by a unique version number (e.g., 10.4 LTS, 11.3, 12.2).
- **Variants**: There are different variants for specific use cases, such as Databricks Runtime for Machine Learning.
- **Release Notes**: Each version comes with detailed release notes that include new features, improvements, and bug fixes.

## Long-Term Support (LTS)

- **LTS Versions**: These versions are supported for an extended period, typically three years, providing stability for long-term projects.
- **Support Lifecycle**: Regular updates and maintenance releases are provided to ensure continued reliability and security

▼ **7.Identify how clusters can be filtered to view those that are accessible by the user**

To filter clusters in Databricks to view those that are accessible by a specific user, you can use the following methods:

## Using the Databricks UI

1. **Navigate to the Clusters Page**: Go to the Clusters section in the Databricks workspace.
2. **Filter by Permissions**: Use the filter options to display clusters based on your access permissions. You can filter clusters where you have specific permissions like "Can Attach To," "Can Restart," or "Can Manage."

▼ **8.Describe how clusters are terminated and the impact of terminating a cluster.**

## Cluster Termination in Databricks

**Cluster termination** can occur either manually or automatically. Here's how it works and its impact:

## How Clusters are Terminated

1. **Manual Termination**:

   - **User-Initiated**: Users can manually terminate a cluster via the Databricks UI, CLI, or API.

   - **Immediate Effect**: The cluster stops running immediately, and all associated resources are released.

2. **Automatic Termination**:

   - **Auto-Termination Settings**: Clusters can be configured to auto-terminate after a specified period of inactivity.

   - **Job Completion**: Job clusters automatically terminate once the job completes.

## Impact of Terminating a Cluster

1. **Resource Release**:

   - **Compute Resources**: All virtual machines (VMs) and other compute resources are deallocated, stopping any ongoing processes.

   - **Cost Savings**: Reduces costs by stopping charges for compute resources and Databricks Units (DBUs).

2. **Data and State**:

   - **Data Persistence**: Data stored in external storage (e.g., S3, Azure Blob Storage) remains unaffected.

   - **Ephemeral Data**: Any data stored on the cluster's local storage is lost.

3. **Running Processes**:

   - **Job Interruption**: Any running jobs or processes are terminated, which can lead to incomplete tasks.

   - **Session Termination**: Active user sessions are closed, and any unsaved work is lost.

▼ **9.Identify a scenario in which restarting the cluster will be useful.**

Restarting a cluster can be particularly useful in several scenarios. Here's one common example:

## Scenario: Resolving Performance Issues

**Context**: Over time, a cluster might experience performance degradation due to memory leaks, resource contention, or other issues.

**Action**: Restarting the cluster can help clear out any accumulated issues, refresh the environment, and restore optimal performance.

**Impact**:

- **Resource Cleanup**: Frees up memory and other resources that might be tied up by long-running processes.

- **Configuration Updates**: Applies any configuration changes or updates that were made but not yet applied.

- **Job Stability**: Improves the stability and performance of ongoing and future jobs.

## Example

Imagine you have a cluster running multiple Spark jobs, and you notice that the jobs are taking longer to complete than usual. By restarting the cluster, you can:

- Clear any residual data or processes that might be causing slowdowns.

- Ensure that the cluster is running with the latest configurations and updates.

- Provide a fresh environment for your jobs, potentially improving their execution time and reliability.

▼ **10.Describe how to use multiple languages within the same notebook.**

In Databricks notebooks, you can seamlessly use multiple languages such as Python, SQL, Scala, and R within the same notebook. Here's how you can do it:

## Using Multiple Languages in Databricks Notebooks

1. **Language Magic Commands**:

   - Each cell in a Databricks notebook can be designated to a specific language using magic commands.

   - The magic command for each language is:

     - **Python**: `%python`

     - **SQL**: `%sql`

     - **Scala**: `%scala`

     - **R**: `%r`

▼ **11.Identify how to run one notebook from within another notebook.**

   - **Using the** `%run` **Command**:

     - The `%run` command allows you to include another notebook within a notebook. This is useful for modularizing your code.

     - Example:

```
%run "./path/to/your/notebook"
```

- This command runs the entire notebook inline, and any functions or variables defined in the called notebook become available in the calling notebook.

- **Using** `dbutils.notebook.run()`:

  - The `dbutils.notebook.run()` function is part of the Databricks Utilities and allows you to run a notebook as a separate job.

  - Example:

```
result = dbutils.notebook.run("/path/to/your/notebook", timeout_seco
nds=60, arguments={"arg1": "value1"})
```

  - This method is useful for passing parameters to the notebook and handling return values. It also allows for more complex workflows and dependencie

▼ **12.Identify how notebooks can be shared with others.**

Sharing notebooks in Databricks is straightforward and can be done in a few different ways. Here are the steps:

1. **Share a Notebook**:

   - Click the **Share** button at the top of the notebook.

   - In the Sharing dialog, you can select who to share the notebook with and what level of access they have. You can choose from the following permission levels:
     - **Can Read**: Allows the user to view the notebook.
     - **Can Run**: Allows the user to run the notebook.
     - **Can Edit**: Allows the user to edit the notebook.
     - **Can Manage**: Allows the user to manage permissions and settings for the notebook.

2. **Manage Notebook Permissions with Folders**:

   - You can manage notebook permissions by adding the notebook to folders. Notebooks in a folder inherit all permission settings of that folder. For example, if a user has **Can Run** permission on a folder, they will have **Can Run** permission on all notebooks within that folder.

3. **Code Comments**:

   - You can have discussions with collaborators using the comments feature. To add a comment to code, highlight the code section and click the comment bubble icon. You

can then add your comment and click **Comment**. This is useful for collaborative development and review

▼ **13.Describe how Databricks Repos enables CI/CD workflows in Databrick**

Databricks Repos significantly enhances CI/CD (Continuous Integration/Continuous Deployment) workflows by integrating with Git repositories and enabling seamless collaboration and automation. Here's how it works:

1. **Version Control Integration**:

   - Databricks Repos allows you to clone Git repositories into your Databricks workspace. This integration ensures that your notebooks, libraries, and other files are version-controlled, enabling collaborative development and tracking of changes.

2. **Branching and Merging**:

   - You can create branches for different features or tasks, work on them independently, and merge them back into the main branch when ready. This helps in managing parallel development and integrating changes smoothly.

3. **Automated Testing and Deployment**:

   - By integrating with CI/CD tools like GitHub Actions, Azure DevOps, or Jenkins, you can automate the testing and deployment of your Databricks workflows. For example, you can set up pipelines that automatically run tests on your notebooks and deploy them to production environments upon successful completion.

4. **Databricks Asset Bundles (DABs)**:

   - Databricks Asset Bundles enable programmatic management of Databricks workflows. You can define your workflows in YAML files, validate them, and deploy them using the Databricks CLI. This approach ensures consistency and repeatability in your deployment processes.

5. **Collaboration and Code Review**:

   - With Git integration, team members can collaborate on the same project, review each other's code, and provide feedback through pull requests. This fosters a collaborative environment and ensures code quality.

By leveraging these features, Databricks Repos streamlines the development, testing, and deployment processes, making it easier to maintain high-quality code and accelerate the delivery of data engineering and data science projects.

▼ **14.Identify Git operations available via Databricks Repos.**

Databricks Repos supports a variety of Git operations, enabling seamless integration with your Git workflows. Here are some of the key Git operations you can perform:

1. **Clone a Repository**:

- You can clone a remote Git repository into your Databricks workspace. This allows you to work with the repository's contents directly within Databricks.

2. **Branch Management**:

- **Create a Branch**: You can create new branches for feature development or other tasks.

- **Switch Branches**: Easily switch between different branches to work on various parts of your project.

- **Merge Branches**: Merge changes from one branch into another, facilitating collaborative development.

3. **Commit and Push Changes**:

- **Commit Changes**: Save your changes to the local repository with a commit message.

- **Push Changes**: Push your committed changes to the remote repository, ensuring your work is backed up and shared with your team.

4. **Pull Changes**:

- **Pull Changes**: Fetch and integrate changes from the remote repository into your local branch. This keeps your local repository up-to-date with the latest changes from your team.

5. **Rebase and Resolve Conflicts**:

- **Rebase**: Rebase your branch on top of another branch to integrate changes smoothly.

- **Resolve Conflicts**: Handle merge conflicts that arise during rebasing or merging, ensuring a clean integration of changes.

6. **View Diffs**:

- **Visual Comparison**: Compare differences between commits, branches, or the working directory to understand changes and review code.

These operations make it easy to manage your code and collaborate with others directly within Databricks

▼ **15.Identify limitations in Databricks Notebooks version control functionality relative to Repos**

Databricks Notebooks have some limitations in version control functionality compared to Databricks Repos. Here are the key differences:

1. **Granularity of Version Control**:

- **Notebooks**: Version control is limited to individual notebooks. You can track changes and revert to previous versions, but this is done on a per-notebook basis.

- **Repos**: Allows for version control at the repository level, enabling you to manage multiple notebooks and other files together. This supports more complex project structures and workflows.

2. **Branching and Merging**:

- **Notebooks**: Do not support branching and merging. You cannot create branches for different features or merge changes from different branches.

- **Repos**: Fully supports Git branching and merging, allowing for parallel development and integration of changes from multiple branches.

3. **Collaboration**:

- **Notebooks**: Collaboration is limited to sharing individual notebooks. Multiple users can edit a notebook, but there is no built-in mechanism for handling conflicts or reviewing changes.

- **Repos**: Facilitates collaboration through Git workflows, including pull requests, code reviews, and conflict resolution.

4. **Integration with CI/CD**:

- **Notebooks**: Limited integration with CI/CD pipelines. You can manually export notebooks and integrate them into CI/CD workflows, but this process is not seamless.

- **Repos**: Designed to integrate with CI/CD tools, enabling automated testing, deployment, and continuous integration directly from the repository.

5. **File Types and Structure**:

- **Notebooks**: Primarily focused on notebook files. Managing other file types (e.g., scripts, data files) within the same project can be cumbersome.

- **Repos**: Supports a variety of file types and allows you to organize them within a single repository, making it easier to manage complex projects.

These limitations highlight the advantages of using Databricks Repos for more advanced version control and collaboration needs.

# Section 2: ELT with Apache Spark(29%)

▼ **16.Extract data from a single file and from a directory of files**

To extract data from a single file or a directory of files in Databricks, you can use SQL or PySpark. Here are some examples to help you get started:

### Extracting Data from a Single File

- You can use SQL to query a single file directly. For example, if you have a JSON file, you can use:

```sql
SELECT * FROM json.`/path/to/your/file.json`
```

Make sure to use backticks around the file path.

### Extracting Data from a Directory of Files

- If you want to query all files in a directory, you can specify the directory path. For example, for JSON files

```sql
SELECT * FROM json.`/path/to/your/directory/`
```

### Using PySpark

- You can also use PySpark to read data from files. Here's how you can do it:

### Single File

```python
df = spark.read.format("json").load("/path/to/your/file.json")
df.show()
```

### Directory of Files

```python
df = spark.read.format("json").load("/path/to/your/directory/")
df.show()
```

These methods allow you to easily extract and work with data stored in various file formats on Databricks

▼ **17.Identify the prefix included after the FROM keyword as the data type**

The prefix included after the `FROM` keyword indicates the data type or format of the files being queried. Here are the prefixes and their corresponding data types:

- `json` : This prefix indicates that the data is in JSON format.

- `csv` : This prefix would indicate that the data is in CSV format.

- `parquet` : This prefix would indicate that the data is in Parquet format.

For example:

```
SELECT * FROM json.`/path/to/your/file.json`
```

In this SQL query, `json` is the prefix that specifies the data type of the file being queried.

▼ **18.Create a view, a temporary view, and a CTE as a reference to a file**

## Creating a View

A view is a saved query that you can reference like a table. Here's how to create a view:

```
CREATE VIEW my_view AS
SELECT * FROM json.`/path/to/your/file.json`;
```

## Creating a Temporary View

A temporary view is similar to a regular view but is only available within the current session. Here's how to create a temporary view:

```
CREATE OR REPLACE TEMP VIEW my_temp_view AS
SELECT * FROM json.`/path/to/your/file.json`;
```

## Creating a Common Table Expression (CTE)

A CTE is a temporary result set that you can reference within a `SELECT` , `INSERT` , `UPDATE` , or `DELETE` statement. Here's how to create a CTE:

```
WITH my_cte AS (
  SELECT * FROM json.`/path/to/your/file.json`
)
SELECT * FROM my_cte;
```

These examples use JSON files, but you can replace `json` with other formats like `csv` or `parquet` depending on your data type

▼ **19.Identify that tables from external sources are not Delta Lake tables.**

- To identify that tables from external sources are not Delta Lake tables in Databricks, you can look for specific characteristics that distinguish Delta Lake tables from other types of

tables. Here are some key points:

1. **Metadata and Log Files**: Delta Lake tables have a `_delta_log` directory that stores transaction logs. If this directory is absent, the table is not a Delta Lake table.

2. **Table Description**: You can use the `DESCRIBE FORMATTED` command to check the table's metadata. Delta Lake tables will have specific metadata entries related to Delta Lake.

3. **Programmatic Check**: You can use a Python function in a Databricks notebook to programmatically determine if a table is a Delta table. Here's an example:

```python
def is_delta_table(table_name: str) -> bool:
    desc_table = spark.sql(f"DESCRIBE FORMATTED {table_name}").colle
ct()
    location = [row[1] for row in desc_table if row[0] == 'Locatio
n'][0]
    try:
        dbutils.fs.ls(f"{location}/_delta_log")
        return True
    except Exception:
        return False


result = is_delta_table("your_table_name")
if result:
    print("Yes, it is a Delta table.")
else:
    print("No, it is not a Delta table.")
```

4. **Performance Guarantees**: Tables from external sources do not provide the same performance guarantees as Delta Lake tables. Delta Lake tables offer ACID transactions, scalable metadata handling, and unification of streaming and batch data processing

▼ **20.Create a table from a JDBC connection and from an external CSV file**

- you can create a table from a JDBC connection and from an external CSV file in Databricks:

## Creating a Table from a JDBC Connection

You can use the following steps to create a table from a JDBC connection:

1. **Set up the JDBC connection**:

```
jdbc_url = "jdbc:your_database_url"
jdbc_table = "your_table_name"
jdbc_properties = {
    "user": "your_username",
    "password": "your_password",
    "driver": "your_jdbc_driver"
}
```

2. **Read data from the JDBC source**:

```
df = spark.read.jdbc(url=jdbc_url, table=jdbc_table, properties=j
dbc_properties)
```

3. **Create a table from the DataFrame**:

```
df.write.saveAsTable("your_databricks_table_name")
```

## Creating a Table from an External CSV File

You can use the following steps to create a table from an external CSV file:

1. **Read the CSV file into a DataFrame**:

```
df = spark.read.format("csv").option("header", "true").load("/pat
h/to/your/file.csv")
```

2. **Create a table from the DataFrame**:

```
df.write.saveAsTable("your_databricks_table_name")
```

Alternatively, you can use SQL to create a table directly from a CSV file:

```
CREATE TABLE your_table_name
USING csv
OPTIONS (path "/path/to/your/file.csv", header "true");
```

- These methods will help you create tables from both JDBC connections and external CSV files in databricks

▼ **21.Identify how the count_if function and the count where x is null can be used**

- you can use the `count_if` function and the `count` function with a condition where a column is `NULL` in Databricks:

## Using `count_if` Function

The `count_if` function counts the number of true values for a given Boolean expression. Here are some examples:

1. **Count even numbers in a column**:

   ```
   SELECT count_if(col % 2 = 0) FROM your_table;
   ```

2. **Count non-null values in a column**:

   ```
   SELECT count_if(col IS NOT NULL) FROM your_table;
   ```

3. **Count distinct even numbers in a column**:

   ```
   SELECT count_if(DISTINCT col % 2 = 0) FROM your_table;
   ```

## Using `count` with a Condition where `x` is `NULL`

To count the number of rows where a specific column is `NULL`, you can use the `count` function with a `CASE` statement or a `FILTER` clause:

1. **Using `CASE` statement**:

   ```
   SELECT count(CASE WHEN col IS NULL THEN 1 END) FROM your_table;
   ```

2. **Using `FILTER` clause**:

   ```
   SELECT count(*) FILTER (WHERE col IS NULL) FROM your_table;
   ```

These methods allow you to count specific conditions within your data effectively

▼ **22.Identify how the count(row) skips NULL values.**

The `count` function in SQL, including Databricks SQL, skips `NULL` values when counting rows for a specific column. Here's how it works:

## How `count` Skips `NULL` Values

1. **Counting Non-NULL Values**:

   - When you use `count(column_name)`, it counts only the rows where `column_name` is not `NULL`.
   - Example:This query will count all rows where `col` is not `NULL`.

```
SELECT count(col) FROM your_table;
```

2. **Counting All Rows**:

- When you use `count(*)` or `count(1)`, it counts all rows, including those with `NULL` values in any column.

- Example:This query will count all rows in the table, regardless of `NULL` values.

```
SELECT count(*) FROM your_table;
```

## Example

Consider a table with the following data:

| id | value |
|----|-------|
| 1  | 10    |
| 2  | NULL  |
| 3  | 20    |
| 4  | NULL  |

- Using `count(value)`:

```
SELECT count(value) FROM your_table;
```

This will return `2` because there are two non-`NULL` values in the `value` column.

- Using `count(*)`:

```
SELECT count(*) FROM your_table;
```

This will return `4` because it counts all rows, including those with `NULL` values in the `value` column

▼ **23.Deduplicate rows from an existing Delta Lake table.**

To deduplicate rows from an existing Delta Lake table in Databricks, you can use PySpark or SQL. Here are the steps for both methods

## Using PySpark

1. **Load the Delta Table into a DataFrame**:

```
df = spark.table("your_table_name")
```

2. **Drop Duplicate Rows**:

```
df_deduped = df.dropDuplicates()
```

3. **Overwrite the Original Table with the Deduplicated DataFrame**:

```
df_deduped.write.format("delta").mode("overwrite").saveAsTable("your
_table_name")
```

## Using SQL

1. **Create a Temporary View with Distinct Rows**:

```
CREATE OR REPLACE TEMP VIEW deduped_view AS
SELECT DISTINCT * FROM your_table_name;
```

2. **Overwrite the Existing Delta Table with the Distinct Rows**:

```
INSERT OVERWRITE TABLE your_table_name
SELECT * FROM deduped_view;
```

## Example

Here's a complete example using PySpark:

```
# Load the Delta table
df = spark.table("your_table_name")

# Drop duplicates
df_deduped = df.dropDuplicates()

# Overwrite the original table
df_deduped.write.format("delta").mode("overwrite").saveAsTable("your_ta
ble_name")
```

And using SQL:

```
-- Create a temporary view with distinct rows
CREATE OR REPLACE TEMP VIEW deduped_view AS
SELECT DISTINCT * FROM your_table_name;

-- Overwrite the existing Delta table with the distinct rows
INSERT OVERWRITE TABLE your_table_name
SELECT * FROM deduped_view;
```

These methods will help you remove duplicate rows from your Delta Lake table efficiently

▼ **24.Create a new table from an existing table while removing duplicate rows.**

To create a new table from an existing table while removing duplicate rows in Databricks, you can use either PySpark or SQL. Here are the steps for both methods:

## Using PySpark

1. **Load the Existing Table into a DataFrame**:

```
df = spark.table("existing_table_name")
```

2. **Drop Duplicate Rows**:

```
df_deduped = df.dropDuplicates()
```

3. **Create a New Table from the Deduplicated DataFrame**:

```
df_deduped.write.format("delta").saveAsTable("new_table_name")
```

## Using SQL

1. **Create a Temporary View with Distinct Rows**:

```
CREATE OR REPLACE TEMP VIEW deduped_view AS
SELECT DISTINCT * FROM existing_table_name;
```

2. **Create a New Table from the Temporary View**:

```
CREATE TABLE new_table_name AS
SELECT * FROM deduped_view;
```

## Example

Here's a complete example using PySpark:

```
# Load the existing table
df = spark.table("existing_table_name")

# Drop duplicates
df_deduped = df.dropDuplicates()
```

```
# Create a new table from the deduplicated DataFrame
df_deduped.write.format("delta").saveAsTable("new_table_name")
```

And using SQL:

```
-- Create a temporary view with distinct rows
CREATE OR REPLACE TEMP VIEW deduped_view AS
SELECT DISTINCT * FROM existing_table_name;

-- Create a new table from the temporary view
CREATE TABLE new_table_name AS
SELECT * FROM deduped_view;
```

These methods will help you create a new table from an existing one while ensuring that duplicate rows are removed

▼ **25.Deduplicate a row based on specific column**

To deduplicate rows based on a specific column in Databricks, you can use PySpark or SQL. Here's how you can do it:

## Using PySpark

1. **Load the Existing Table into a DataFrame**:

```
df = spark.table("existing_table_name")
```

2. **Drop Duplicate Rows Based on a Specific Column**:

```
df_deduped = df.dropDuplicates(["specific_column"])
```

3. **Create a New Table or Overwrite the Existing Table**:

```
df_deduped.write.format("delta").mode("overwrite").saveAsTable("exis
ting_table_name")
```

## Using SQL

1. **Create a Temporary View with Distinct Rows Based on a Specific Column**:

```
CREATE OR REPLACE TEMP VIEW deduped_view AS
SELECT DISTINCT * FROM existing_table_name
WHERE specific_column IS NOT NULL;
```

2. **Overwrite the Existing Table with the Distinct Rows**:

```
INSERT OVERWRITE TABLE existing_table_name
SELECT * FROM deduped_view;
```

## Example

Here's a complete example using PySpark:

```python
# Load the existing table
df = spark.table("existing_table_name")

# Drop duplicates based on a specific column
df_deduped = df.dropDuplicates(["specific_column"])

# Overwrite the original table with the deduplicated DataFrame
df_deduped.write.format("delta").mode("overwrite").saveAsTable("existing_table_name")
```

And using SQL:

```sql
-- Create a temporary view with distinct rows based on a specific column
CREATE OR REPLACE TEMP VIEW deduped_view AS
SELECT DISTINCT * FROM existing_table_name
WHERE specific_column IS NOT NULL;

-- Overwrite the existing table with the distinct rows
INSERT OVERWRITE TABLE existing_table_name
SELECT * FROM deduped_view;
```

These methods will help you remove duplicate rows based on a specific column efficiently

▼ **26. Validate that the primary key is unique across all rows.**

- To validate that the primary key is unique across all rows in a Databricks Delta table using PySpark, you can use the following code:

```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, count

# Initialize Spark session
spark = SparkSession.builder.appName("PrimaryKeyValidation").getOrCreate()
```

```
# Load your Delta table
df = spark.read.format("delta").load("path_to_your_delta_table")

# Group by the primary key column and count occurrences
duplicate_keys = df.groupBy("primary_key_column").agg(count("*").alias("cc
                    .filter(col("count") > 1)

# Show duplicate keys if any
duplicate_keys.show()
```

- To validate that the primary key is unique across all rows in a database table, you can use a SQL query. Here's a simple way to check for duplicates:

```
SELECT primary_key_column, COUNT(*)
FROM table_name
GROUP BY primary_key_column
HAVING COUNT(*) > 1;
```

▼ **27.Validate that a field is associated with just one unique value in another field.**

To validate that a field is associated with just one unique value in another field using PySpark, you can follow these steps:

1. **Group by the first field**.

2. **Count the distinct values of the second field**.

3. **Filter the results to find any groups with more than one distinct value**.

Here's an example in PySpark:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, countDistinct

# Initialize Spark session
spark = SparkSession.builder.appName("FieldAssociationValidation").getO
rCreate()

# Load your Delta table
df = spark.read.format("delta").load("path_to_your_delta_table")

# Group by the first field and count distinct values of the second fiel
d
```

```
validation_df = df.groupBy("field1").agg(countDistinct("field2").alias
("distinct_count")).filter(col("distinct_count") > 1)

# Show results if any
validation_df.show()
```

In this script:

- Replace `"field1"` with the name of the field you want to check for unique associations.
- Replace `"field2"` with the name of the field that should have unique values for each value in `"field1"`.
- Replace `"path_to_your_delta_table"` with the actual path to your Delta table.

This will show any instances where a value in `field1` is associated with more than one unique value in `field2`.

▼ **28. Validate that a value is not present in a specific field**

To validate that a specific value is not present in a particular field using PySpark, you can use the following approach:

1. **Filter the DataFrame** to check for the presence of the value.

2. **Count the number of rows** that match the condition.

Here's an example:

```
from pyspark.sql import SparkSession

# Initialize Spark session
spark = SparkSession.builder.appName("ValueValidation").getOrCreate()

# Load your Delta table
df = spark.read.format("delta").load("path_to_your_delta_table")

# Define the value and field to check
value_to_check = "value_to_check"
field_to_check = "field_to_check"

# Filter the DataFrame to find rows with the specific value
filtered_df = df.filter(col(field_to_check) == value_to_check)

# Count the number of matching rows
count = filtered_df.count()

# Validate if the value is not present
```

```
if count == 0:
    print(f"The value '{value_to_check}' is not present in the field
'{field_to_check}'.")
else:
    print(f"The value '{value_to_check}' is present in the field '{fiel
d_to_check}' {count} times.")
```

In this script:

- Replace `"value_to_check"` with the value you want to validate.
- Replace `"field_to_check"` with the name of the field you want to check.
- Replace `"path_to_your_delta_table"` with the actual path to your Delta table.

This will check if the specified value is present in the given field and print the result accordingly.

▼ **29.Cast a column to a timestamp.**

To cast a column to a timestamp in Databricks notebooks, you can use the `to_timestamp` function or the `cast` function. Here are examples of both methods:

## Using `to_timestamp` Function

The `to_timestamp` function converts a string expression to a timestamp using an optional format.

```
SELECT to_timestamp(column_name, 'yyyy-MM-dd HH:mm:ss') AS new_column_n
ame
FROM table_name;
```

## Using `cast` Function

The `cast` function can also be used to convert a string to a timestamp.

```
SELECT cast(column_name AS TIMESTAMP) AS new_column_name
FROM table_name;
```

## Example in PySpark

If you're using PySpark, you can achieve the same with the following code:

```
from pyspark.sql.functions import to_timestamp

df = df.withColumn('new_column_name', to_timestamp(df['column_name'],
'yyyy-MM-dd HH:mm:ss'))
```

These methods should help you convert your column to a timestamp format in Databricks

▼ **30.Extract calendar data from a timestamp.**

To extract calendar data from a timestamp in Databricks, you can use the `extract` function. This function allows you to retrieve specific parts of a timestamp, such as the year, month, day, hour, minute, or second. Here are some examples:

## Using SQL in Databricks

```sql
SELECT
    extract(YEAR FROM timestamp_column) AS year,
    extract(MONTH FROM timestamp_column) AS month,
    extract(DAY FROM timestamp_column) AS day,
    extract(HOUR FROM timestamp_column) AS hour,
    extract(MINUTE FROM timestamp_column) AS minute,
    extract(SECOND FROM timestamp_column) AS second
FROM
    table_name;
```

## Using PySpark

If you're working with PySpark, you can use the `year`, `month`, `dayofmonth`, `hour`, `minute`, and `second` functions from `pyspark.sql.functions`:

```python
from pyspark.sql.functions import year, month, dayofmonth, hour, minute, second

df = df.withColumn('year', year(df['timestamp_column']))
        .withColumn('month', month(df['timestamp_column']))
        .withColumn('day', dayofmonth(df['timestamp_column']))
        .withColumn('hour', hour(df['timestamp_column']))
        .withColumn('minute', minute(df['timestamp_column']))
        .withColumn('second', second(df['timestamp_column']))
```

These methods should help you extract the specific calendar data you need from your timestamp column

▼ **31.Extract a specific pattern from an existing string column**

To extract a specific pattern from an existing string column in Databricks, you can use the `regexp_extract` function. This function allows you to use regular expressions to find and extract patterns from strings. Here's how you can do it:

## Using SQL in Databricks

```
SELECT regexp_extract(column_name, 'your_regex_pattern', 1) AS extracte
d_pattern
FROM table_name;
```

- `column_name` : The name of the column you want to extract the pattern from.
- `your_regex_pattern` : The regular expression pattern you want to match.
- `1` : The group index of the pattern you want to extract. Use `0` to match the entire pattern.

## Example

If you have a column with email addresses and you want to extract the domain part of the email, you can use:

```
SELECT regexp_extract(email_column, '@([^\\s]+)', 1) AS domain
FROM users;
```

## Using PySpark

In PySpark, you can achieve the same with the `regexp_extract` function from `pyspark.sql.functions` :

```
from pyspark.sql.functions import regexp_extract

df = df.withColumn('extracted_pattern', regexp_extract(df['column_nam
e'], 'your_regex_pattern', 1))
```

## Example

For extracting the domain from an email column:

```
df = df.withColumn('domain', regexp_extract(df['email_column'], '@([^
\\s]+)', 1))
```

These methods should help you extract specific patterns from your string columns in Databricks

▼ **32.Utilize the dot syntax to extract nested data field**

To extract nested data fields in Spark using dot syntax, you can follow these steps. This is particularly useful when dealing with complex data types like structs in JSON or Parquet files.

## Using PySpark

In PySpark, you can use dot notation to access nested fields directly. Here's an example:

```python
from pyspark.sql import SparkSession

# Initialize SparkSession
spark = SparkSession.builder.appName("NestedFieldsExample").getOrCreate
()

# Sample data with nested structure
data = [
    ("John", {"city": "New York", "state": "NY"}),
    ("Alice", {"city": "Los Angeles", "state": "CA"})
]

# Define schema
schema = "name STRING, address STRUCT<city: STRING, state: STRING>"

# Create DataFrame
df = spark.createDataFrame(data, schema)

# Show DataFrame
df.show(truncate=False)

# Select nested fields using dot notation
df.select("address.city", "address.state").show()
```

## Using SQL in Databricks

If you are using SQL in Databricks, you can also use dot notation to access nested fields:

```sql
SELECT
    address.city AS city,
    address.state AS state
FROM
    table_name;
```

## Handling Columns with Dots in Their Names

If your column names contain dots and you need to reference them, you should use backticks to escape the column names:

```python
# Using backticks to escape column names with dots
df.select("`address.city`", "`address.state`").show()
```

These methods should help you efficiently extract nested fields in Spark

▼ **33.Identify the benefits of using array function**

Using the `array` function in Databricks and Spark offers several benefits, especially when dealing with complex data structures. Here are some key advantages:

## 1. Simplifies Data Manipulation

The `array` function allows you to group multiple values into a single array, making it easier to handle and manipulate complex data. This is particularly useful when working with nested data structures.

## 2. Enhanced Readability and Maintainability

By using arrays, you can make your code more readable and maintainable. Instead of managing multiple columns separately, you can work with a single array column, which simplifies your queries and transformations.

## 3. Efficient Data Processing

Arrays enable efficient data processing by allowing you to perform operations on entire collections of elements at once. This can lead to performance improvements, especially when combined with higher-order functions like `transform` , `filter` , and `aggregate` .

## 4. Flexibility in Data Modeling

Arrays provide flexibility in data modeling, allowing you to represent complex relationships and hierarchies within your data. This can be particularly useful for applications that require handling of multi-valued attributes or nested data.

## 5. Support for Higher-Order Functions

Databricks and Spark support higher-order functions that operate on arrays, such as `transform` , `filter` , and `aggregate` . These functions allow you to apply custom logic to each element of an array, enabling powerful and expressive data transformations

## Example in PySpark

Here's a simple example of creating and using an array in PySpark:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import array, col

# Initialize SparkSession
spark = SparkSession.builder.appName("ArrayExample").getOrCreate()

# Sample data
```

```
data = [(1, "a"), (2, "b"), (3, "c")]

# Create DataFrame
df = spark.createDataFrame(data, ["id", "value"])

# Create an array column
df = df.withColumn("array_col", array(col("id"), col("value")))

# Show DataFrame
df.show(truncate=False)
```

This example demonstrates how to create an array column from existing columns, making it easier to manage and process the data.

▼ **34.Parse JSON strings into structs.**

To parse JSON strings into structs in Databricks, you can use the `from_json` function. This function converts a JSON string into a struct (or other complex data types) based on a specified schema. Here's how you can do it:

## Using SQL in Databricks

```
SELECT from_json(column_name, 'STRUCT<field1: STRING, field2: INT>') AS
parsed_column
FROM table_name;
```

## Example

Suppose you have a JSON string column named `json_column` and you want to parse it into a struct with fields `name` (string) and `age` (integer):

```
SELECT from_json(json_column, 'STRUCT<name: STRING, age: INT>') AS pars
ed_data
FROM table_name;
```

## Using PySpark

In PySpark, you can achieve the same using the `from_json` function from `pyspark.sql.functions` :

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import from_json
from pyspark.sql.types import StructType, StructField, StringType, Inte
gerType

# Initialize SparkSession
```

```
spark = SparkSession.builder.appName("ParseJSONExample").getOrCreate()

# Sample data
data = [("John", '{"name": "John", "age": 30}'), ("Alice", '{"name": "A
lice", "age": 25}')]

# Define schema
schema = StructType([
    StructField("name", StringType(), True),
    StructField("age", IntegerType(), True)
])

# Create DataFrame
df = spark.createDataFrame(data, ["id", "json_column"])

# Parse JSON column
df = df.withColumn("parsed_data", from_json(df["json_column"], schema))

# Show DataFrame
df.show(truncate=False)
```

### Explanation

- **SQL**: The `from_json` function takes the JSON string column and the schema definition as arguments, converting the JSON string into a struct.
- **PySpark**: The `from_json` function is used similarly, but you need to define the schema using `StructType` and `StructField`.

▼ **35.Identify which result will be returned based on a join query.**

- The result of a join query in Databricks (or any SQL-based system) depends on the type of join used and the data in the tables being joined. Here's a brief overview of different join types and the results they produce:

### 1. Inner Join

Returns only the rows that have matching values in both tables.

```
SELECT *
FROM table1
INNER JOIN table2
ON table1.id = table2.id;
```

**Result**: Rows where `table1.id` matches `table2.id`.

## 2. Left (Outer) Join

Returns all rows from the left table, and the matched rows from the right table. If no match is found, NULL values are returned for columns from the right table.

```
SELECT *
FROM table1
LEFT JOIN table2
ON table1.id = table2.id;
```

**Result**: All rows from `table1`, with matching rows from `table2` or NULL if no match.

## 3. Right (Outer) Join

Returns all rows from the right table, and the matched rows from the left table. If no match is found, NULL values are returned for columns from the left table.

```
SELECT *
FROM table1
RIGHT JOIN table2
ON table1.id = table2.id;
```

**Result**: All rows from `table2`, with matching rows from `table1` or NULL if no match.

## 4. Full (Outer) Join

Returns all rows when there is a match in either table. Rows without a match in one of the tables will have NULLs for columns from that table.

```
SELECT *
FROM table1
FULL JOIN table2
ON table1.id = table2.id;
```

**Result**: All rows from both tables, with NULLs where there is no match.

## 5. Left Semi Join

Returns only the rows from the left table that have a match in the right table.

```
SELECT *
FROM table1
LEFT SEMI JOIN table2
ON table1.id = table2.id;
```

**Result**: Rows from `table1` that have a matching row in `table2`.

## 6. Left Anti Join

Returns only the rows from the left table that do not have a match in the right table.

```
SELECT *
FROM table1
LEFT ANTI JOIN table2
ON table1.id = table2.id;
```

**Result**: Rows from `table1` that do not have a matching row in `table2`.

## 7. Cross Join

Returns the Cartesian product of the two tables, i.e., all possible combinations of rows.

```
SELECT *
FROM table1
CROSS JOIN table2;
```

**Result**: Every row from `table1` combined with every row from `table2`.

## Example Scenario

Suppose you have two tables, `employees` and `departments`:

**employees**:

| id | name | dept_id |
|----|------|---------|
| 1 | John | 10 |
| 2 | Alice | 20 |
| 3 | Bob | 30 |

**departments**:

| dept_id | dept_name |
|---------|-----------|
| 10 | HR |
| 20 | Engineering |
| 40 | Sales |

## Inner Join Example:

```
SELECT employees.name, departments.dept_name
FROM employees
INNER JOIN departments
ON employees.dept_id = departments.dept_id;
```

**Result**:

| name | dept_name |
|------|-----------|
| John | HR |
| Alice | Engineering |

This query returns only the employees who have a matching department.

▼ **36.Identify a scenario to use the explode function versus the flatten function**

The `explode` and `flatten` functions in Databricks serve different purposes and are used in different scenarios. Here's a detailed explanation of when to use each:

## `explode` Function

The `explode` function is used to transform an array or map into a set of rows. Each element of the array or each key-value pair of the map becomes a separate row. This is particularly useful when you need to work with individual elements of an array or map as separate rows.

## Scenario for `explode`

Suppose you have a DataFrame with a column containing arrays of values, and you want to analyze each value individually. For example, you have a DataFrame with a column `tags` that contains arrays of tags for different articles, and you want to count the occurrence of each tag.

**Example**:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import explode

# Initialize SparkSession
spark = SparkSession.builder.appName("ExplodeExample").getOrCreate()

# Sample data
data = [("Article1", ["tag1", "tag2", "tag3"]), ("Article2", ["tag2",
"tag3", "tag4"])]

# Create DataFrame
df = spark.createDataFrame(data, ["article", "tags"])

# Explode the tags column
df_exploded = df.withColumn("tag", explode(df["tags"]))

# Show the result
df_exploded.show()
```

**Result**:

```
+---------+------------+-----+
| article | tags       | tag |
+---------+------------+-----+
| Article1| [tag1, tag2, tag3]| tag1|
| Article1| [tag1, tag2, tag3]| tag2|
| Article1| [tag1, tag2, tag3]| tag3|
| Article2| [tag2, tag3, tag4]| tag2|
| Article2| [tag2, tag3, tag4]| tag3|
| Article2| [tag2, tag3, tag4]| tag4|
+---------+------------+-----+
```

## `flatten` Function

The `flatten` function is used to transform an array of arrays into a single array. This is useful when you have nested arrays and you want to combine them into a single array.

## Scenario for `flatten`

Suppose you have a DataFrame with a column containing nested arrays, and you want to merge these nested arrays into a single array. For example, you have a DataFrame with a column `nested_tags` that contains arrays of arrays of tags, and you want to create a single array of tags for each row.

**Example**:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import flatten

# Initialize SparkSession
spark = SparkSession.builder.appName("FlattenExample").getOrCreate()

# Sample data
data = [("Article1", [["tag1", "tag2"], ["tag3"]]), ("Article2", [["tag2", "tag3"], ["tag4"]])]

# Create DataFrame
df = spark.createDataFrame(data, ["article", "nested_tags"])

# Flatten the nested_tags column
df_flattened = df.withColumn("tags", flatten(df["nested_tags"]))
```

```
# Show the result
df_flattened.show(truncate=False)
```

**Result:**

```
+---------+-----------------+-----------------+
| article | nested_tags     | tags            |
+---------+-----------------+-----------------+
| Article1| [[tag1, tag2], [tag3]]| [tag1, tag2, tag3]|
| Article2| [[tag2, tag3], [tag4]]| [tag2, tag3, tag4]|
+---------+-----------------+-----------------+
```

## Summary

- **Use** `explode` when you need to transform an array or map into multiple rows, each representing an element of the array or a key-value pair of the map.

- **Use** `flatten` when you need to merge nested arrays into a single array.

▼ **37.Identify the PIVOT clause as a way to convert data from a long format to a wide format.**

The `PIVOT` clause in SQL is a powerful tool for transforming data from a long format to a wide format. This is particularly useful when you need to summarize and reorganize data for reporting or analysis.

## What is the PIVOT Clause?

The `PIVOT` clause rotates rows into columns, allowing you to aggregate data and create a more readable and structured format. This is especially helpful when you have categorical data that you want to spread across columns.

## Example Scenario

Suppose you have a table `sales` with the following data:

| year | quarter | region | sales |
|------|---------|--------|-------|
| 2018 | 1 | east | 100 |
| 2018 | 2 | east | 20 |
| 2018 | 3 | east | 40 |
| 2018 | 4 | east | 40 |
| 2019 | 1 | east | 120 |
| 2019 | 2 | east | 110 |
| 2019 | 3 | east | 80 |
| 2019 | 4 | east | 60 |
| 2018 | 1 | west | 105 |

| 2018 | 2 | west | 25 |
| 2018 | 3 | west | 45 |
| 2018 | 4 | west | 45 |
| 2019 | 1 | west | 125 |
| 2019 | 2 | west | 115 |
| 2019 | 3 | west | 85 |
| 2019 | 4 | west | 65 |

You want to convert this data into a wide format where each quarter's sales are in separate columns.

## Using the PIVOT Clause

Here's how you can use the `PIVOT` clause to achieve this:

```
SELECT year, region, q1, q2, q3, q4
FROM sales
PIVOT (
    SUM(sales) AS sales
    FOR quarter IN (1 AS q1, 2 AS q2, 3 AS q3, 4 AS q4)
);
```

**Result**:

| year | region | q1 | q2 | q3 | q4 |
|------|--------|-----|-----|-----|-----|
| 2018 | east | 100 | 20 | 40 | 40 |
| 2019 | east | 120 | 110 | 80 | 60 |
| 2018 | west | 105 | 25 | 45 | 45 |
| 2019 | west | 125 | 115 | 85 | 65 |

## Explanation

- `SUM(sales) AS sales`: Aggregates the sales data.

- `FOR quarter IN (1 AS q1, 2 AS q2, 3 AS q3, 4 AS q4)`: Specifies the pivot columns and their new names.

## Benefits of Using PIVOT

- **Simplifies Data Analysis**: Makes it easier to compare data across different categories.

- **Improves Readability**: Transforms data into a more intuitive format.

- **Reduces Complexity**: Eliminates the need for multiple `CASE` statements or complex `GROUP BY` queries.

The `PIVOT` clause is a versatile tool for data transformation in SQL, making it easier to work with and analyze

▼ **38.Define a SQL UDF**

A SQL User-Defined Function (UDF) allows you to extend the capabilities of SQL by defining custom functions that can be used in SQL queries. These functions can perform complex calculations, transformations, or custom data manipulations that are not available with built-in SQL functions.

## Creating a SQL UDF in Databricks

Here's how you can define a SQL UDF in Databricks:

## Example: Creating a Simple SQL UDF

Suppose you want to create a UDF that calculates the length of a string. You can define this function as follows:

```
-- Create a SQL UDF to calculate the length of a string
CREATE OR REPLACE FUNCTION get_string_length(input_string STRING)
RETURNS INT
RETURN LENGTH(input_string);
```

- `CREATE OR REPLACE FUNCTION` : This statement is used to create a new UDF or replace an existing one.

- `get_string_length` : The name of the UDF.

- `input_string STRING` : The parameter for the UDF, which is a string in this case.

- `RETURNS INT` : Specifies that the function returns an integer.

- `RETURN LENGTH(input_string)` : The logic of the function, which calculates the length of the input string.

## Using the SQL UDF

Once the UDF is created, you can use it in your SQL queries just like any other function:

```
-- Use the UDF in a SQL query
SELECT name, get_string_length(name) AS name_length
FROM your_table;
```

## Benefits of SQL UDFs

- **Reusability**: UDFs can be reused across multiple queries, making your code more modular and maintainable.

- **Abstraction**: They provide a layer of abstraction, simplifying complex logic and making SQL queries more readable.
- **Performance**: SQL UDFs are optimized by the SQL engine, making them more efficient than external UDFs written in other languages.

## Example in PySpark

If you prefer to define a UDF using PySpark, you can do so as follows:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import udf
from pyspark.sql.types import IntegerType

# Initialize SparkSession
spark = SparkSession.builder.appName("UDFExample").getOrCreate()

# Define the UDF
@udf(returnType=IntegerType())
def get_string_length(input_string):
    return len(input_string)

# Sample DataFrame
data = [("Alice",), ("Bob",), ("Carol",)]
df = spark.createDataFrame(data, ["name"])

# Use the UDF in a DataFrame query
df = df.withColumn("name_length", get_string_length(df["name"]))
df.show()
```

This example demonstrates how to create and use a UDF in PySpark to calculate the length of a string

▼ **39.Identify the location of a function.**

To identify the location of a function in Databricks, you can use the "Go to definition" and "Peek definition" features available in Databricks notebooks. These features help you quickly navigate to the definition of a function within your code.

## Steps to Use "Go to Definition" and "Peek Definition"

1. **Go to Definition**:
   - Right-click on the function name in your notebook.
   - Select "Go to definition" or press `F12` .

- This will take you directly to the location where the function is defined.

2. **Peek Definition**:

   - Right-click on the function name.

   - Select "Peek definition".

   - This will show a small window with the function's definition without navigating away from your current location.

## Example

Suppose you have a function `my_function` defined in your notebook:

```
def my_function(x):
    return x * 2


# Using the function
result = my_function(5)
```

- To find where `my_function` is defined, right-click on `my_function` and choose "Go to definition" or "Peek definition".

## Troubleshooting

If these features do not work, ensure that:

- Your notebook is attached to a cluster.

- You are using a supported Databricks Runtime version (12.2 LTS and above)

▼ **40.Describe the security model for sharing SQL UDFs.**

The security model for sharing SQL User-Defined Functions (UDFs) in Databricks is designed to ensure that these functions are used safely and securely across different users and environments. Here are the key aspects of this security model:

### 1. Access Control

SQL UDFs can be controlled using Access Control Language (ACL). This allows administrators to define who can create, modify, and execute UDFs. By setting appropriate permissions, you can ensure that only authorized users have access to sensitive functions

1

.

### 2. Unity Catalog Integration

Databricks integrates with Unity Catalog, which provides fine-grained access controls for data and UDFs. UDFs registered in Unity Catalog can have associated permissions, ensuring that

only users with the necessary privileges can use them

2

.

## 3. Session-Based and Registered UDFs

- **Session-Based UDFs**: These are scoped to the current SparkSession and are not shared across sessions. They are useful for temporary or ad-hoc analyses.
- **Registered UDFs**: These are registered in Unity Catalog and can be shared across multiple queries, sessions, and users. They come with access controls to manage who can use them.

    2

## 4. Isolation and Sandboxing

To prevent malicious code execution, SQL UDFs are isolated from the core system. This isolation ensures that UDFs cannot perform unauthorized actions or access restricted resources

1

.

## 5. Auditing and Monitoring

Databricks provides auditing and monitoring capabilities to track the usage of UDFs. This helps in identifying any unauthorized access or misuse of UDFs, ensuring compliance with security policies

1

.

## Example of Creating a Secure SQL UDF

Here's an example of creating a SQL UDF with access control:

```sql
-- Create a SQL UDF to calculate the length of a string
CREATE OR REPLACE FUNCTION get_string_length(input_string STRING)
RETURNS INT
RETURN LENGTH(input_string);

-- Grant execute permission to a specific user
GRANT EXECUTE ON FUNCTION get_string_length TO user_name;
```

This example demonstrates how to create a UDF and grant execute permissions to a specific user, ensuring that only authorized users can run the function.

## Summary

The security model for sharing SQL UDFs in Databricks includes access control, integration with Unity Catalog, session-based and registered UDFs, isolation and sandboxing, and auditing and monitoring. These measures ensure that UDFs are used securely and efficiently across different environments

▼ **41.Use CASE/WHEN in SQL code.**

The `CASE` statement in SQL is used to create conditional logic within your queries. It allows you to perform different actions based on different conditions, similar to an `if-else` statement in programming languages. Here's how you can use the `CASE` statement in SQL:

## Basic Syntax

```sql
SELECT
    CASE
        WHEN condition1 THEN result1
        WHEN condition2 THEN result2
        ...
        ELSE default_result
    END AS new_column_name
FROM
    table_name;
```

▼ **42.Leverage CASE/WHEN for custom control flow**

Using the `CASE` statement in Databricks SQL allows you to implement custom control flow within your queries. This is particularly useful for creating conditional logic to handle different scenarios based on your data. Here's how you can leverage `CASE/WHEN` for custom control flow in Databricks:

## Example Scenario

Suppose you have a table `sales` with columns `year`, `quarter`, and `revenue`, and you want to categorize the revenue into different performance levels.

**Table: sales**

| year | quarter | revenue |
|------|---------|---------|
| 2023 | 1 | 10000 |
| 2023 | 2 | 15000 |
| 2023 | 3 | 20000 |
| 2023 | 4 | 25000 |

## Using CASE/WHEN

You can use the `CASE` statement to create a new column `performance` that categorizes the revenue:

```sql
SELECT
    year,
    quarter,
    revenue,
    CASE
        WHEN revenue >= 20000 THEN 'High'
        WHEN revenue >= 15000 THEN 'Medium'
        ELSE 'Low'
    END AS performance
FROM
    sales;
```

**Result**:

| year | quarter | revenue | performance |
|------|---------|---------|-------------|
| 2023 | 1       | 10000   | Low         |
| 2023 | 2       | 15000   | Medium      |
| 2023 | 3       | 20000   | High        |
| 2023 | 4       | 25000   | High        |

## Explanation

- `CASE` : Starts the conditional logic.

- `WHEN condition THEN result` : Specifies the condition and the result if the condition is true.

- `ELSE default_result` : Specifies the default result if none of the conditions are true.

- `END` : Ends the `CASE` statement.

## Nested CASE Statements

For more complex logic, you can nest `CASE` statements:

```sql
SELECT
    year,
    quarter,
    revenue,
    CASE
        WHEN revenue >= 20000 THEN 'High'
        WHEN revenue >= 15000 THEN
            CASE
                WHEN revenue >= 18000 THEN 'Medium-High'
```

```
            ELSE 'Medium'
        END
    ELSE 'Low'
  END AS performance
FROM
    sales;
```

## Benefits of Using CASE/WHEN

- **Flexibility**: Allows you to implement complex conditional logic directly within your SQL queries.

- **Readability**: Makes your queries more readable and easier to understand.

- **Efficiency**: Reduces the need for multiple queries or complex joins to achieve the same result.

## Practical Use Cases

- **Data Categorization**: Categorize numerical data into different ranges (e.g., low, medium, high).

- **Conditional Aggregation**: Apply different aggregation rules based on conditions.

- **Data Cleaning**: Replace or modify values based on specific conditions.

# Section 3: Incremental Data Processing(22%)

▼ **43.Identify where Delta Lake provides ACID transactions**

Delta Lake provides ACID transactions in the context of data lakes, ensuring reliable and consistent data management. Here are the key areas where Delta Lake implements ACID transactions:

1. **Atomicity**: Ensures that all operations within a transaction are completed successfully or none at all. This prevents partial updates that could lead to data corruption.

2. **Consistency**: Guarantees that transactions only bring the database from one valid state to another, maintaining data integrity.

3. **Isolation**: Ensures that concurrent transactions do not interfere with each other, providing a consistent view of the data.

4. **Durability:** Ensures that once a transaction is committed, it remains so, even in the event of a system failure.

**Delta Lake achieves these properties through its transaction log**, which records all changes and ensures that data operations are reliable and consistent

▼ **44.Identify the benefits of ACID transaction**

ACID transactions offer several key benefits that are crucial for maintaining reliable and consistent databases:

1. **Data Integrity**: ACID transactions ensure that data remains **accurate and consistent**, even in the event of errors, system failures, or power outages. This is essential for applications where data accuracy is critical, such as financial systems.

2. **Simplified Concurrency Control**: They **allow multiple transactions to occur simultaneously** without interfering with each other. This is particularly important in environments where many users or processes need to access and modify the same data concurrently.

3. **Intuitive Data Access Logic**: ACID transactions simplify the **logic required to manage data operations**. Developers don't need to write complex code to handle partial updates or rollbacks, as the database system ensures that all operations within a transaction are completed successfully or not at all.

4. **Future-Proofing Database Needs**: By ensuring that transactions are handled reliably, ACID properties help future-proof databases against potential issues that could arise from data corruption or inconsistencies.

These benefits make ACID transactions a cornerstone of modern database management, providing a robust framework for data reliability and consistency

▼ **45.Identify whether a transaction is ACID-compliant.**

To determine if a transaction is ACID-compliant, you need to check if it adheres to the following four properties:

1. **Atomicity**:

2. **Consistency**:

3. **Isolation**:

4. **Durability**:

If a transaction meets all these criteria, it can be considered ACID-compliant. This ensures reliable and consistent data management, which is crucial for applications where data integrity is paramount.

## ▼ 46.Compare and contrast data and metadata.

comparison between data and metadata:

## Data

- **Definition**: Data refers to raw facts and figures that can be processed to extract meaningful information. It can include numbers, text, images, and more.

- **Examples**: Sales figures, customer names, transaction records, sensor readings.

- **Purpose**: Used to perform analysis, make decisions, and derive insights.

- **Processing**: Data may or may not be processed before it is stored.

## Metadata

- **Definition**: Metadata is data about data. It provides information that helps describe, manage, and understand the data.
- **Examples**: File size, creation date, author, data type, and tags.
- **Purpose**: Helps in organizing, finding, and managing data. It provides context and makes data easier to use.
- **Processing**: Metadata is always processed information.

## Key Differences

- **Informative Value**: Data can be raw and unprocessed, while metadata is always informative and processed.
- **Usage**: Data is used for analysis and decision-making, whereas metadata helps in understanding and managing the data.
- **Storage**: Data can be stored in its raw form, while metadata is stored as processed information.

In summary, while data represents the actual content, metadata provides essential information about that content, making it easier to manage and utilize effectively.

## ▼ 47.Compare and contrast managed and external tables.

## Managed Tables

- **Definition**: Managed tables, also known as internal tables, are fully controlled by the database management system (DBMS). The system manages both the data and the

metadata.

- **Data Storage:** Data is stored within the database's storage system (e.g., HDFS for Hadoop, a specific cloud storage service).

- **Data Management:** When a managed table is dropped, both the data and the metadata are deleted by the DBMS.

- **Use Case:** Ideal for scenarios where the DBMS should have full control over the data lifecycle, including storage, management, and deletion.

## External Tables

- **Definition:** External tables store metadata within the DBMS, but the actual data resides outside the database, typically in an external storage system (e.g., Amazon S3, Azure Blob Storage).

- **Data Storage:** Data is stored externally, outside the database's storage system.

- **Data Management:** When an external table is dropped, only the metadata is removed; the actual data remains intact in the external storage.

- **Use Case:** Suitable for scenarios where data needs to be shared across multiple systems or where the data lifecycle is managed outside the DBMS.

## Key Differences

- **Control:** Managed tables give the DBMS full control over data and metadata, while external tables only allow the DBMS to manage metadata.

- **Data Deletion:** Dropping a managed table deletes both data and metadata, whereas dropping an external table only deletes the metadata.

- **Flexibility:** External tables offer more flexibility in terms of data storage and sharing across different systems

▼ **48.Identify a scenario to use an external table.**

A great scenario to use an external table is when you need to integrate and query large datasets stored in a data lake without moving the data into your database. Here are some specific examples:

1. **Data Lakes Integration:** If your organization uses a data lake to store vast amounts of raw, unstructured data, external tables allow you to query this data directly without the need to load it into your data warehouse. This is efficient and cost-effective.

2. **Real-time Data Analysis:** When dealing with streaming data or real-time analytics, external tables enable you to query the data as it arrives, providing up-to-date insights without the latency of data movement.

3. **Cost-Effective Storage:** Storing large volumes of data in a data warehouse can be expensive. External tables allow you to keep the data in cheaper storage solutions (like

cloud storage) and query it as needed, reducing storage costs.

4. **Data Sharing Across Systems**: If you need to share data across multiple systems or platforms, external tables can be used to provide access to the data without duplicating it.

These scenarios highlight the flexibility and efficiency of using external tables, especially in environments where data is large, dynamic, and doesn't require heavy transformations before analysis.

▼ **49.Create a managed table.**

Creating a managed table involves defining the table structure and specifying the data it will contain. Here's a basic example using SQL:

## SQL Example

```
CREATE TABLE my_catalog.my_schema.my_managed_table (
    id INT,
    name STRING,
    age INT,
    created_at TIMESTAMP
);
```

In this example:

- `my_catalog` : The catalog that will contain the table.

- `my_schema` : The schema that will contain the table.

- `my_managed_table` : The name of the table.

- **Columns**: `id` , `name` , `age` , and `created_at` with their respective data types.

## Steps to Create a Managed Table

1. **Define the Table Structure**: Specify the columns and their data types.

2. **Execute the SQL Command**: Run the `CREATE TABLE` command in your SQL environment.

## Example in Databricks

If you're using Databricks, you can create a managed table using the following command:

```
CREATE TABLE my_catalog.my_schema.my_managed_table (
    id INT,
    name STRING,
    age INT,
    created_at TIMESTAMP
);
```

This command will create a managed table where Databricks manages both the data and the metadata

▼ **50.Identify the location of a table.**

To identify the location of a table, you can use various methods depending on the database or data management system you're using. Here are some common approaches:

## SQL Databases

1. **Query the System Catalog**: Most SQL databases have system catalogs or information schemas that store metadata about tables. You can query these to find the location of a table.

```
SELECT table_schema, table_name
FROM information_schema.tables
WHERE table_name = 'your_table_name';
```

2. **Database-Specific Commands**: Some databases have specific commands to show table locations. For example, in PostgreSQL:

```
\dt your_table_name
```

## Data Lakes (e.g., Delta Lake)

1. **Check the Metadata**: In systems like Delta Lake, you can check the metadata logs to find the location of a table.

```
DESCRIBE DETAIL your_table_name;
```

▼ **51.Inspect the directory structure of Delta Lake file.**

Delta Lake organizes its files in a specific directory structure to manage data and maintain ACID transactions. Here's an overview of the key components:

## Directory Structure

1. **Root Directory**: This is the main directory where your Delta table is stored. It contains all the data files and subdirectories.

   - Example: `/path/to/delta-table/`

2. **Data Files**: These are the actual data files stored in formats like Parquet. They contain the table's data.

   - Example: `/path/to/delta-table/part-00000-xxxx.snappy.parquet`

3. **delta_log Directory**: This subdirectory contains the transaction logs. These logs record all changes made to the Delta table, ensuring ACID compliance.

- Example: `/path/to/delta-table/_delta_log/`

- Inside `_delta_log` , you'll find:

  - **JSON Files**: Each transaction is recorded in a JSON file.

    - Example: `/path/to/delta-table/_delta_log/00000000000000000010.json`

  - **Checkpoint Files**: These Parquet files provide a snapshot of the table's state at a specific point in time, improving query performance.

    - Example: `/path/to/delta-table/_delta_log/00000000000000000010.checkpoint.parquet`

## Example Structure

```
/path/to/delta-table/
    ├── part-00000-xxxx.snappy.parquet
    ├── part-00001-xxxx.snappy.parquet
    ├── _delta_log/
         ├── 00000000000000000000.json
         ├── 00000000000000000001.json
         ├── 00000000000000000001.checkpoint.parquet
```

## Key Points

- **Transaction Logs**: The `_delta_log` directory is crucial for maintaining the ACID properties of Delta Lake. It logs every transaction, ensuring data consistency and reliability.

- **Data Organization**: The data files are stored in a structured manner, allowing efficient data retrieval and management.

This directory structure helps Delta Lake provide robust data management and transactional guarantees

▼ **52.Identify who has written previous versions of a table. ‖ 53.Review a history of table transactions.**

To identify who has written previous versions of a Delta Lake table, you can use the `DESCRIBE HISTORY` command. This command provides detailed information about each operation performed on the table, including the user who executed the operation. Here's how you can do it:

## Using SQL

```
DESCRIBE HISTORY your_table_name;
```

This command will return a history of all operations on the table, including:

- **version**: The version number of the table.

- **timestamp**: When the operation was committed.

- **userId**: The ID of the user who performed the operation.

- **userName**: The name of the user who performed the operation.

- **operation**: The type of operation (e.g., INSERT, DELETE).

- **operationParameters**: Parameters of the operation.

- **job**: Details of the job that ran the operation.

- **notebook**: Details of the notebook from which the operation was run.

- **clusterId**: The ID of the cluster on which the operation ran

## Additional Information

- **Retention**: The history of table transactions is retained for 30 days by default, but this can be configured using the `delta.logRetentionDuration` setting.

    1

- **Time Travel**: You can also use Delta Lake's time travel feature to query the table at specific points in time, which can be useful for auditing and rollback purposes

▼ **54.Roll back a table to a previous version.**

To roll back a Delta Lake table to a previous version in Databricks, you can use the `RESTORE` command. This command allows you to revert the table to a specific version or a specific timestamp. Here's how you can do it:

## Using SQL

## Restore to a Specific Version

First, identify the version you want to restore to by using the `DESCRIBE HISTORY` command:

```
DESCRIBE HISTORY your_table_name;
```

Once you have the desired version number, you can restore the table:

```
RESTORE TABLE your_table_name TO VERSION AS OF version_number;s
```

## Restore to a Specific Timestamp

If you prefer to restore the table to a specific point in time, use the timestamp:

```
RESTORE TABLE your_table_name TO TIMESTAMP AS OF 'yyyy-MM-dd HH:mm:ss';
```

## Example

Here's an example of restoring a table named `employee` to version 5:

```
RESTORE TABLE employee TO VERSION AS OF 5;
```

Or restoring to a specific timestamp:

```
RESTORE TABLE employee TO TIMESTAMP AS OF '2024-12-01 00:00:00';
```

### Additional Information

- **Retention**: Ensure that the data files for the version you want to restore to are still available. Data files may be deleted by the `VACUUM` command if they are older than the retention period.
- **Validation**: After restoring, you can validate the table's state by querying it to ensure it has been correctly reverted.

▼ **55.Identify that a table can be rolled back to a previous version.| 56.Query a specific version of a table**

Delta Lake table in Databricks can be rolled back to a previous version. This is done using the `RESTORE` command, which allows you to revert the table to a specific version or a specific timestamp.

### How to Roll Back a Table

### Restore to a Specific Version

First, identify the version you want to restore to by using the `DESCRIBE HISTORY` command:

```
DESCRIBE HISTORY your_table_name;
```

Once you have the desired version number, you can restore the table:

```
RESTORE TABLE your_table_name TO VERSION AS OF version_number;
```

### Restore to a Specific Timestamp

If you prefer to restore the table to a specific point in time, use the timestamp:

```
RESTORE TABLE your_table_name TO TIMESTAMP AS OF 'yyyy-MM-dd HH:mm:ss';
```

### Example

Here's an example of restoring a table named `employee` to version 5:

```
RESTORE TABLE employee TO VERSION AS OF 5;
```

Or restoring to a specific timestamp:

```
RESTORE TABLE employee TO TIMESTAMP AS OF '2024-12-01 00:00:00';
```

## Additional Information

- **Retention**: Ensure that the data files for the version you want to restore to are still available. Data files may be deleted by the `VACUUM` command if they are older than the retention period.

  1

- **Validation**: After restoring, you can validate the table's state by querying it to ensure it has been correctly reverted.

▼ **57.Identify why Z-ordering is beneficial to Delta Lake table**

Z-ordering is a powerful technique used in Delta Lake to optimize query performance by organizing data in a way that enhances data skipping. Here are the key benefits of Z-ordering:

1. **Improved Query Performance**: Z-ordering **colocates related information in the same set of files**. This co-locality allows Delta Lake's data-skipping algorithms to dramatically reduce the amount of data that needs to be read during queries. For example, if you frequently query a table by a specific column, Z-ordering by that column can make those queries much faster.

2. **Efficient Data Skipping**: By organizing data based on the values of one or more columns, Z-ordering helps the query engine skip over large portions of data that do not match the query criteria. This is particularly useful for columns with high cardinality (many unique values).

3. **Reduced I/O Operations**: Since Z-ordering minimizes the number of files that need to be read, it reduces the input/output operations required to execute a query. This leads to faster query execution times and lower resource consumption.

4. **Optimized Storage**: Z-ordering can also help in compacting the data files, reducing the number of small files and improving storage efficiency.

## Example Scenario

Suppose you have a Delta table with a billion rows and you frequently query it by a specific column, such as `date`. By Z-ordering the table by the `date` column, you ensure that rows with similar dates are stored close together. This allows the query engine to skip over large chunks of data that fall outside the queried date range, significantly speeding up the query.

In summary, Z-ordering is beneficial for enhancing query performance, improving data skipping, reducing I/O operations, and optimizing storage in Delta Lake tables.

▼ **58.Identify how vacuum commits deletes.**

Delta Lake, the `VACUUM` command is used to clean up and permanently delete data files that are no longer referenced by the Delta table. This helps in managing storage costs and maintaining the efficiency of the data lake. Here's how `VACUUM` commits deletes:

## How `VACUUM` Works

1. **Identifies Unused Files**: The `VACUUM` command scans the Delta table's transaction log to identify data files that are no longer referenced by any active version of the table.

2. **Retention Period**: By default, `VACUUM` only removes files that are older than a specified retention period (usually 7 days). This retention period ensures that recent changes can still be accessed for time travel queries.

3. **Deletes Files**: Once the unused files are identified and the retention period is checked, `VACUUM` deletes these files from the storage system.

4. **Commit Operation**: The deletion of files is committed, meaning the files are permanently removed and cannot be recovered through Delta Lake's time travel feature.

## Example Command

```
VACUUM your_table_name RETAIN 168 HOURS;
```

This command will remove all files older than 7 days (168 hours) from the Delta table.

## Important Considerations

- **Time Travel**: After running `VACUUM`, you lose the ability to time travel to versions of the table that relied on the deleted files.

- **Dry Run**: You can use the `DRY RUN` option to see which files would be deleted without actually removing them:

```
VACUUM your_table_name DRY RUN;
```

- **Storage Optimization**: Regularly running `VACUUM` helps in optimizing storage by removing obsolete files, thus reducing storage costs.

By understanding and using the `VACUUM` command, you can effectively manage your Delta Lake storage and maintain optimal performance.

## Example Command

```
VACUUM your_table_name DRY RUN;
```

This command will output a list of files that would be deleted if the `VACUUM` command were executed without the `DRY RUN` option.

## Example Output

When you run the `DRY RUN` command, you might see an output like this:

```
Found 10 files and directories in a total of 2 directories that are safe to delete.
```

▼ **59.Identify the kind of files Optimize compact.**

The `OPTIMIZE` command in Delta Lake is used to compact small files into larger ones. This process is known as small file compaction or bin-packing. Here's a detailed look at the kinds of files that `OPTIMIZE` targets:

## Small Files

- **Definition**: Small files are data files that are significantly smaller than the optimal size for efficient querying and storage. These files can be created due to frequent small batch writes or streaming data ingestion.
- **Problem**: Small files cause high I/O overhead, slow query performance, and large metadata transaction logs, which can lead to slower planning times.

## How `OPTIMIZE` Works

- **Compaction**: The `OPTIMIZE` command rewrites these small files into fewer, larger files. This reduces the number of files that need to be read during query execution, improving performance.
- **File Layout**: For tables with partitions, the compaction and data layout are performed within each partition to maintain data locality and improve query efficiency.

## Example Command

```
OPTIMIZE '/path/to/delta/table';
```

This command will compact the small files in the specified Delta table.

## Benefits

- **Improved Query Performance**: By reducing the number of small files, the `OPTIMIZE` command helps in speeding up read queries.
- **Reduced I/O Operations**: Larger files mean fewer files to read, which reduces the I/O operations required for queries.

- **Efficient Storage**: Compaction helps in better utilizing storage by reducing the overhead associated with managing many small files.

In summary, the `OPTIMIZE` command in Delta Lake targets small files to improve query performance, reduce I/O operations, and optimize storage efficiency

▼ **60.Identify CTAS as a solution.**

In Databricks, the **Create Table As Select (CTAS)** command is a versatile solution for various data management tasks. Here's how CTAS can be effectively used in Databricks:

## Scenarios Where CTAS is Useful

1. **Data Transformation and Aggregation**:

   - **Scenario**: You need to transform or aggregate data from an existing table and store the results in a new table.

   - **Solution**: Use CTAS to create a new table with the transformed or aggregated data.

   ```
   CREATE TABLE new_table AS
   SELECT column1, SUM(column2) AS total
   FROM existing_table
   GROUP BY column1;
   ```

2. **Data Backup and Archiving**:

   - **Scenario**: You want to create a backup or archive of a table at a specific point in time.

   - **Solution**: Use CTAS to create a snapshot of the table.

   ```
   CREATE TABLE backup_table AS
   SELECT * FROM original_table;
   ```

3. **Performance Optimization**:

   - **Scenario**: You need to optimize query performance by creating a pre-aggregated or pre-joined table.

   - **Solution**: Use CTAS to create a new table that stores the results of complex queries.

   ```
   CREATE TABLE optimized_table AS
   SELECT a.*, b.*
   FROM table_a a
   JOIN table_b b ON a.id = b.id;
   ```

4. **Data Subsetting**:

   - **Scenario**: You need to create a subset of data for testing or analysis.

- **Solution**: Use CTAS to create a new table with a subset of the data.

```
CREATE TABLE subset_table AS
SELECT * FROM large_table
WHERE condition;
```

## Benefits of Using CTAS in Databricks

- **Efficiency**: CTAS combines the creation of a table and the insertion of data into a single operation, making it efficient.

- **Simplicity**: It simplifies the process of creating new tables based on complex queries.

- **Flexibility**: CTAS can be used for a wide range of data management tasks, from simple data copying to complex transformations.

## Example in Databricks

Here's an example of using CTAS in Databricks to create a new Delta table:

```
CREATE TABLE delta.`/mnt/delta/new_table` AS
SELECT * FROM existing_table;
```

This command creates a new Delta table at the specified location with the data from `existing_table`.

In summary, CTAS is a powerful and flexible tool in Databricks for creating new tables based on the results of `SELECT` queries, making it a valuable solution for data transformation, backup, optimization, and subsetting.

▼ **61.Create a generated column.**

Creating a generated column in Databricks involves using Delta Lake, which supports generated columns. These columns are automatically computed based on a user-specified function over other columns in the table. Here's how you can create a generated column in Databricks:

## Example: Creating a Generated Column

Let's say you have a table `people` with a `birthDate` column, and you want to create a `dateOfBirth` column that extracts the date from the `birthDate` timestamp.

1. **Create the Table with a Generated Column:**

```
CREATE TABLE default.people (
    id INT,
    firstName STRING,
```

```
    middleName STRING,
    lastName STRING,
    gender STRING,
    birthDate TIMESTAMP,
    dateOfBirth DATE GENERATED ALWAYS AS (CAST(birthDate AS DATE)),
    ssn STRING,
    salary INT
);
```

In this example, the `dateOfBirth` column is a generated column that automatically computes its value by casting the `birthDate` timestamp to a date

## Key Points

- **Generated columns** can use most built-in SQL functions in Spark.

- They are stored as normal columns, meaning they occupy storage space.

- You can create Delta tables with generated columns using SQL, Scala, Java, or Python APIs

▼ **62.Add a table comment.**

Adding a comment to a table in Databricks is straightforward. You can use the `COMMENT ON` statement to add or modify comments on tables. Here's how you can do it:

## Example: Adding a Comment to a Table

1. **Add a Comment to a Table:**

```
COMMENT ON TABLE my_table IS 'This is my table comment';
```

1. **Remove a Comment from a Table:**

```
COMMENT ON TABLE my_table IS NULL;
```

## Key Points

- **Syntax:**

```
COMMENT ON TABLE table_name IS 'your_comment';
```

- **Removing a Comment:**

```
COMMENT ON TABLE table_name IS NULL;
```

- You need appropriate privileges to add or modify comments on a table.

## Additional Information

- You can also add comments to other objects like catalogs, schemas, and columns using similar syntax.
- For AI-generated comments, Databricks provides features to automatically generate and manage comments for tables and columns

▼ **63.Use CREATE OR REPLACE TABLE and INSERT OVERWRITE**

Databricks, you can use the `CREATE OR REPLACE TABLE` statement to create a new table or replace an existing one, and the `INSERT OVERWRITE` statement to overwrite data in a table. Here's how you can use both:

## Example: Using `CREATE OR REPLACE TABLE`

1. **Create or Replace a Table:**

```
CREATE OR REPLACE TABLE my_table (
    id INT,
    name STRING,
    age INT
);
```

This statement will create a new table named `my_table` if it doesn't exist, or replace the existing table with the same name.

## Example: Using `INSERT OVERWRITE`

1. **Insert Overwrite Data in a Table:**

```
INSERT OVERWRITE TABLE my_table
SELECT id, name, age
FROM another_table
WHERE age > 30;
```

This statement will overwrite the data in `my_table` with the results of the `SELECT` query from `another_table` where the age is greater than 30

## Key Points

- `CREATE OR REPLACE TABLE` : This command is useful for ensuring that a table is created if it doesn't exist or replaced if it does.
- `INSERT OVERWRITE` : This command replaces the existing data in the table with the new data specified by the query. It can be used with or without partitions.

## ▼ 64.Compare and contrast CREATE OR REPLACE TABLE and INSERT OVERWRITE

compare and contrast the `CREATE OR REPLACE TABLE` and `INSERT OVERWRITE` statements in Databricks:

`CREATE OR REPLACE TABLE`

**Purpose:**

- Creates a new table if it doesn't exist or replaces an existing table with the same name.

**Usage:**

- Used when you want to define the schema and structure of a table, and optionally populate it with initial data.

**Syntax:**

```
CREATE OR REPLACE TABLE table_name (
    column1 data_type,
    column2 data_type,
    ...
);
```

**Key Points:**

- **Schema Definition:** Allows you to define the table schema, including column names and data types.
- **Table Replacement:** If the table already exists, it is dropped and recreated, which means all existing data is lost.
- **Flexibility:** Useful for creating tables from scratch or redefining the schema of existing tables.

`INSERT OVERWRITE`

**Purpose:**

- Overwrites the existing data in a table or partition with new data.

**Usage:**

- Used when you want to replace the data in a table or specific partitions without changing the table schema.

**Syntax:**

```
INSERT OVERWRITE TABLE table_name
SELECT column1, column2, ...
FROM source_table
WHERE condition;
```

**Key Points:**

- **Data Replacement:** Replaces the existing data in the table or specified partitions with the results of a `SELECT` query.

- **Schema Preservation:** The table schema remains unchanged; only the data is overwritten.

- **Efficiency:** Useful for updating large datasets or partitions with new data without altering the table structure.

## Comparison

| Feature | CREATE OR REPLACE TABLE | INSERT OVERWRITE |
|---|---|---|
| **Primary Function** | Create or replace a table | Overwrite data in a table or partition |
| **Schema Definition** | Yes | No |
| **Data Overwrite** | Yes (all data is lost if table exists) | Yes (only data is overwritten) |
| **Use Case** | Defining or redefining table schema | Replacing data in existing tables/partitions |
| **Flexibility** | High (can change schema) | Moderate (schema remains unchanged) |

## Example Use Cases

- `CREATE OR REPLACE TABLE` : When you need to redefine the schema of a table, such as adding or removing columns, and you are okay with losing existing data.

- `INSERT OVERWRITE` : When you need to update the data in a table or partition, such as refreshing a dataset with new values, while keeping the table schema intact.

▼ **65.Identify a scenario in which MERGE should be used.**

The `MERGE` statement in Databricks is particularly useful for scenarios where you need to perform **upserts** (a combination of updates and inserts) or synchronize data between tables. Here's a common scenario where `MERGE` should be used:

## Scenario: Data Synchronization

Imagine you have a **target Delta table** that stores customer information and a **source table** that contains new and updated customer records. You want to update existing customer records in the target table with the new information from the source table and insert new customer records that do not exist in the target table.

## Example: Using `MERGE` for Data Synchronization

1. **Target Table (** `customers` **):**

   - Contains existing customer records.

2. **Source Table (** `new_customers` **):**

   - Contains new and updated customer records.

3. **MERGE Statement**:

```
MERGE INTO customers AS target
USING new_customers AS source
ON target.customer_id = source.customer_id
WHEN MATCHED THEN
  UPDATE SET
    target.name = source.name,
    target.email = source.email,
    target.phone = source.phone
WHEN NOT MATCHED THEN
  INSERT (customer_id, name, email, phone)
  VALUES (source.customer_id, source.name, source.email, source.phone);
```

## Key Points

- **WHEN MATCHED**: Updates existing records in the target table with the data from the source table.

- **WHEN NOT MATCHED**: Inserts new records from the source table into the target table.

- **Data Consistency**: Ensures that the target table is synchronized with the latest data from the source table.

## Benefits of Using `MERGE`

- **Efficiency**: Combines multiple operations (insert, update, delete) into a single statement, reducing the complexity and improving performance.

- **Atomicity**: Ensures that the operations are performed atomically, maintaining data integrity.

- **Flexibility**: Supports complex conditions and actions, making it suitable for various data integration and synchronization tasks

▼ **66.Identify MERGE as a command to deduplicate data upon writing.**

The `MERGE` statement in Databricks is indeed a powerful command for deduplicating data upon writing. It allows you to combine insert, update, and delete operations into a single statement, making it ideal for maintaining data integrity and avoiding duplicates.

## Scenario: Deduplicating Data Upon Writing

Imagine you have a **target Delta table** that may contain duplicate records, and you want to ensure that only unique records are maintained when new data is added. You can use the `MERGE` statement to achieve this by specifying conditions to identify duplicates and handle them appropriately.

## Example: Using `MERGE` to Deduplicate Data

1. **Target Table (** `customers` **):**

   - Contains existing customer records, potentially with duplicates.

2. **Source Table (** `new_customers` **):**

   - Contains new customer records that need to be merged into the target table.

3. **MERGE Statement:**

```
MERGE INTO customers AS target
USING new_customers AS source
ON target.customer_id = source.customer_id
WHEN MATCHED THEN
  UPDATE SET
    target.name = source.name,
    target.email = source.email,
    target.phone = source.phone
WHEN NOT MATCHED THEN
  INSERT (customer_id, name, email, phone)
  VALUES (source.customer_id, source.name, source.email, source.phone);
```

## Key Points

- **Deduplication:** The `MERGE` statement ensures that if a record with the same `customer_id` exists in both the target and source tables, the target table is updated with the new data from the source table. If no matching record is found, a new record is inserted.

- **Efficiency:** Combines multiple operations into a single statement, reducing the complexity and improving performance.

- **Data Integrity:** Maintains data integrity by ensuring that only unique records are present in the target table.

## Benefits of Using `MERGE` for Deduplication

- **Atomic Operations:** Ensures that the deduplication process is atomic, meaning all changes are applied consistently.

- **Flexibility:** Allows for complex conditions and actions, making it suitable for various data integration and deduplication tasks.

- **Scalability:** Efficiently handles large datasets, making it ideal for big data environments.

▼ **67.Describe the benefits of the MERGE command**

he `MERGE` command in Databricks offers several significant benefits, making it a powerful tool for data management and analysis. Here are some key advantages:

## 1. Consolidation of Data

- **Combines Multiple Operations**: The `MERGE` command allows you to perform updates, insertions, and deletions in a single statement. This is particularly useful for consolidating data from different sources into a single dataset.
- **Simplifies Data Integration**: By merging data from various sources, it eliminates the need for multiple steps and manual data manipulation, streamlining the data integration process.

## 2. Efficiency and Performance

- **Optimized Performance**: Databricks has an optimized implementation of `MERGE` that reduces the number of shuffle operations, improving performance for common workloads.
- **Low Shuffle Merge**: This feature processes unmodified rows separately, reducing the amount of shuffled data and enhancing performance. It also helps maintain the data layout, which benefits subsequent operations.

## 3. Data Integrity and Consistency

- **Atomic Operations**: The `MERGE` command ensures that all operations (insert, update, delete) are performed atomically, maintaining data integrity and consistency.
- **Accurate Data Matching**: It performs matching and combining of records based on specified conditions, ensuring that data is accurately updated or inserted.

## 4. Flexibility and Scalability

- **Complex Conditions**: Supports complex conditions and actions, making it suitable for various data synchronization and deduplication tasks.
- **Scalability**: Efficiently handles large datasets, making it ideal for big data environments and ensuring that operations are scalable.

## Use Cases

- **Data Synchronization**: Keeping a target table in sync with a source table by updating existing records and inserting new ones.
- **Deduplication**: Ensuring that only unique records are maintained in a table by merging new data and removing duplicates.
- **Data Warehousing**: Integrating and updating data from multiple sources into a data warehouse.

▼ **68.Identify why a COPY INTO statement is not duplicating data in the target table.**

The `COPY INTO` statement in Databricks is designed to be idempotent, meaning it should not duplicate data in the target table if used correctly. Here are some reasons why `COPY INTO` might not be duplicating data:

## Reasons for No Data Duplication

1. **Idempotent Operation**:

   - `COPY INTO` is designed to skip files that have already been loaded, even if they have been modified since the last load. This ensures that each file is only loaded once, preventing duplication.

2. **File Metadata Tracking**:

   - Databricks tracks the metadata of files that have been loaded into the table. If a file with the same name and path is encountered again, it is skipped.

3. **Correct Usage**:

   - If the initial data load and subsequent incremental loads are done correctly using `COPY INTO`, the command ensures that no duplicate records are inserted.

## Example Usage

```
COPY INTO target_table
FROM 'path/to/source/files'
FILEFORMAT = PARQUET;
```

## Common Issues and Solutions

- **Overlapping Data**: If incremental load files contain overlapping data with previously loaded files, ensure that the files are mutually exclusive to avoid potential duplicates.

- **Initial Load Method**: Ensure that the initial data load is also done using `COPY INTO` rather than other methods like `CREATE TABLE AS SELECT`, which might not track file metadata in the same way.

## Key Points

- **Idempotency**: Ensures that files are only loaded once.

- **Metadata Tracking**: Prevents reloading of the same files.

- **Proper Usage**: Ensures no duplication if used correctly from the start.

▼ **69.Identify a scenario in which COPY INTO should be used.**

The `COPY INTO` statement in Databricks is particularly useful for efficiently loading data from external file sources into Delta tables. Here's a scenario where `COPY INTO` should be used:

## Scenario: Incremental Data Loading from Cloud Storage

Imagine you have a data pipeline that regularly receives new data files in a cloud storage location (e.g., AWS S3, Azure Data Lake Storage, Google Cloud Storage). You want to incrementally load these new files into a Delta table in Databricks without duplicating data.

## Example: Using `COPY INTO` for Incremental Data Loading

1. **Source Files**:

   - New data files are regularly added to a cloud storage location, such as `s3://my-bucket/new-data/`.

2. **Target Table (`my_table`)**:

   - A Delta table in Databricks where the data will be loaded.

3. **COPY INTO Statement**:

```
COPY INTO my_table
FROM 's3://my-bucket/new-data/'
FILEFORMAT = PARQUET
COPY_OPTIONS ('mergeSchema' = 'true');
```

## Key Points

- **Idempotent Operation**: `COPY INTO` ensures that files already loaded are skipped, preventing data duplication.

- **Schema Evolution**: The `mergeSchema` option allows the schema to evolve automatically as new columns are added in the source files.

- **File Format Support**: Supports various file formats such as Parquet, CSV, JSON, and more.

- **Cloud Storage Integration**: Easily integrates with cloud storage services, making it ideal for cloud-based data pipelines.

## Benefits

- **Efficiency**: Simplifies the process of loading new data incrementally without manual intervention.

- **Scalability**: Handles large volumes of data efficiently, making it suitable for big data environments.
- **Data Integrity**: Ensures that each file is processed exactly once, maintaining data integrity.

▼ **70.Use COPY INTO to insert data.**

`COPY INTO` statement to insert data into a Delta table in Databricks, follow these steps:

## Example: Using `COPY INTO` to Insert Data

1. **Create the Target Table** (if it doesn't already exist):

```
CREATE TABLE IF NOT EXISTS my_table (
    id INT,
    name STRING,
    age INT
);
```

1. **Use `COPY INTO` to Load Data from a Source Location**:

```
COPY INTO my_table
FROM 's3://my-bucket/new-data/'
FILEFORMAT = PARQUET
COPY_OPTIONS ('mergeSchema' = 'true');
```

## Key Points

- **Source Location**: The `FROM` clause specifies the path to the source files. This can be a cloud storage location like S3, Azure Data Lake Storage, or Google Cloud Storage.
- **File Format**: The `FILEFORMAT` clause specifies the format of the source files (e.g., PARQUET, CSV, JSON).
- **Schema Evolution**: The `COPY_OPTIONS` clause with `mergeSchema` set to `true` allows the schema to evolve automatically as new columns are added in the source files.

## Benefits

- **Idempotent Operation**: Ensures that files already loaded are skipped, preventing data duplication.
- **Efficiency**: Simplifies the process of loading new data incrementally without manual intervention.
- **Scalability**: Handles large volumes of data efficiently, making it suitable for big data environments.

▼ **71.Identify the components necessary to create a new DLT pipeline.**

Creating a new Delta Live Tables (DLT) pipeline in Databricks involves several key components and steps. Here's a breakdown of what you need:

## Components Necessary to Create a New DLT Pipeline

1. **Databricks Workspace**:

   - Ensure you have access to a Databricks workspace where you can create and manage your DLT pipelines.

2. **Cluster**:

   - You need permission to create clusters or access to a cluster policy that defines a Delta Live Tables cluster. The DLT runtime creates a cluster before running your pipeline.

3. **Data Source**:

   - Identify the source data you want to process. This could be data stored in cloud storage (e.g., AWS S3, Azure Data Lake Storage) or other databases.

4. **Pipeline Configuration**:

   - **Pipeline Name**: Provide a unique name for your pipeline.

   - **Serverless Option**: Optionally, you can check the box for serverless to simplify resource management.

   - **Catalog and Schema**: Select a catalog to publish data and specify a schema within that catalog.

5. **ETL Code**:

   - Write the ETL (Extract, Transform, Load) code using Python or SQL. This code will define how data is ingested, transformed, and loaded into the target tables.

6. **Notebooks**:

   - Create notebooks to write and manage your ETL code. These notebooks will be used to declare materialized views and streaming tables.

7. **Permissions**:

   - Ensure you have the necessary permissions to create schemas, tables, and manage clusters within your Databricks workspace.

## Steps to Create a New DLT Pipeline

1. **Create a Pipeline**:

   - Navigate to Delta Live Tables in the Databricks sidebar and click on "Create Pipeline".

   - Provide the necessary configuration details such as pipeline name, serverless option, catalog, and schema.

2. **Write ETL Code**:

- Use Databricks notebooks to write your ETL code in Python or SQL. This code will define the data transformations and loading processes.

3. **Start the Pipeline**:

   - Once the pipeline is configured and the ETL code is ready, start the pipeline to begin processing data.

## Example Configuration

```sql
CREATE TABLE IF NOT EXISTS my_table (
    id INT,
    name STRING,
    age INT
);


COPY INTO my_table
FROM 's3://my-bucket/new-data/'
FILEFORMAT = PARQUET
COPY_OPTIONS ('mergeSchema' = 'true');
```

This example shows how to create a table and use the `COPY INTO` statement to load data from an S3 bucket.

▼ **72.Identify the purpose of the target and of the notebook libraries in creating a pipeline.**

When creating a Delta Live Tables (DLT) pipeline in Databricks, the **target** and **notebook libraries** play crucial roles in the pipeline's functionality and organization.

## Purpose of the Target

The **target** in a DLT pipeline specifies where the tables created by your pipeline will be published. This is particularly important when you move beyond development and testing phases. Here's why it's essential:

- **Data Availability**: Publishing tables to a target makes them available for querying and analysis elsewhere in your Databricks environment.

- **Organization**: Helps in organizing your data assets by specifying a catalog and schema, ensuring that your tables are stored in a structured and accessible manner.

- **Consistency**: Ensures that all tables created by the pipeline are consistently published to the same location, which is crucial for data governance and management.

## Purpose of Notebook Libraries

**Notebook libraries** are used to write and manage the ETL (Extract, Transform, Load) code that defines the data transformations and loading processes in your pipeline. Here's why they are important:

- **Development and Validation**: Notebooks allow you to develop and validate your pipeline code interactively. You can test your transformations and logic step-by-step before deploying them.

- **Modularity**: By using notebooks, you can modularize your ETL code, making it easier to manage, debug, and maintain.

- **Documentation**: Notebooks can include markdown cells for documentation, making it easier to understand the purpose and functionality of different parts of your pipeline.

- **Flexibility**: You can use notebooks to write your ETL code in Python or SQL, depending on your preference and the requirements of your pipeline.

## Example Workflow

1. **Define the Target**:

   - Specify the catalog and schema where your tables will be published.

   ```
   CREATE TABLE IF NOT EXISTS my_catalog.my_schema.my_table (
       id INT,
       name STRING,
       age INT
   );
   ```

2. **Develop ETL Code in Notebooks**:

   - Write your data transformation logic in a Databricks notebook.

   ```
   # Example Python code in a notebook
   df = spark.read.format("csv").option("header", "true").load("/path/to/source.csv")
   df_transformed = df.withColumn("age", df["age"].cast("int"))
   df_transformed.write.format("delta").saveAsTable("my_catalog.my_schema.my_table")
   ```

By understanding and utilizing the target and notebook libraries effectively, you can create robust and maintainable data pipelines in Databricks.

▼ **73.Compare and contrast triggered and continuous pipelines in terms of cost and latency**

Let's compare and contrast **triggered** and **continuous** pipelines in Databricks, focusing on cost and latency:

## Triggered Pipelines

**Cost:**

- **Lower Cost**: Triggered pipelines generally incur lower costs because the cluster runs only for the duration needed to process the data and then stops. This reduces the overall compute time and resource usage.

**Latency:**

- **Higher Latency**: Since triggered pipelines process data at scheduled intervals (e.g., every 10 minutes, hourly, or daily), there is a delay between data arrival and processing. This makes them less suitable for scenarios requiring real-time data updates.

## Continuous Pipelines

**Cost:**

- **Higher Cost**: Continuous pipelines require an always-running cluster, which increases the cost due to constant resource usage. This is necessary to process data as it arrives, ensuring that the pipeline is always up-to-date.

**Latency:**

- **Lower Latency**: Continuous pipelines process data in near real-time, with updates occurring as frequently as every few seconds to minutes. This makes them ideal for applications that need immediate data freshness and low-latency processing.

## Comparison Table

| Feature | Triggered Pipelines | Continuous Pipelines |
|---|---|---|
| **Cost** | Lower (cluster runs only during updates) | Higher (cluster runs continuously) |
| **Latency** | Higher (scheduled intervals) | Lower (real-time processing) |
| **Use Case** | Batch processing, periodic updates | Real-time analytics, streaming data |
| **Resource Usage** | Efficient, runs only when needed | Constant, always running |

## Use Cases

- **Triggered Pipelines**: Suitable for batch processing tasks where data freshness is not critical, such as daily reports or periodic data aggregation.
- **Continuous Pipelines**: Ideal for real-time analytics, monitoring systems, and applications that require immediate data updates, such as fraud detection or live dashboards.

▼ **74.Identify which source location is utilizing Auto Loader.**

Auto Loader in Databricks can load data files from several cloud storage sources. Here are the supported source locations:

1. **Amazon S3** (`s3://`)

2. **Azure Data Lake Storage Gen2** (`abfss://`)

3. **Google Cloud Storage** (`gs://`)

4. **Azure Blob Storage** (`wasbs://`)

5. **Azure Data Lake Storage Gen1** (`adl://`) - Note that this is being deprecated.

6. **Databricks File System** (`dbfs:/`)

Auto Loader is designed to incrementally and efficiently process new data files as they arrive in these cloud storage locations, making it ideal for handling large volumes of data in near real-time.

▼ **75.Identify a scenario in which Auto Loader is beneficial**

Auto Loader in Databricks is particularly beneficial in scenarios where you need to efficiently and incrementally process large volumes of data from cloud storage. Here's a common scenario where Auto Loader shines:

## Scenario: Real-Time Data Ingestion for IoT Devices

Imagine you are managing a fleet of IoT devices that continuously generate data and store it in a cloud storage location, such as AWS S3 or Azure Data Lake Storage. You need to process this data in near real-time to monitor device performance, detect anomalies, and trigger alerts.

## Benefits of Using Auto Loader

1. **Incremental Data Processing**:

   - Auto Loader can automatically detect and process new files as they arrive in the cloud storage, ensuring that your data pipeline is always up-to-date without manual intervention.

2. **Scalability**:

   - It can handle billions of files and scale to support near real-time ingestion of millions of files per hour, making it ideal for high-volume data sources like IoT devices.

3. **Schema Evolution**:

   - Auto Loader supports schema inference and evolution, allowing it to adapt to changes in the data schema over time without requiring manual updates.

4. **Fault Tolerance**:

   - It provides exactly-once processing guarantees by tracking the ingestion progress and ensuring that each file is processed only once, even in the event of failures.

## Example Workflow

1. **Configure Auto Loader**:

   - Set up Auto Loader to monitor the cloud storage location where IoT data is stored.

   ```
   df = (spark.readStream
         .format("cloudFiles")
         .option("cloudFiles.format", "json")
         .load("s3://my-bucket/iot-data/"))
   ```

2. **Process Data**:

   - Define the transformations and actions to process the incoming data.

   ```
   processed_df = df.withColumn("timestamp", current_timestamp())
   ```

3. **Write to Delta Table**:

   - Write the processed data to a Delta table for further analysis.

   ```
   (processed_df.writeStream
     .format("delta")
     .option("checkpointLocation", "/path/to/checkpoint")
     .start("/path/to/delta-table"))
   ```

By using Auto Loader, you can ensure that your data pipeline is robust, scalable, and capable of handling real-time data ingestion efficiently.

▼ **76.Identify why Auto Loader has inferred all data to be STRING from a JSON source**

Auto Loader in Databricks infers all data to be `STRING` from a JSON source by default because JSON does not inherently encode data types. Here are the key reasons:

1. **Default Behavior**:

   - For formats like JSON and CSV, Auto Loader infers all columns as strings by default. This includes nested fields in JSON files.

2. **Schema Inference Configuration**:

   - The default setting for schema inference in Auto Loader is to treat all columns as strings unless explicitly configured otherwise. This is controlled by the `cloudFiles.inferColumnTypes` option.

## How to Address This

To ensure that Auto Loader infers the correct data types, you can enable the `cloudFiles.inferColumnTypes` option:

```
df = (spark.readStream
    .format("cloudFiles")
    .option("cloudFiles.format", "json")
    .option("cloudFiles.inferColumnTypes", "true")
    .load("s3://my-bucket/json-data/"))
```

## Key Points

- **Schema Evolution**: Enabling `cloudFiles.inferColumnTypes` allows Auto Loader to infer the actual data types of columns, which can help in maintaining the correct schema as new data arrives.

- **Flexibility**: This option provides flexibility in handling different data types and ensures that your data is accurately represented in the Delta table.

▼ **77.Identify the default behavior of a constraint violation**

In Databricks, the default behavior when a constraint is violated is to fail the transaction with an error. This ensures that data integrity is maintained by preventing invalid data from being written to the table.

## Types of Constraints and Their Enforcement

1. **NOT NULL Constraint**:

   - Ensures that a column cannot have null values.

   - If a null value is inserted into a column with a NOT NULL constraint, the transaction fails.

2. **CHECK Constraint**:

   - Ensures that a specified condition is true for each row in the table.

   - If a row violates the CHECK constraint, the transaction fails.

## Example

```
CREATE TABLE people (
    id INT NOT NULL,
    name STRING,
    age INT CHECK (age > 0)
);

-- Attempting to insert invalid data
INSERT INTO people (id, name, age) VALUES (1, 'John Doe', -5);
-- This will fail because the age is not greater than 0
```

In this example, the insertion fails because the `age` value violates the CHECK constraint.

## ▼ 78.Identify the impact of ON VIOLATION DROP ROW and ON VIOLATION FAIL UPDATE for a constraint violation

In Databricks, the `ON VIOLATION` clause can be used to specify the action to take when a data quality constraint is violated. Here's a comparison of the impacts of `ON VIOLATION DROP ROW` and `ON VIOLATION FAIL UPDATE`:

`ON VIOLATION DROP ROW`

**Impact:**

- **Data Handling**: Rows that violate the specified constraint are excluded from the target table. This means that any invalid records are simply dropped and not written to the table.

- **Data Quality**: Ensures that only valid data is written to the table, maintaining high data quality.

- **Metrics**: The number of dropped rows is recorded as part of the data quality metrics, which can be monitored.

- **Use Case**: Useful when you want to ensure that only clean data is loaded, and you can afford to lose invalid records.

`ON VIOLATION FAIL UPDATE`

**Impact:**

- **Data Handling**: The entire update operation fails if any row violates the specified constraint. No data is written to the table until the violation is resolved.

  2

- **Data Quality**: Ensures strict adherence to data quality rules by preventing any invalid data from being written.

- **Error Handling**: Requires manual intervention to correct the data or the constraint before the update can be retried.

- **Use Case**: Suitable for scenarios where data integrity is critical, and you cannot afford to have any invalid data in the table.

## Example

```
-- Using ON VIOLATION DROP ROW
CREATE TABLE people (
    id INT,
    name STRING,
```

```
    age INT,
    CONSTRAINT valid_age CHECK (age > 0) ON VIOLATION DROP ROW
);

-- Using ON VIOLATION FAIL UPDATE
CREATE TABLE people (
    id INT,
    name STRING,
    age INT,
    CONSTRAINT valid_age CHECK (age > 0) ON VIOLATION FAIL UPDATE
);
```

## Summary

- `ON VIOLATION DROP ROW` : Drops invalid rows, ensuring only valid data is written, but potentially losing some data.

- `ON VIOLATION FAIL UPDATE` : Fails the entire update operation if any row is invalid, ensuring strict data integrity but requiring manual intervention to resolve issues.

▼ **79.Explain change data capture and the behavior of APPLY CHANGES INTO**

## Change Data Capture (CDC)

**Change Data Capture (CDC)** is a process that identifies and captures incremental changes (inserts, updates, and deletes) in a database. This allows you to track changes in real-time and apply them to another system, ensuring data consistency and enabling real-time analytics. CDC is crucial for scenarios like synchronizing databases, updating data warehouses, and maintaining data lakes

## Behavior of `APPLY CHANGES INTO`

In Databricks, the `APPLY CHANGES INTO` statement is used within Delta Live Tables (DLT) to simplify the CDC process. Here's how it works:

1. **Source Data**: It processes changes from a change data feed (CDF) or snapshots of the source data.

2. **Sequencing**: You must specify a column in the source data to sequence records, ensuring proper ordering.

3. **Handling Out-of-Order Data**: Automatically manages out-of-sequence records, eliminating the need for complex logic to handle such cases.

4. **SCD Types**: Supports Slowly Changing Dimensions (SCD) Type 1 and Type 2:
   - **SCD Type 1**: Updates records directly without retaining history.

- **SCD Type 2**: Retains a history of records, either for all updates or specific columns.

## Example Usage

```
APPLY CHANGES INTO target_table
FROM source_table
KEYS (id)
SEQUENCE BY timestamp
STORED AS SCD TYPE 2;
```

## Benefits

- **Simplifies CDC**: Reduces the complexity of implementing CDC by handling out-of-order data and providing built-in support for SCD types.

- **Real-Time Processing**: Enables near real-time data updates, making it ideal for applications that require up-to-date information.

- **Data Integrity**: Ensures accurate and consistent data synchronization between source and target tables

▼ **80.Query the events log to get metrics, perform audit log in, examine lineage.**

To query the events log for metrics, perform audit logging, and examine lineage in Databricks, you can use the following steps:

## 1. Querying the Events Log for Metrics

Databricks provides system tables that you can query to get metrics related to your Delta Live Tables (DLT) pipelines. These tables include information about pipeline progress, data quality checks, and other metrics.

**Example Query:**

```
SELECT * FROM delta_live_tables.event_log
WHERE pipeline_id = 'your_pipeline_id';
```

This query retrieves all events related to a specific pipeline, allowing you to monitor its progress and performance

## 2. Performing Audit Logging

Audit logs in Databricks capture detailed usage patterns and events, which are essential for monitoring and compliance. To access and query audit logs, you need to configure audit log delivery to a storage location like DBFS (Databricks File System).

**Steps to Configure Audit Logging:**

1. **Enable Audit Logging**: Follow the instructions to enable audit logging in your Databricks workspace.

2. **Query Audit Logs**: Once enabled, you can query the audit logs stored in DBFS.

**Example Query:**

```
SELECT * FROM dbfs:/path/to/audit/logs
WHERE event_type = 'user_login';
```

This query retrieves all user login events from the audit logs

## 3. Examining Lineage

Data lineage in Databricks helps you understand the flow of data through your pipelines, including transformations and dependencies. The Delta Live Tables event log contains lineage information.

**Example Query:**

```
SELECT * FROM delta_live_tables.lineage_log
WHERE pipeline_id = 'your_pipeline_id';
```

This query retrieves lineage information for a specific pipeline, allowing you to trace data transformations and dependencies

## Summary

- **Metrics**: Use the `event_log` table to query pipeline metrics.
- **Audit Logging**: Enable and query audit logs to monitor usage and compliance.
- **Lineage**: Use the `lineage_log` table to examine data lineage.

▼ **81.Troubleshoot DLT syntax: Identify which notebook in a DLT pipeline produced an error, identify the need for LIVE in create statement, identify the need for STREAM in from clause.**

To troubleshoot Delta Live Tables (DLT) syntax issues, you can follow these steps:

## 1. Identify Which Notebook Produced an Error

When a DLT pipeline encounters an error, you can use the Databricks UI to identify the specific notebook and cell that caused the issue:

- **View Pipeline Events**: Navigate to the Delta Live Tables UI and check the event log for your pipeline. This log provides detailed information about each step, including errors.

- **Notebook Connection**: If your notebook is connected to the pipeline, you can view the pipeline's dataflow graph and event log directly from the notebook.

## 2. Identify the Need for `LIVE` in `CREATE` Statement

In Delta Live Tables, the `LIVE` keyword is used to reference tables and views managed by the pipeline. It ensures that the pipeline correctly tracks dependencies and updates.

**Example:**

```
CREATE LIVE TABLE my_table AS
SELECT * FROM source_table;
```

- **Purpose**: The `LIVE` keyword indicates that `my_table` is a managed table within the DLT pipeline.

### 3. Identify the Need for `STREAM` in `FROM` Clause

When working with streaming data, you need to use the `STREAM` keyword in the `FROM` clause to indicate that the source is a streaming table.

**Example:**

```
CREATE LIVE TABLE my_streaming_table AS
SELECT * FROM STREAM(LIVE.source_stream);
```

- **Purpose**: The `STREAM` keyword tells the pipeline to treat `source_stream` as a streaming source, enabling continuous data processing.

## Summary

- **Error Identification**: Use the Delta Live Tables UI or notebook connection to identify which notebook and cell produced an error.

- `LIVE` **Keyword**: Required in `CREATE` statements to reference managed tables and views within the pipeline.

- `STREAM` **Keyword**: Required in the `FROM` clause to indicate streaming sources.

# Section 4: Production Pipelines(16%)

▼ **82.Identify the benefits of using multiple tasks in Job**

Using multiple tasks in a Databricks Job offers several benefits, enhancing both efficiency and manageability of data workflows. Here are the key advantages:

## 1. Improved Workflow Organization

- **Modularity:** Breaking down a complex job into smaller, manageable tasks allows for better organization and modularity. Each task can focus on a specific part of the workflow, making it easier to develop, debug, and maintain.

- **Clear Dependencies**: Tasks can be organized in a Directed Acyclic Graph (DAG), clearly defining dependencies and execution order. This ensures that tasks are executed in the

correct sequence.

## 2. Enhanced Performance and Efficiency

- **Parallel Execution**: Multiple tasks can run in parallel, significantly reducing the overall execution time of the job. This is particularly beneficial for large-scale data processing and machine learning workflows.

- **Cluster Reuse**: Tasks within a job can share the same cluster, reducing the overhead associated with spinning up new clusters for each task. This leads to faster job execution and lower costs.

## 3. Cost Savings

- **Resource Optimization**: By reusing clusters across tasks, you can optimize resource utilization and reduce the costs associated with cluster startup and underutilization.

- **Efficient Scaling**: Databricks Jobs can scale efficiently by leveraging multiple tasks, ensuring that resources are used effectively without unnecessary expenditure.

## 4. Simplified Monitoring and Management

- **Unified Monitoring**: The Databricks UI provides a unified view of all tasks within a job, making it easier to monitor progress, identify bottlenecks, and troubleshoot issues.

- **Automated Alerts**: You can set up automated alerts for task failures or completions, ensuring that you are promptly notified of any issues that need attention.

## Example Use Case

Imagine you have a data pipeline that involves extracting data from multiple sources, transforming it, and then loading it into a data warehouse. By using multiple tasks, you can:

1. **Extract Data**: Separate tasks for extracting data from different sources.

2. **Transform Data**: Tasks for cleaning and transforming the data.

3. **Load Data**: Tasks for loading the transformed data into the data warehouse.

This modular approach ensures that each step is handled efficiently and can be monitored independently.

▼ **83. Set up a predecessor task in Jobs.**

To set up a predecessor task in Databricks Jobs, you need to define task dependencies. This ensures that a task runs only after its predecessor tasks have successfully completed. Here's how you can do it:

## Steps to Set Up a Predecessor Task

1. **Create or Configure a Job**:

   - Navigate to the Databricks workspace and click on **Workflows** in the sidebar.

- Click on the job name to open the job configuration.

2. **Add Tasks**:

   - Click on the **Tasks** tab to view the task graph.

   - Add a new task by clicking the **+** button.

3. **Define Task Dependencies**:

   - When configuring a task, use the `depends_on` field to specify predecessor tasks.

## Example Configuration

Here's an example of how to set up a job with two tasks where the second task depends on the first:

```yaml
jobs:
  - name: my-job
    tasks:
      - task_key: task_1
        notebook_task:
          notebook_path: /path/to/notebook1
      - task_key: task_2
        depends_on:
          - task_key: task_1
        notebook_task:
          notebook_path: /path/to/notebook2
```

In this example:

- **Task 1** ( `task_1` ) runs a notebook located at `/path/to/notebook1` .

- **Task 2** ( `task_2` ) runs a notebook located at `/path/to/notebook2` and depends on the successful completion of `task_1` .

## Benefits of Using Task Dependencies

- **Sequential Execution**: Ensures tasks are executed in the correct order.

- **Error Handling**: Prevents downstream tasks from running if a predecessor task fails.

- **Modularity:** Allows you to break down complex workflows into smaller, manageable tasks.

▼ **84.Identify a scenario in which a predecessor task should be set up.**

Setting up a predecessor task in Databricks Jobs is essential when you need to ensure that certain tasks are completed before others can start. Here's a scenario where this is particularly useful:

## Scenario: ETL Pipeline with Data Validation

Imagine you have an ETL (Extract, Transform, Load) pipeline that involves multiple steps, including data extraction, transformation, validation, and loading into a data warehouse. You want to ensure that data is validated before it is loaded into the final destination.

## Steps in the Pipeline

1. **Data Extraction**:

   - Extract data from various sources (e.g., databases, APIs, cloud storage).

2. **Data Transformation**:

   - Clean and transform the extracted data to match the required schema.

3. **Data Validation**:

   - Validate the transformed data to ensure it meets quality standards (e.g., no missing values, correct data types).

4. **Data Loading**:

   - Load the validated data into the data warehouse.

## Setting Up Predecessor Tasks

In this scenario, you would set up predecessor tasks to ensure that each step is completed before the next one begins:

1. **Task 1: Data Extraction**:

   - Extract data from sources.

     ```
     tasks:
       - task_key: extract_data
         notebook_task:
           notebook_path: /path/to/extract_notebook
     ```

2. **Task 2: Data Transformation**:

   - Transform the extracted data.

     ```
     tasks:
       - task_key: transform_data
         depends_on:
           - task_key: extract_data
         notebook_task:
           notebook_path: /path/to/transform_notebook
     ```

3. **Task 3: Data Validation**:

   - Validate the transformed data.

```
tasks:
  - task_key: validate_data
    depends_on:
      - task_key: transform_data
    notebook_task:
      notebook_path: /path/to/validate_notebook
```

4. **Task 4: Data Loading**:

   - Load the validated data into the data warehouse.

```
tasks:
  - task_key: load_data
    depends_on:
      - task_key: validate_data
    notebook_task:
      notebook_path: /path/to/load_notebook
```

## Benefits

- **Data Integrity**: Ensures that only validated data is loaded into the data warehouse, maintaining data quality.

- **Error Handling**: If any task fails (e.g., validation fails), the subsequent tasks do not run, preventing invalid data from being loaded.

- **Modularity**: Each task can be developed, tested, and maintained independently, improving the overall manageability of the pipeline.

▼ **85.Review a task's execution history.**

To review a task's execution history in Databricks, you can use the Databricks UI to access detailed information about past runs. Here's how you can do it:

## Steps to Review Task Execution History

1. **Navigate to Workflows**:

   - In the Databricks workspace, click on **Workflows** in the sidebar.

2. **Select the Job**:

   - Click on the name of the job you want to review. This will open the job details page.

3. **View Job Runs**:

   - On the job details page, go to the **Runs** tab. This tab shows a list of all the runs for the selected job, including both active and completed runs.

4. **Select a Task**:

- Click on a specific run to view its details. This will open the job run details page.

- On the job run details page, you can see a list of tasks that were executed as part of the job run.

5. **View Task Run History**:

- Click on the task you are interested in. This will open the task run details page.

- On the task run details page, you can view the history of all runs for that task, including successful and unsuccessful runs.

## Example

Here's a quick example of what you might see:

- **Job Name**: Data Processing Job

- **Run ID**: 12345

- **Task Name**: Data Transformation

- **Task Run History**:

  - **Run ID**: 12345-1

    - **Status**: Success

    - **Start Time**: 2024-12-16 10:00:00

    - **End Time**: 2024-12-16 10:15:00

  - **Run ID**: 12345-2

    - **Status**: Failed

    - **Start Time**: 2024-12-17 10:00:00

    - **End Time**: 2024-12-17 10:05:00

## Benefits

- **Detailed Insights**: Provides detailed information about each task run, including start and end times, status, and any error messages.

- **Troubleshooting**: Helps in identifying and troubleshooting issues by examining the history of task executions.

- **Performance Monitoring**: Allows you to monitor the performance and reliability of your tasks over time.

▼ **86.Identify CRON as a scheduling opportunity.**

CRON is a powerful scheduling tool that you can use in Databricks to automate the execution of jobs at specified times or intervals. Here's how CRON can be utilized as a scheduling opportunity:

## CRON Scheduling in Databricks

**CRON Syntax**:

- CRON expressions allow you to define complex schedules using a simple string format. This format specifies the exact times and intervals at which a job should run.

**Example CRON Expression**:

- `0 0 * * *`: This expression schedules a job to run at midnight every day.

## Setting Up CRON Scheduling in Databricks

1. **Navigate to Workflows**:

   - In the Databricks workspace, click on **Workflows** in the sidebar.

2. **Create or Edit a Job**:

   - Click on the job name to open the job configuration or create a new job.

3. **Add a Schedule**:

   - In the job details panel, click **Add trigger**.

   - Select **Scheduled** as the trigger type.

   - Choose **Advanced** to use CRON syntax for more control over the schedule.

4. **Define the CRON Expression**:

   - Enter your CRON expression in the provided field. Optionally, you can select the **Show Cron Syntax** checkbox to edit the schedule in Quartz Cron Syntax.

## Benefits of Using CRON Scheduling

- **Flexibility**: CRON expressions provide fine-grained control over job scheduling, allowing you to specify exact times, days, and intervals.

- **Automation**: Automates repetitive tasks, reducing the need for manual intervention and ensuring timely execution of jobs.

- **Efficiency**: Helps in optimizing resource usage by scheduling jobs during off-peak hours or at specific intervals.

## Example Use Case

Imagine you have a job that processes daily sales data. You can use a CRON expression to schedule this job to run at 2 AM every day, ensuring that the data is processed and ready for analysis by the start of the business day.

```
schedule:
  cron: "0 2 * * *"
```

## ▼ 87.Debug a failed task.

To debug a failed task in Databricks, follow these steps to identify the cause of the failure and resolve the issue:

## Steps to Debug a Failed Task

1. **Identify the Failed Task**:

   - Navigate to **Workflows** in the Databricks sidebar.

   - Click on **Job Runs** to see a list of all job runs.

   - Find the job run that contains the failed task. The **Runs** tab shows a history of runs, including successful and failed ones.

   - Click on the specific job run to view details.

2. **View Task Details**:

   - In the job run details page, locate the failed task in the task list.

   - Click on the failed task to open the **Task run details** page. This page provides detailed information about the task, including error messages, logs, and metadata.

3. **Examine Error Messages and Logs**:

   - Review the error message displayed on the **Task run details** page to understand the nature of the failure.

   - Check the **Driver logs** and **Executor logs** for more detailed information. These logs can provide insights into what went wrong during the task execution.

4. **Use the Apache Spark UI**:

   - If the task involves a Spark job, use the Spark UI to debug the application. The Spark UI provides detailed information about the job's execution, including stages, tasks, and performance metrics.

   - Access the Spark UI by clicking on the **View Spark UI** link in the task details.

5. **Fix the Cause of Failure**:

   - Based on the error messages and logs, identify the root cause of the failure. Common issues include data quality problems, misconfigurations, or insufficient compute resources.

   - Update the task configuration if needed. For example, you might need to adjust the cluster settings, increase resource quotas, or correct data paths.

6. **Re-run the Task**:

   - After fixing the issue, re-run the failed task. You can do this from the job run details page by clicking on **Re-run**.

   - Monitor the task to ensure it completes successfully.

## Example Scenario

Imagine a task failed due to insufficient memory. The error message in the logs might indicate an `OutOfMemoryError`. To resolve this, you could increase the memory allocation for the cluster:

1. **Edit Task Configuration**:

   - Click on **Edit task** in the job details panel.

   - Under **Compute**, click **Configure** to adjust the cluster settings.

   - Increase the memory allocation and save the changes.

2. **Re-run the Task**:

   - Re-run the task to see if the issue is resolved.

By following these steps, you can effectively debug and resolve issues with failed tasks in Databricks.

▼ **88.Set up a retry policy in case of failure.**

Setting up a retry policy in Databricks ensures that tasks are automatically retried in case of failure, improving the robustness of your workflows. Here's how you can configure a retry policy for a job:

## Steps to Set Up a Retry Policy

1. **Navigate to Workflows**:

   - In the Databricks workspace, click on **Workflows** in the sidebar.

2. **Create or Edit a Job**:

   - Click on the job name to open the job configuration or create a new job.

3. **Configure Task Settings**:

   - In the job details panel, click on the task you want to configure.

   - Scroll down to the **Retry policy** section.

4. **Set Retry Policy**:

   - **Max Retries**: Specify the maximum number of retry attempts.

   - **Min Retry Interval**: Define the minimum interval between retries.

   - **Retry on Timeout**: Enable this option if you want the task to be retried in case of a timeout.

## Example Configuration

Here's an example of how to set up a retry policy for a task:

```
tasks:
  - task_key: my_task
```

```
    notebook_task:
      notebook_path: /path/to/notebook
    max_retries: 3
    min_retry_interval_millis: 60000  # 1 minute
    retry_on_timeout: true
```

## Benefits of a Retry Policy

- **Increased Resilience**: Automatically retries tasks that fail due to transient issues, reducing the need for manual intervention.

- **Improved Reliability**: Ensures that temporary failures (e.g., network issues, resource contention) do not cause the entire job to fail.

- **Cost Efficiency**: By specifying a retry interval, you can avoid immediate retries that might fail again due to the same transient issue.

▼ **89.Create an alert in the case of a failed task.**

To create an alert for a failed task in Databricks, you can set up email or system notifications. Here's how you can do it:

## Steps to Create an Alert for a Failed Task

1. **Navigate to Workflows**:

   - In the Databricks workspace, click on **Workflows** in the sidebar.

2. **Create or Edit a Job**:

   - Click on the job name to open the job configuration or create a new job.

3. **Add Notifications**:

   - In the job details panel, click **Edit notifications**.

4. **Add Email Notification**:

   - Click **Add Notification** and select **Email address** in the **Destination** field.

   - Enter the email address where you want to receive notifications.

   - Check the box for **Failure** to receive alerts when a task fails.

   - Click **Confirm** to save the notification settings.

     1

## Example Configuration

Here's an example of setting up an email notification for a failed task:

```
tasks:
  - task_key: my_task
```

```
notebook_task:
  notebook_path: /path/to/notebook
notifications:
  - email: "your-email@example.com"
    on_failure: true
```

## System Notifications

You can also set up system notifications to integrate with tools like Slack, Microsoft Teams, PagerDuty, or any webhook-based service:

1. **Add System Notification**:

   - Click **Add Notification** and select **System destination** in the **Destination** field.

   - Choose the system destination (e.g., Slack, Microsoft Teams).

   - Check the box for **Failure** to receive alerts when a task fails.

   - Click **Confirm** to save the notification settings.

## Benefits

- **Immediate Alerts**: Receive instant notifications when a task fails, allowing you to take prompt action.

- **Integration**: Seamlessly integrate with popular communication and incident management tools.

- **Customization**: Configure notifications for different events (e.g., job start, success, failure) and multiple destinations.

▼ **90.Identify that an alert can be sent via email.**

you can send alerts via email using Databricks. Databricks allows you to configure email notifications for various events, such as job starts, completions, failures, or when a job exceeds a certain duration. You can set up these notifications to be sent to one or more email addresses

To configure email notifications in Databricks:

1. **Navigate to the Job Details Panel**: Go to the job you want to monitor.

2. **Edit Notifications**: Click on "Edit notifications" in the job details panel.

3. **Add Email Notification**: Click on "+ Add" next to Notifications and select "Email" as the destination.

4. **Configure Notification Types**: Choose the types of events (e.g., job start, completion, failure) for which you want to receive email alerts.

5. **Save Settings**: Confirm and save your notification settings.

Additionally, you can use custom scripts within Databricks notebooks to send emails programmatically using services like AWS SNS or SMTP servers

# Section 5: Data Governance(9%)

▼ **91.Identify one of the four areas of data governance.**

One of the four key areas of data governance in a Databricks environment is **unifying data and AI management**. This involves establishing a comprehensive governance process that ensures the availability, usability, integrity, and security of data and AI assets. By unifying management, organizations can streamline data operations, improve data quality, and enhance decision-making capabilities

▼ **92.Compare and contrast meta stores and catalogs.**

Meta stores and catalogs serve different purposes in managing and organizing data, particularly in environments like Databricks.

## Meta Stores

- **Definition**: A meta store is a centralized repository that stores metadata about data assets. It includes information such as schema definitions, table locations, and data types.

- **Functionality**: Meta stores manage metadata for various data sources, ensuring consistency and accessibility. They are crucial for data discovery, data lineage, and governance.

- **Example**: In Databricks, the Hive Metastore is commonly used to manage metadata for tables and databases.

## Catalogs

- **Definition**: Catalogs are collections of data assets, such as tables and views, organized within a specific namespace. They provide a structured way to manage and access data.

- **Functionality**: Catalogs help in organizing data assets into logical groupings, making it easier to manage permissions, access controls, and data governance policies.

- **Example**: Databricks Unity Catalog is an example that provides a unified governance solution for data and AI assets, allowing for fine-grained access control and auditing.

## Comparison

- **Scope**: Meta stores focus on metadata management, while catalogs organize actual data assets.

- **Usage**: Meta stores are used for maintaining metadata consistency and supporting data governance, whereas catalogs are used for structuring and managing data access.
- **Integration**: Catalogs often rely on meta stores to retrieve metadata about the data assets they manage.

▼ **93.Identify Unity Catalog securable.**

Unity Catalog in Databricks secures data through a hierarchical model of securable objects, each with specific privileges that can be granted to users, groups, or service principals. Here are the main securable objects in Unity Catalog:

1. **Metastore**: The top-level container for metadata. Privileges on the metastore can be granted to manage catalogs within it.

2. **Catalog**: Organizes data assets and can contain schemas, tables, and views. Privileges on a catalog can be inherited by all objects within it.

3. **Schema**: Also known as databases, schemas contain tables and views. Privileges on a schema can be inherited by all tables and views within it.

4. **Table**: The lowest level in the hierarchy, tables can be managed or external. Privileges on tables control access to the data they contain.

5. **View**: Read-only objects created from queries on tables. Privileges on views control access to the data they present.

6. **Volume**: Storage containers for data, either managed or external. Privileges on volumes control access to the data stored within them.

7. **Function**: User-defined functions or MLflow registered models contained within a schema.

Privileges in Unity Catalog are hierarchical, meaning that granting a privilege at a higher level (like a catalog) automatically grants that privilege to all lower-level objects (like schemas and tables) within that catalog

▼ **94.Define a service principal**

A **service principal** in Databricks is an identity created for use with automated tools, jobs, and applications. It provides API-only access to Databricks resources, enhancing security by avoiding the use of user or group credentials for automation purposes.

Service principals are particularly useful for:

- **Automated Jobs**: Running scheduled tasks without relying on user credentials.
- **CI/CD Pipelines**: Integrating Databricks with continuous integration and deployment workflows.
- **Scripts and Applications**: Allowing scripts and external applications to interact with Databricks securely.

Service principals can be granted specific roles and permissions, similar to users, to control their access to resources within Databricks

▼ **95.Identify the cluster security modes compatible with Unity Catalog.**

Unity Catalog in Databricks supports the following cluster security modes:

1. **Single User Access Mode**: This mode is designed for clusters used by a single user. It provides the highest level of security by ensuring that only the cluster owner can access the data and resources.

2. **Shared Access Mode**: This mode allows multiple users to share the same cluster. It is recommended for most workloads as it supports fine-grained access control and user isolation, ensuring that each user's code runs in a secure, isolated environment.

3. **No Isolation Shared Mode**: This is a legacy mode that does not support Unity Catalog. It is primarily used for backward compatibility with the Hive Metastore

▼ **96.Create a UC-enabled all-purpose cluster.**

To create a Unity Catalog (UC)-enabled all-purpose cluster in Databricks, follow these steps:

1. **Navigate to Compute**:

   - In the Databricks workspace sidebar, click on **Compute**.

2. **Create Cluster**:

   - Click the **Create Cluster** button.

3. **Cluster Configuration**:

   - **Cluster Name**: Enter a name for your cluster.

   - **Cluster Mode**: Select **Standard**.

   - **Databricks Runtime Version**: Choose a runtime version that supports Unity Catalog (e.g., Databricks Runtime 10.0 or above).

   - **Access Mode**: Select **Single User** or **Shared** to ensure compatibility with Unity Catalog.

   - **Enable Unity Catalog**: Ensure that the cluster is configured to use Unity Catalog by selecting the appropriate options in the cluster configuration.

4. **Advanced Options**:

   - **Libraries**: If needed, you can add libraries or JAR files.

   - **Tags**: Add any tags for organizational purposes.

   - **Autoscaling**: Enable autoscaling if required.

5. **Create**:

   - Click the **Create Cluster** button to spin up your UC-enabled all-purpose cluster.

Your cluster will start and be ready to use shortly. You can now leverage Unity Catalog for fine-grained access control and data governance.

▼ **97.Create a DBSQL warehouse.**

To create a Databricks SQL (DBSQL) warehouse, follow these steps:

1. **Navigate to SQL Warehouses**:

   - In the Databricks workspace sidebar, click on **SQL Warehouses**.

2. **Create SQL Warehouse**:

   - Click the **Create SQL Warehouse** button.

3. **Configure Warehouse**:

   - **Name**: Enter a name for your SQL warehouse.

   - **Cluster Size**: Choose the size of the cluster based on your workload requirements.

   - **Auto Stop**: Set the idle time after which the warehouse should automatically stop to save costs.

   - **Scaling**: Configure the minimum and maximum number of clusters to handle concurrent queries.

4. **Advanced Options** (Optional):

   - You can configure additional settings such as tags, libraries, and network configurations if needed.

5. **Create**:

   - Click the **Create** button to finalize the setup.

Your SQL warehouse will start automatically and be ready for use. You can now run SQL queries, create dashboards, and connect to BI tools using this warehouse

▼ **98.Identify how to query a three-layer namespace.**

To query a three-layer namespace in Databricks, particularly when using Unity Catalog, you need to follow the format: `<catalog-name>.<schema-name>.<table-name>`. Here's how you can do it:

1. **Select the Catalog**:

   - Ensure you have the necessary permissions to use the catalog. For example, `USE CATALOG my_catalog;`.

2. **Select the Schema**:

   - Similarly, ensure you have permissions to use the schema. For example, `USE SCHEMA my_catalog.my_schema;`.

3. **Query the Table**:

   - Use the fully qualified name to query the table. For example:

   ```
   SELECT * FROM my_catalog.my_schema.my_table;
   ```

## Example

If you have a catalog named `sales`, a schema named `transactions`, and a table named `orders`, your query would look like this:

```
SELECT * FROM sales.transactions.orders;
```

## Permissions

Ensure you have the following permissions:

- **USE_CATALOG** on the catalog.
- **USE_SCHEMA** on the schema.
- **SELECT** on the table.

These permissions are necessary to access and query the data within the three-layer namespace

▼ **99.Implement data object access control**

To implement data object access control in Databricks, you can use Access Control Lists (ACLs) to manage permissions for various workspace objects. Here's a step-by-step guide:

## Step-by-Step Guide to Implement Data Object Access Control

1. **Enable Access Control**:
   - Ensure your Databricks workspace is on the Premium plan or above, as access control features require this level.
   - Navigate to the **Admin Console** and enable access control for your workspace.

2. **Define Roles and Permissions**:
   - Identify the roles and permissions needed for your organization. Common roles include Admin, Data Engineer, Data Scientist, and Business Analyst.
   - Use predefined roles or create custom roles based on your requirements.

3. **Assign Permissions**:
   - **Workspace Objects**: Assign permissions to users or groups for notebooks, clusters, jobs, and other workspace objects.
     - Example: To grant a user read access to a notebook, navigate to the notebook, click on **Permissions**, and add the user with the **CAN READ** permission.
   - **Data Objects**: Use Unity Catalog to manage permissions for catalogs, schemas, tables, and views.
     - Example: To grant a user select access to a table, use the following SQL command:

```
GRANT SELECT ON TABLE my_catalog.my_schema.my_table TO
 `user@example.com`;
```

4. **Fine-Grained Access Control**:

   - Implement column-level and row-level security using dynamic views.

   - Example: Create a view that filters rows based on the current user:

     ```
     CREATE VIEW my_catalog.my_schema.filtered_table AS
     SELECT * FROM my_catalog.my_schema.my_table
     WHERE user_id = current_user();
     ```

5. **Audit and Monitor**:

   - Regularly audit permissions and access logs to ensure compliance and security.

   - Use Databricks audit logs to track changes and access patterns.

## Example: Granting Permissions

Here's an example of how to grant different levels of access to a table in Unity Catalog:

```
-- Grant SELECT permission to a user
GRANT SELECT ON TABLE my_catalog.my_schema.my_table TO `user@example.co
m`;


-- Grant INSERT permission to a group
GRANT INSERT ON TABLE my_catalog.my_schema.my_table TO `data_engineers
`;


-- Revoke DELETE permission from a user
REVOKE DELETE ON TABLE my_catalog.my_schema.my_table FROM `user@exampl
e.com`;
```

By following these steps, you can effectively manage data object access control in Databricks, ensuring that users have the appropriate permissions to access and manipulate data securely

▼ **100.Identify colocating meta stores with a workspace as best practice.**

Colocating meta stores with a workspace is considered a best practice in Databricks for several reasons:

1. **Performance Optimization**: By colocating the metastore with the workspace, you reduce latency and improve query performance. This is because the metadata retrieval and data access operations are faster when they are in the same region.

2. **Simplified Management**: Managing a single metastore for a workspace simplifies administrative tasks. It ensures that all metadata is centralized, making it easier to maintain consistency and apply governance policies.

3. **Cost Efficiency**: Colocating helps in reducing data transfer costs between regions. When the metastore and workspace are in the same region, you avoid cross-region data transfer charges.

4. **Enhanced Security**: Keeping the metastore and workspace in the same region can enhance security by reducing the attack surface and ensuring that data governance policies are consistently applied

▼ **101. Identify using service principals for connections as best practice.**

Using service principals for connections is considered a best practice for several reasons:

1. **Enhanced Security**: Service principals provide a secure way to authenticate applications and automated tools without using user credentials. This reduces the risk of credential exposure and unauthorized access.

2. **Granular Access Control**: You can assign specific roles and permissions to service principals, ensuring that they have only the necessary access to perform their tasks. This follows the principle of least privilege.

3. **Automated Credential Management**: Service principals support automated credential management, including rotation and expiration, which helps maintain security without manual intervention.

4. **Audit and Monitoring**: Service principals allow for better tracking and auditing of actions performed by automated processes. This helps in identifying and responding to potential security incidents.

5. **Scalability**: Using service principals simplifies the management of permissions and access for large-scale deployments and integrations, making it easier to scale your infrastructure securely

▼ **102. Identify the segregation of business units across catalog as best practice.**

Segregating business units across catalogs is considered a best practice in Databricks for several reasons:

1. **Data Isolation:** By assigning each business unit its own catalog, you ensure that data is isolated and access is controlled at a granular level. This helps prevent unauthorized access and maintains data privacy.

2. **Simplified Governance**: Managing permissions and governance policies becomes easier when each business unit has its own catalog. You can apply specific policies and controls tailored to the needs of each unit.

3. **Improved Performance**: Segregating data into different catalogs can enhance performance by reducing the complexity of queries and metadata management. This is particularly

beneficial for large organizations with diverse data needs.

4. **Organizational Alignment**: Catalogs can mirror the organizational structure, making it intuitive for users to find and manage data relevant to their business unit. This alignment helps in maintaining clarity and efficiency in data operations.

## Example

If you have multiple business units such as Sales, Marketing, and Finance, you can create separate catalogs for each:

- `sales_catalog`

- `marketing_catalog`

- `finance_catalog`

Each catalog can then contain schemas and tables specific to the respective business unit, ensuring clear separation and management.