

What is difference between OLAP and OLTP?

	OLAP	OLTP
	Online Analytical Processing	Online Transaction Processing
Purpose	OLAP databases are designed for complex analytical and ad-hoc querying. They are used for data analysis, business intelligence, and decision support. OLAP systems are optimized for retrieving and processing large volumes of historical or aggregated data to help users gain insights and make informed decisions.	OLTP databases are designed for transactional processing. They are used for recording and managing day-to-day operational data, such as sales transactions, inventory management, and customer interactions. OLTP systems are optimized for quick and precise retrieval and updating of individual records.
Data Structure	OLAP databases typically store historical, summarized, and multidimensional data. They use a star or snowflake schema and often involve data warehousing techniques. Data is denormalized and aggregated for efficient querying.	OLTP databases store normalized transactional data with a focus on maintaining data integrity. The data is typically organized into tables with relationships to minimize redundancy and ensure accuracy.
Query Complexity	OLAP queries are often complex and involve aggregations, grouping, and calculations. Users perform data analysis to answer strategic questions or gain insights into trends.	OLTP queries are straightforward and primarily involve basic CRUD operations (Create, Read, Update, Delete). The focus is on recording and retrieving individual transactions efficiently.
Data Volume	OLAP databases handle large volumes of historical data and are optimized for read-heavy workloads.	OLTP databases handle a high volume of concurrent, short, and simple transactions. They are optimized for write-heavy workloads.
Indexing	OLAP databases often have a limited number of indexes optimized for query performance.	OLTP databases typically have many indexes to ensure the speedy retrieval of individual records.
Concurrency	OLAP systems usually have low concurrency requirements because analytical queries are less frequent.	OLTP systems require high levels of concurrency to support multiple simultaneous users performing various transactions.
Data Modification	Data modification operations (inserts, updates, deletes) are infrequent in OLAP systems and usually occur during the data loading process.	Data modification operations are frequent and form the core of OLTP database functionality.

## What is Data Warehouse?

A data warehouse is a specialized and centralized repository for storing, managing, and consolidating large volumes of data from various sources within an organization. It is designed to facilitate efficient querying and analysis of data to support business intelligence and decision-making processes. Data warehouses are a core component of data management and analytics in many organizations. Data warehouses play a critical role in enabling organizations to gain insights from their data, support strategic decision-making, and improve overall business performance. Here are the key characteristics and functions of a data warehouse:

1. **Data Integration:** Data warehouses gather data from multiple sources, which can include operational databases, external data feeds, spreadsheets, and more. The data is integrated and transformed to ensure consistency and quality.
2. **Data Storage:** Data is stored in a structured format within the data warehouse. This format is often a relational database management system (RDBMS) with tables, columns, and rows. The data is organized for ease of retrieval and analysis.
3. **Data Modeling:** Data in a data warehouse is typically organized using a specific schema that is optimized for query performance. Common schemas include the star schema and snowflake schema. These schemas define fact tables (containing numerical measures) and dimension tables (describing attributes).
4. **Data History:** Data warehouses often maintain historical data, allowing organizations to analyze trends and track changes over time. This historical perspective is essential for business intelligence and reporting.
5. **Data Quality:** Data quality is a priority in data warehousing. Systems often include mechanisms for data cleansing, validation, and consistency checks to ensure the accuracy and reliability of the stored data.
6. **Data Security:** Data warehouses incorporate security features to protect sensitive and confidential data. These features include access control, authentication, encryption, and auditing.
7. **Query and Analysis:** Data warehouses are designed to support complex querying and reporting. Tools and technologies for data analysis, such as business intelligence (BI) software, can be used to derive insights from the data.
8. **Scalability:** Data warehouses need to be scalable to handle the growing volume of data and increasing user demands. This scalability may involve hardware and software scaling, as well as data partitioning and distribution.
9. **Data Access:** Users can access the data warehouse through various means, including SQL queries, reporting tools, dashboards, and analytics applications.
10. **Data Loading:** Data is periodically loaded into the data warehouse, often in batch processes scheduled at regular intervals.

## What is DAG and it's benefit?

A DAG, or Directed Acyclic Graph, is a data structure commonly used in various computational and mathematical applications, including computer science, task scheduling, and data processing. It is a collection of nodes connected by directed edges or arrows in a way that does not form any cycles. Here's an explanation of what a DAG is and its benefits:

Directed Acyclic Graph (DAG):

1. **Directed:** In a DAG, each edge has a direction, which means it goes from one node (or vertex) to another. This direction indicates a relationship or dependency between the two nodes. Node A points to node B if there's a directed edge from A to B.
2. **Acyclic:** The term "acyclic" means that there are no cycles or closed loops in the graph. In other words, you cannot start at one node and follow directed edges to return to the same node. This property is essential in various applications, especially when modeling dependencies, to avoid infinite loops.

Benefits of DAGs:

1. **Modeling Dependencies:** DAGs are widely used to model and represent dependencies between tasks, operations, or events. For example, in project scheduling, tasks can be represented as nodes, and dependencies between tasks as directed edges in a DAG. This helps in visualizing and understanding the order in which tasks should be performed.
2. **Task Scheduling:** DAGs are crucial for task scheduling in various domains, such as job scheduling in data processing frameworks like Apache Spark or Apache Airflow. Each task is represented as a node, and dependencies between tasks are defined using directed edges. This ensures that tasks are executed in the correct order to meet dependencies.
3. **Optimizing Computations:** In computational tasks like dynamic programming and optimization problems, DAGs are used to store and organize intermediate results. By storing and reusing intermediate computations, redundant work can be avoided, leading to more efficient algorithms.
4. **Parallel Processing:** In distributed computing and parallel processing, DAGs are used to represent tasks or data transformations that can be executed concurrently. This allows for the efficient utilization of multiple processors or computing resources.
5. **Data Flow:** In data processing and ETL (Extract, Transform, Load) pipelines, DAGs represent the flow of data from source to destination. Each node in the DAG represents a transformation or processing step, and edges indicate the data flow direction.
6. **Dependency Resolution:** In software dependency management, DAGs help resolve and manage dependencies between software components or libraries. This is common in package managers like npm (Node Package Manager) and pip (Python package manager).
7. **Reducing Re-computation:** By storing intermediate results in a DAG, redundant computations can be minimized. This is particularly beneficial in scenarios where calculations are expensive, such as in scientific simulations or complex data analysis.

In summary, DAGs provide a structured way to represent and manage dependencies, making them a valuable tool in various computational and scheduling tasks. They help improve efficiency, minimize redundancy, and ensure that tasks or operations are executed in the correct order, which is crucial in fields like project management, data processing, and distributed computing.

### What is difference Temp view vs Global View in Data bricks?

	Temporary View (Temp View):	Global View (Global Temp View)
Scope	Temp views are limited to a single session or notebook in Databricks. They are local to the specific session or notebook where they are created.	Global views have a broader scope and are available across multiple sessions or notebooks in Databricks. They are global to the entire workspace.
Lifetime	Temp views are temporary and exist for the duration of the session or notebook in which they were created. They are automatically removed when the session ends, or the notebook is closed.	Global views are temporary like temp views but persist as long as the Databricks workspace is active. They are removed only when explicitly dropped or if the workspace is deleted.
Usage	Temp views are typically used for short-term or ad-hoc data analysis within a specific notebook. They are useful when you need to query or analyse data within a specific context.	Global views are useful for sharing data across multiple notebooks or collaborating with others within the same Databricks workspace. They allow you to define a dataset once and access it from different notebooks or sessions.

In summary, the key difference between Temp Views and Global Views in Databricks is their scope and lifetime. Temp Views are local to a specific session or notebook and have a short-term existence, while Global Views are accessible across the entire Databricks workspace and persist until explicitly dropped or the workspace is deleted. The choice between the two depends on your specific use case and whether you need the data to be available globally or just within a specific session or notebook.

### What is difference between coalesce and re-partition?

**coalesce** and **repartition** are two operations in Apache Spark used for controlling the partitioning of data in a Resilient Distributed Dataset (RDD) or a DataFrame. While both operations are related to partitioning data, they serve different purposes and have some key differences:

Coalesce	Repartition
<ul style="list-style-type: none"><li>• <b>coalesce</b> is an operation that reduces the number of partitions in an RDD or DataFrame.</li><li>• It is used to merge smaller partitions into larger ones, which can be beneficial when you have too many small partitions, and you want to reduce the overhead of managing them.</li><li>• <b>coalesce</b> typically results in data shuffling, but it tries to minimize data movement by coalescing adjacent partitions within the same executor node.</li><li>• It's often used to optimize the number of partitions to match the available resources (CPU cores, memory) and minimize task scheduling overhead.</li><li>• You can optionally specify the number of partitions you want after coalescing, and Spark will try to coalesce the data into that number of partitions.</li><li>• Syntax to Reduce the number of partitions to 4 <b>df = df.coalesce(4)</b></li></ul>	<ul style="list-style-type: none"><li>• <b>repartition</b> is used to change the number of partitions in an RDD or DataFrame, either by increasing or decreasing the number of partitions.</li><li>• It can be used to redistribute data across partitions to achieve better parallelism or to balance the data distribution.</li><li>• <b>repartition</b> often involves data shuffling, and it allows you to explicitly specify the number of partitions you want in the resulting RDD or DataFrame.</li><li>• You can also use <b>repartition</b> to change the partitioning column or expression, which is particularly useful when you want to change the way data is distributed across partitions based on a specific column.</li><li>• Syntax to Repartition to 8 partitions <b>df = df.repartition(8)</b></li></ul>

## What is difference between coalesce and partition?

The terms "coalesce" and "partition" are often used in the context of data processing, particularly in distributed computing frameworks like Apache Spark, but they have different meanings and purposes. Let me clarify the distinction between the two:

Coalesce	Partition
<p>Coalesce is primarily used to reduce the number of partitions or data shuffling in a distributed dataset or table. It involves combining existing partitions into a smaller number of partitions without performing a full data shuffle. This operation is typically used to optimize the distribution of data and reduce the overhead of managing a large number of partitions.</p>	<p>Partitioning refers to the way data is physically organized or distributed within a dataset or table. It involves dividing data into smaller subsets based on specific criteria, such as a column value or a hashing algorithm. Partitioning is often used to improve query performance, as it allows for efficient data retrieval based on the partitioning key.</p>
<p>Key points about coalesce:</p> <ul style="list-style-type: none"><li>• Decreases the number of partitions.</li><li>• Minimizes data shuffling.</li><li>• Typically used to improve query performance by reducing the overhead of managing many partitions.</li><li>• Does not change the data's content; it only changes the organization of partitions.</li></ul> <ul style="list-style-type: none"><li>• In the context of Apache Spark, "coalesce" refers to an operation that reduces the number of partitions in a Resilient Distributed Dataset (RDD) or a DataFrame. The purpose of coalesce is to merge smaller partitions into larger ones, thereby reducing the overhead of managing a large number of partitions.</li><li>• Coalesce tries to minimize data shuffling by coalescing adjacent partitions within the same executor node.</li><li>• It is typically used to optimize the number of partitions to match the available resources (CPU cores, memory) and minimize task scheduling overhead.</li><li>• Syntax to Reduce the number of partitions to 3 <b>new_rdd = old_rdd.coalesce(3)</b></li></ul>	<p>Key points about partitioning:</p> <ul style="list-style-type: none"><li>• Divides data into smaller subsets (partitions) based on specific criteria (e.g., a column's value).</li><li>• Improves query performance by allowing for partition pruning, where only relevant partitions are scanned during queries.</li><li>• Commonly used in databases and distributed data processing frameworks.</li><li>• May involve data redistribution and can be an expensive operation if not done carefully.</li></ul> <ul style="list-style-type: none"><li>• In the context of data storage and databases, "partition" typically refers to the practice of dividing and organizing data in a data store (e.g., a table in a relational database) into smaller, manageable units based on certain criteria, such as date ranges or specific values of a column.</li><li>• Data partitioning helps improve data management, retrieval, and query performance. Queries can be executed on specific partitions, which can reduce the amount of data that needs to be scanned.</li><li>• In the context of Spark, when we talk about "partitioning," it often relates to how data is distributed across partitions within RDDs or DataFrames.</li><li>• Syntax to create Partition at table level: <b>CREATE TABLE sales (</b>     <b>sale_date DATE,</b>     <b>amount DECIMAL</b> <b>)</b>     <b>USING delta</b>     <b>PARTITION BY RANGE (sale_date);</b>      <b>OR</b> df.write.format("delta")/     .partitionBy("state")/     .save("/FileStore/tables/sales")</li></ul>

## How to create spark session?

Creating a Spark session is a fundamental step when working with Apache Spark, a popular open-source big data processing framework. A Spark session provides an entry point to interact with Spark and allows you to configure various settings and work with data in a distributed and parallel manner. Here's how you can create a Spark session using the Python API (PySpark):

```
from pyspark.sql import SparkSession
```

```
# Create a Spark session
```

```
spark = SparkSession.builder \
```

```
.appName("YourAppName") \ # Set your application name
```

```
.master("local[*]") \ # Set the Spark master URL (local[*] for local mode)
```

```
.config("key", "value") \ # Additional configurations (if needed)
```

```
.getOrCreate()
```

# Once you have created the Spark session, you can use it to work with Spark DataFrames, perform transformations, and run Spark jobs.

Let's break down the steps:

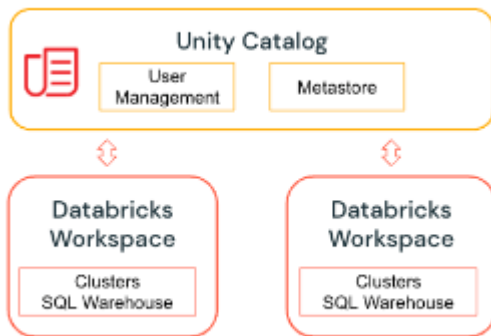
1. Import the **SparkSession** class from **pyspark.sql**.
2. Create a **SparkSession** instance using the **builder** pattern.
3. Set the application name with **.appName("YourAppName")**, where "YourAppName" is the name you want to give to your Spark application.
4. Set the master URL using **.master("local[\*]")**. The master URL specifies the cluster manager or execution mode. In this example, we're using "local[\*]" for local mode. In a production environment, you would specify the cluster manager's URL (e.g., "yarn", "mesos", "spark://host:port").
5. Optionally, you can set additional configuration options using **.config("key", "value")**. For example, you can configure the number of CPU cores, memory, or other Spark settings here.
6. Finally, call **.getOrCreate()** to either get an existing Spark session or create a new one if it doesn't exist.

Once you've created the Spark session, you can use it to perform various data processing tasks, including reading and writing data, running SQL queries, and applying transformations to Spark DataFrames.

Remember that this example is for PySpark (Python API for Spark). If you're using Spark with a different language API (such as Scala or Java), the process for creating a Spark session will be similar, but the code syntax may vary.

## What is Unity Catalog in Azure Databricks and it's Key Features?

Unity Catalog provides centralized access control, auditing, lineage, and data discovery capabilities across Azure Databricks workspaces.



Key features of Unity Catalog include:

- **Define once, secure everywhere:** Unity Catalog offers a single place to administer data access policies that apply across all workspaces.
- **Standards-compliant security model:** Unity Catalog's security model is based on standard ANSI SQL and allows administrators to grant permissions in their existing data lake using familiar syntax, at the level of catalogs, databases (also called schemas), tables, and views.
- **Built-in auditing and lineage:** Unity Catalog automatically captures user-level audit logs that record access to your data. Unity Catalog also captures lineage data that tracks how data assets are created and used across all languages.
- **Data discovery:** Unity Catalog lets you tag and document data assets, and provides a search interface to help data consumers find data.
- **System tables (Public Preview):** Unity Catalog lets you easily access and query your account's operational data, including audit logs, billable usage, and lineage.



## What is lazy transformation in databricks?

Databricks and Apache Spark, "lazy transformation" refers to a fundamental concept related to how data transformations are processed in a lazy and optimized manner. Lazy transformation is a key part of the Spark's execution model, and it contributes to the platform's efficiency and performance. Here's an explanation of lazy transformation in Databricks and Apache Spark:

### Lazy Transformation in Apache Spark:

1. **Transformations:** In Spark, transformations are operations that modify a Resilient Distributed Dataset (RDD) or a DataFrame to produce a new RDD or DataFrame. Examples of transformations include **map**, **filter**, **groupBy**, and **join**. Transformations are performed on distributed data, and they represent a sequence of operations to be executed on that data.
2. **Lazy Evaluation:** Spark employs a lazy evaluation strategy for transformations. This means that when you apply a transformation to an RDD or DataFrame, Spark doesn't immediately execute the transformation. Instead, it records the transformation in a logical execution plan (a directed acyclic graph or DAG).
3. **Benefits of Lazy Evaluation:**
  - **Optimization:** Lazy evaluation allows Spark to perform various optimizations on the execution plan. It can reorder, combine, or eliminate certain operations to improve performance.
  - **Efficiency:** Spark avoids unnecessary work by only executing the transformations that are required to compute the final result. This saves computation and memory resources.
  - **Fault Tolerance:** Spark can recover from failures more efficiently because it has the original data and can recompute only the lost or missing transformations.
4. **Action Triggers Execution:** Lazy evaluation means that transformations are not executed until an "action" is called. Actions are operations that return results to the driver program or write data to an external storage system. Common actions include **count**, **collect**, **saveAsTextFile**, and **show**.

#### Example:

```
# Create an RDD and apply transformations (lazy)

data = [1, 2, 3, 4, 5]

rdd = sc.parallelize(data)

filtered_rdd = rdd.filter(lambda x: x % 2 == 0)

mapped_rdd = filtered_rdd.map(lambda x: x * 2)

# Perform an action to trigger execution

result = mapped_rdd.collect()
```

In this example, the **filter** and **map** transformations are applied to the RDD, but the actual execution doesn't happen until the **collect** action is called. The lazy evaluation mechanism allows Spark to optimize and execute only the necessary computations to produce the final result.

Lazy transformation is a core principle of Spark's design, enabling it to be efficient, fault-tolerant, and capable of optimizing workloads for distributed data processing.



## What is ACID transaction in Databricks?

ACID is an acronym that stands for Atomicity, Consistency, Isolation, and Durability. ACID transactions are a set of properties that guarantee reliable processing of database transactions. These properties ensure that database transactions are processed reliably, and they are often associated with relational database management systems (RDBMS). While Databricks itself is not a database management system, it can be used in conjunction with various data storage systems, including databases, data lakes, and more.

Let's break down what each of the ACID properties means:

1. **Atomicity:** This property ensures that a transaction is treated as a single, indivisible unit of work. Either all the changes made within a transaction are committed (applied) to the database, or none of them are. If any part of the transaction fails, the entire transaction is rolled back, and the database remains in its original state.
2. **Consistency:** Consistency ensures that a transaction brings the database from one consistent state to another. In other words, it guarantees that the database adheres to a set of integrity constraints (e.g., foreign key relationships, unique constraints) before and after the transaction.
3. **Isolation:** Isolation guarantees that multiple transactions can be executed concurrently without interfering with each other. It ensures that each transaction is isolated from others until it's completed. The levels of isolation, often referred to as isolation levels (e.g., Read Uncommitted, Read Committed, Repeatable Read, Serializable), define the degree to which transactions are isolated from each other.
4. **Durability:** Durability ensures that once a transaction is committed, its changes are permanent and will survive any system failures, such as power outages or crashes. These changes are typically written to durable storage (e.g., disk or solid-state drive) and will be available even if the database system restarts.

In the context of Databricks, the term "ACID transaction" is commonly associated with data lakes and distributed processing frameworks like Delta Lake. Delta Lake is an open-source storage layer that brings ACID transaction capabilities to data lakes. With Delta Lake, you can perform ACID transactions on data stored in distributed file systems like Hadoop Distributed File System (HDFS) or cloud-based storage systems like Amazon S3 or Azure Data Lake Storage.

Databricks, when used with Delta Lake, provides the capability to work with ACID transactions, ensuring that data operations (inserts, updates, deletes) are atomic, consistent, isolated, and durable, even in a distributed and parallel processing environment.

In summary, ACID transactions in Databricks, when combined with Delta Lake, offer reliable and consistent data processing and management in data lakes, providing the same level of transactional integrity that traditional relational databases offer.

## What is Azure Active Directory and it's key feature?

Azure Active Directory (Azure AD) is Microsoft's cloud-based identity and access management service. It is designed to provide secure and seamless authentication and authorization services for applications, services, and resources running in the Microsoft Azure cloud environment and other Microsoft services. Azure AD is a fundamental component of Microsoft's cloud ecosystem, including Azure, Office 365, and various other Microsoft cloud services.

Key features and aspects of Azure Active Directory include:

1. **Identity Management:** Azure AD allows organizations to manage user identities and access control. It includes features such as user provisioning, single sign-on (SSO), multi-factor authentication (MFA), and self-service password reset.
2. **Single Sign-On (SSO):** Azure AD enables users to sign in once with their credentials and gain access to various applications and services without the need to enter their credentials repeatedly. This simplifies the user experience and enhances security.
3. **Multi-Factor Authentication (MFA):** Organizations can enforce additional security measures by requiring users to provide multiple forms of verification, such as a password and a mobile app approval, before gaining access to resources.
4. **Application Access Management:** Azure AD allows organizations to control and manage access to applications, whether they are hosted in the cloud or on-premises. It supports thousands of pre-integrated applications and services.
5. **Identity Federation:** Organizations can establish trust relationships between their on-premises Active Directory and Azure AD, enabling seamless access to cloud resources using existing on-premises credentials.
6. **B2B and B2C Identity Scenarios:** Azure AD supports both business-to-business (B2B) and business-to-consumer (B2C) identity scenarios. This means it can be used for external partner collaboration and customer-facing applications.
7. **Conditional Access:** Conditional Access policies can be defined to control access based on specific conditions, including device compliance, location, and more.
8. **Security and Monitoring:** Azure AD provides security features like Identity Protection, which helps detect and mitigate security risks. It also offers reporting and auditing capabilities to monitor access and identity-related activities.
9. **Integration with Azure Services:** Azure AD is tightly integrated with other Microsoft Azure services, making it easier to secure and manage identity for Azure-hosted applications and resources.
10. **Developer Integration:** Azure AD provides APIs and tools for developers to build identity and authentication features into their applications, whether they are web applications, mobile apps, or API services.

Azure AD is a critical component for securing access to resources in the Azure cloud environment and Microsoft's broader ecosystem. It plays a central role in identity and access management, contributing to a more secure and efficient cloud computing experience for organizations and users.

## What is data Lineage and it's key component and benefit?

Data lineage is a critical aspect of data management and governance that involves tracking and documenting the flow of data as it moves through various stages of a data pipeline or data ecosystem. It provides a comprehensive view of how data is sourced, transformed, and consumed, and it helps organizations understand the origins, transformations, and destinations of their data. Data lineage is valuable for a variety of reasons, including data quality assurance, compliance, troubleshooting, and impact analysis.

Key components and concepts related to data lineage include:

1. **Data Sources:** Data lineage begins with the identification of data sources. These sources can be databases, files, external systems, or data feeds. Understanding where data originates is crucial for tracking its journey.
2. **Data Transformations:** Data often undergoes various transformations as it is processed. These transformations can include data cleaning, aggregation, filtering, and enrichment. Each transformation should be documented in the data lineage to show how data evolves.
3. **Data Movement:** Data may be moved from one location to another, whether that's within an organization's infrastructure or across systems. Understanding the paths data takes is essential.
4. **Data Storage:** Data lineage also includes information about where data is stored at different stages. This can include data warehouses, data lakes, or specific databases.
5. **Consumers:** Data lineage tracks who consumes the data and how they use it. This includes reports, dashboards, analytics tools, and other applications that rely on the data.

Benefits of Data Lineage:

1. **Data Quality Assurance:** Data lineage helps organizations ensure data accuracy and integrity by tracing data anomalies back to their sources. This is critical for data quality management.
2. **Compliance:** For industries with regulatory requirements (e.g., finance, healthcare), data lineage is essential for demonstrating data traceability and compliance with data governance regulations.
3. **Troubleshooting:** When data issues or errors occur, data lineage can quickly pinpoint the source of the problem, making it easier to address and resolve issues.
4. **Impact Analysis:** Data lineage allows organizations to assess the potential impact of changes to data sources, transformations, or data models before implementing those changes. This helps in planning and risk management.
5. **Documentation:** Data lineage serves as comprehensive documentation for an organization's data infrastructure, making it easier for data engineers, analysts, and other stakeholders to understand and work with data.
6. **Data Governance:** Data lineage is a fundamental component of data governance, as it helps establish data stewardship and data ownership practices.
7. **Data Lifecycle Management:** Understanding data lineage is crucial for managing the lifecycle of data, including data archiving, retention, and deletion.

Tools and technologies, both commercial and open-source, are available to help organizations capture and visualize data lineage. These tools automate the tracking and documentation of data movement and transformations, making it more efficient to manage data across its lifecycle.

## What is Wide and Narrow transformations?

In the context of Apache Spark, a widely used distributed data processing framework, "wide" and "narrow" transformations refer to two categories of operations used in Spark's Resilient Distributed Dataset (RDD) and DataFrame APIs. These terms describe how data is processed and transformed across partitions in a distributed computing environment.

### 1. Narrow Transformations:

- Narrow transformations are transformations that operate on a single partition of data at a time. They don't require shuffling or data exchange between partitions.
- They are executed independently on each partition, and the result for each partition is computed based solely on the data within that partition.
- Examples of narrow transformations include **map**, **filter**, **union**, and **local** operations. These transformations do not require data to be reorganized or moved between partitions.
- Narrow transformations are generally more efficient and faster because they don't incur the overhead of data shuffling.

### 2. Wide Transformations:

- Wide transformations, also known as "shuffle transformations," are transformations that involve data exchange and reorganization between partitions. They require the reshuffling and redistribution of data across partitions, and they often involve grouping, aggregating, or sorting operations that require data from multiple partitions.
- Examples of wide transformations include **groupByKey**, **reduceByKey**, **join**, and **sortByKey**. These transformations involve data shuffling and result in a more significant computational overhead.
- Wide transformations are typically more time-consuming and resource-intensive because they involve moving and redistributing data across the cluster. They can also lead to network and disk I/O bottlenecks.

The choice of whether to use narrow or wide transformations depends on your specific data processing needs and efficiency considerations. In general:

- Narrow transformations are preferred when you can achieve the desired result without reshuffling data between partitions. They are typically faster and more efficient.
- Wide transformations are necessary when you need to aggregate or combine data from multiple partitions, and reshuffling is required to produce the correct result. While wide transformations are slower, they can be essential for complex data processing tasks.

Understanding the distinction between narrow and wide transformations is crucial for optimizing the performance of your Spark applications, as minimizing data shuffling can significantly improve processing speed and resource utilization.

## Difference between delta lake storage and data lake storage?

Delta Lake and Data Lake Storage are two different concepts related to data storage and management in big data and cloud computing. Here are the key differences between Delta Lake and Data Lake Storage:

### 1. Data Lake Storage:

- Data Lake Storage is a storage concept that typically refers to a central repository for storing and managing vast amounts of raw and unprocessed data in its native format. It can store structured data, semi-structured data, and unstructured data.
- Data Lake Storage is often implemented using distributed file systems such as Hadoop Distributed File System (HDFS) on-premises or cloud-based storage solutions like Amazon S3, Azure Data Lake Storage (ADLS), or Google Cloud Storage.
- Data in a Data Lake Storage system is often stored in a variety of file formats like Parquet, Avro, JSON, and others, making it suitable for big data analytics and data processing.

### 2. Delta Lake Storage:

- Delta Lake is a storage layer that can be built on top of Data Lake Storage systems to add transactional and ACID (Atomicity, Consistency, Isolation, Durability) capabilities to data lakes.
- Delta Lake provides features like ACID transactions, schema enforcement, data versioning, and data lineage. It ensures data quality, reliability, and consistency for data stored in Data Lake Storage.
- With Delta Lake, you can perform operations like inserts, updates, deletes, and merges on data in a transactional manner, similar to traditional relational databases.
- It uses a proprietary file format and transaction log to manage data changes and versions.

In summary, Data Lake Storage is a general concept for storing large volumes of raw and unprocessed data in a Data Lake, while Delta Lake is a specific technology or storage layer that adds transactional capabilities and data quality features to Data Lakes. Delta Lake is often used in scenarios where data quality, consistency, and reliability are crucial, and it's compatible with various cloud-based and on-premises Data Lake Storage solutions. Delta Lake can be considered an enhancement to Data Lake Storage, providing ACID transaction capabilities for big data workloads.