



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Q1. When loading a CSV file in Spark, the following steps occur under the hood:

The file is split into smaller chunks called partitions, which are distributed across the nodes in the cluster.

Each partition is read and parsed into a DataFrame, which is a distributed collection of data organized into named columns.

The DataFrames from each partition are combined into a single DataFrame, which is returned to the user.

The data in the DataFrame is then optimized for processing using Spark's query optimizer.

The data is then ready to be transformed, analyzed, and used in machine learning models.

Q2. What cluster manager do in spark:

A cluster manager is responsible for managing the resources of a Spark cluster, including allocating and scheduling resources such as CPU and memory, and launching and monitoring the execution of Spark applications. Spark supports several cluster managers, including:

Standalone: The simplest cluster manager included with Spark that runs on a single machine or a cluster of machines. It is easy to set up and can be used for testing and development.

Apache Mesos: A general-purpose cluster manager that can run Spark, Hadoop, and other distributed applications. It supports dynamic resource allocation and can be used in large-scale production environments.

Apache Hadoop YARN: The resource manager used in Hadoop clusters, it allows Spark to run alongside other Hadoop components such as HDFS and MapReduce.

Kubernetes: An open-source container orchestration system that can run Spark and other distributed applications. It supports scaling and self-healing, and can run on-premises or in the cloud.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

The cluster manager you choose will depend on your specific needs, such as the size of your cluster and the types of applications you are running.

Q3: What is driver program in spark

In Spark, a driver program is the main program that creates the SparkContext, runs the main function, and coordinates the execution of tasks across the cluster.

It also communicates with the cluster manager to acquire resources and schedule tasks.

The driver program also converts the high-level operations specified in the application code into a logical execution plan, known as the directed acyclic graph (DAG) of stages, and then it schedules these stages for execution on the cluster.

The driver program is also the program that the user interacts with and where the results of the computation are returned.

It runs on the driver node and it's where you configure the Spark job and submit it to the cluster manager.

In short, the driver program is responsible for:

Creating and initializing the SparkContext

Converting the high-level operations into a logical execution plan

Coordinating the execution of tasks across the cluster

Communicating with the cluster manager to acquire resources and schedule tasks

Returning the results of the computation to the user

Q4: What is the difference between cluster manager and driver program in spark

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

In Spark, a cluster manager and a driver program are two separate components that work together to execute Spark applications.

A cluster manager is responsible for managing the resources of a Spark cluster, including allocating and scheduling resources such as CPU and memory, and launching and monitoring the execution of Spark applications. It also takes care of the cluster scaling and fault-tolerance.

A driver program, on the other hand, is the main program that creates the SparkContext, runs the main function, and coordinates the execution of tasks across the cluster. It is responsible for creating the RDDs (Resilient Distributed Datasets) and DataFrames, and performing transformations and actions on them. It also communicates with the cluster manager to acquire resources and schedule tasks.

To simplify, a cluster manager is responsible for managing the resources of the cluster and scheduling the execution of Spark applications, while the driver program is responsible for coordinating the execution of tasks and returning the results of the computation to the user.

In short, a cluster manager is responsible for:

Allocating and scheduling resources such as CPU and memory

Launching and monitoring the execution of Spark applications

Scaling and fault-tolerance

A driver program is responsible for:

Creating and initializing the SparkContext

Converting the high-level operations into a logical execution plan

Coordinating the execution of tasks across the cluster

Communicating with the cluster manager to acquire resources and schedule tasks

Returning the results of the computation to the user.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Q5: What are the stages in the spark

In Spark, a directed acyclic graph (DAG) of stages is used to represent the logical execution plan of a Spark application. Each stage represents a set of transformations or actions that can be executed in parallel on the cluster. The stages are connected in a DAG, with the output of one stage serving as the input to another stage.

When a Spark application is executed, the driver program converts the high-level operations specified in the application code into a logical execution plan, known as the DAG of stages, and then it schedules these stages for execution on the cluster. The stages are then executed in parallel on the cluster, with the output of one stage serving as the input to another stage. This allows Spark to perform complex operations on large datasets in a distributed and efficient manner.

Q6: Explain the partitioning in spark and waht is default partitioning and how it works

Partitioning in Spark refers to the process of dividing a large dataset into smaller, more manageable chunks called partitions. These partitions are then distributed across the nodes of a cluster, allowing Spark to perform parallel processing on the data.

The default partitioning in Spark is called Hash Partitioning. It works by using a hash function to map each record in the dataset to a specific partition based on the value of one or more specified columns. The number of partitions can be specified by the user, but if not provided, Spark will use a default value based on the number of cores available on the cluster.

Hash partitioning is useful for distributing data evenly across the cluster, but it can lead to uneven partition sizes if the data is not distributed uniformly across the specified columns. This can cause some partitions to be overloaded with data, while others have less data, leading to poor performance.

Another partitioning method is Range partitioning, where the data is partitioned based on the range of values of a specified column. This is useful when the data is already sorted by the specified column, and can lead to more evenly distributed partitions and better performance.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

In addition, Spark also support bucketing partitioning, where the data is partitioned based on the value of a specified column, and then grouped into smaller, fixed-size buckets. This is useful when the data has a large number of distinct values for the specified column.

In summary, partitioning in Spark is the process of dividing a large dataset into smaller, more manageable chunks called partitions, which are then distributed across the nodes of a cluster. The default partitioning method in Spark is Hash partitioning, which uses a hash function to map each record in the dataset to a specific partition based on the value of one or more specified columns, but other partitioning methods are available like Range and bucketing partitioning.

Q7: When u load the file available in blob storage as data frame in azure databricks what happens? does it get moved from adls to databricks as cluster is in databricks. what happened internally

When you load a file available in blob storage as a DataFrame in Azure Databricks, the file is not physically moved from ADLS (Azure Data Lake Storage) to Databricks. Instead, Databricks creates a reference to the file in ADLS, and reads the data from the file into a DataFrame in memory.

Internally, when you use the `spark.read.csv()` or similar method to read a file from blob storage, the driver program in Databricks creates a Spark job that is executed on the cluster. This job reads the file from ADLS using the Azure Blob Storage connector, which is a library that allows Spark to read data from and write data to Azure Blob Storage. The connector uses the Azure Blob Storage API to read the file, and then converts the data into a DataFrame that can be used in Spark.

The data is then distributed across the nodes of the cluster and cached in memory, so that it can be quickly accessed for subsequent processing. This allows Databricks to perform complex operations on large datasets stored in ADLS in a distributed and efficient manner.

Q8: How spark read a large file (petabyte) when file can not be fit in spark's main memory

Spark can read a large file that cannot fit in its main memory using a technique called "external shuffling". External shuffling allows Spark to read and process large datasets that cannot fit in memory by breaking the data into smaller chunks, and processing them separately.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

When Spark reads a large file, it first splits the file into smaller chunks called partitions. Each partition is then read into memory and processed separately. The processed data is then written to disk, rather than being kept in memory.

When the processing for a partition is complete, Spark retrieves the next partition from disk and loads it into memory for processing. In this way, Spark is able to process large files by reading and processing small chunks of data at a time.

In addition, Spark can also read large files that cannot fit in memory by using a technique called "data locality". This means that Spark reads the data from the nodes of the cluster that are closest to the data, rather than reading the data from a central location, which can help to improve the performance of large data reading and processing.

It's important to mention that there are other factors that can impact the performance of large data reading like the type of storage (e.g. distributed file systems, object storage) and the network infrastructure that connects the cluster with the storage.

In summary, Spark can read a large file (petabyte) that cannot fit in its main memory using a technique called external shuffling, which breaks the data into smaller chunks, processes them separately and writes them to disk. Also, Spark can use the data locality, reading the data from the nodes of the cluster that are closest to the data which can improve the performance.

Q9: is pyspark code get converted to java while executing in spark, please explain in detail?

When you run PySpark code, it is translated into Java code before being executed on the Spark cluster. This is done by the PySpark driver program, which converts the PySpark code into a series of Java instructions that can be executed by the Spark engine.

The process starts when the PySpark code is executed on the driver program. The driver program first converts the PySpark code into a directed acyclic graph (DAG) of stages, where each stage represents a set of operations that need to be executed on a specific dataset. Then, the driver program converts this DAG into a series of Java instructions, which are then passed to the Spark engine for execution.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

The Spark engine then schedules the Java instructions to be executed on the cluster. Each stage of the DAG is executed on a specific set of worker nodes in the cluster, and the results are passed to the next stage of the DAG for further processing. The worker nodes use the Java Virtual Machine (JVM) to execute the Java instructions and process the data.

It's worth mentioning that PySpark uses Py4J library to communicate with the Spark engine and it's responsible for the translation of the python code to JVM instructions.

In summary, PySpark code is translated into Java code before being executed on the Spark cluster. This is done by the PySpark driver program, which converts the PySpark code into a series of Java instructions, then schedules the execution of these instructions on the Spark cluster.

Q10: what is difference between `df.count()` and `df.count`

`df.count()` and `df.count` are both used to count the number of rows in a DataFrame, but they are different in how they are implemented.

`df.count()` is a method that counts the number of rows in the DataFrame by collecting all the data to the driver program and counting the number of rows there. This means that the data is read into memory on the driver program, which can cause performance issues if the DataFrame is very large.

On the other hand, `df.count` is an attribute of the DataFrame, it returns the number of rows without collecting the data to the driver program. Instead, it relies on Spark's ability to perform distributed computation and it returns the result by counting the number of rows in each partition and then summing them up. This can be more efficient when dealing with large DataFrames because it avoids the need to move the data to the driver program.

So, the main difference between `df.count()` and `df.count` is that the former is a method that loads all the data into the driver node, while the latter is an attribute that uses the distributed computing capabilities of Spark to count the rows.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Q11. How to do join in spark?

Joining two or more dataframes in Apache Spark can be done using the join method. The method takes two dataframes and a join expression as its arguments and returns a new dataframe that contains the joined data. The join expression specifies the join type and the columns to join on.

Here's an example of how to perform an inner join between two dataframes in Spark:

```
# create two dataframes
df1 = spark.createDataFrame([(1, "John", 25), (2, "Jane", 30), (3, "Jim", 35)],
                             ["id", "name", "age"])
df2 = spark.createDataFrame([(1, "NYC"), (2, "LA"), (3, "DC")],
                             ["id", "city"])

# perform the inner join
join_df = df1.join(df2, df1.id == df2.id, "inner")

# display the joined dataframe
join_df.show()
```

```
▶ (3) Spark Jobs
▶ df1: pyspark.sql.dataframe.DataFrame = [id: long, name: string ... 1 more field]
▶ df2: pyspark.sql.dataframe.DataFrame = [id: long, city: string]
▶ join_df: pyspark.sql.dataframe.DataFrame = [id: long, name: string ... 3 more fields]
```

```
+---+-----+-----+-----+
| id|name|age| id|city|
+---+-----+-----+
| 1|John| 25| 1| NYC|
| 2|Jane| 30| 2|  LA|
| 3| Jim| 35| 3|  DC|
```

You can also perform left, right, and outer joins using the same method by specifying the join type in the join expression. For example, to perform a left join, you can use the following code:

Q12. How to to left join in spark?

Note that in a left join, all the rows from the left dataframe (df1) are preserved and any matching rows from the right dataframe (df2) are included. If there is no matching row in the right dataframe, the resulting row in the joined dataframe will have null values for the columns from the right dataframe.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

```
#left join
join_df = df1.join(df2, df1.id == df2.id, "left")
```

Q13. How to right join in spark

Note that in a right join, all the rows from the right dataframe (df2) are preserved and any matching rows from the left dataframe (df1) are included. If there is no matching row in the left dataframe, the resulting row in the joined dataframe will have null values for the columns from the left dataframe.

```
#right join
join_df = df1.join(df2, df1.id == df2.id, "right")
```

Q14. How to do join optimization in spark

There are several strategies that you can use to optimize joins in Apache Spark:

1. Broadcast Join: If one of the dataframes is small enough to fit in memory, you can use the `broadcast` method to broadcast it to all the executors, which will significantly reduce the network I/O and improve the performance of the join.

```
join_df = df1.join(broadcast(df2), df1.id == df2.id)
```

Partitioning: By partitioning the dataframes on the join key, you can reduce the amount of data that needs to be shuffled over the network, which can significantly improve the performance of the join. You can use the `repartition` method to change the partitioning of the dataframes

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

```
#Join optimization Partition based on id

from pyspark.sql import SparkSession

# create a Spark session
spark = SparkSession.builder.appName("PartitionedJoinExample").getOrCreate()

# create two dataframes
df1 = spark.createDataFrame([(1, "John", 25), (2, "Jane", 30), (3, "Jim", 35)],
                             ["id", "name", "age"])
df2 = spark.createDataFrame([(1, "NYC"), (2, "LA")],
                             ["id", "city"])

# repartition the dataframes on the join key
df1 = df1.repartition(df1.id)
df2 = df2.repartition(df2.id)

# perform the join
join_df = df1.join(df2, df1.id == df2.id)

# display the joined dataframe
join_df.show()
```

Sort-Merge Join: If the dataframes are sorted on the join key, you can use a sort-merge join, which is more efficient than a broadcast join or a repartitioned join. Sorting the dataframes on the join key allows Spark to efficiently perform the join by merging the sorted data.

```
from pyspark.sql import SparkSession

# create two dataframes
df1 = spark.createDataFrame([(1, "John", 25), (2, "Jane", 30), (3, "Jim", 35)],
                             ["id", "name", "age"])
df2 = spark.createDataFrame([(1, "NYC"), (2, "LA")],
                             ["id", "city"])

# sort the dataframes on the join key
df1 = df1.sortWithinPartitions("id")
df2 = df2.sortWithinPartitions("id")

# perform the join
join_df = df1.join(df2, df1.id == df2.id)

# display the joined dataframe
join_df.show()
```

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Q15: explain sorwithinPartion in spark in detail

`sortWithinPartitions` is a method in Apache Spark that is used to sort the data within each partition of a dataframe. Sorting the data within each partition can improve the performance of certain operations, such as joins and aggregations, by reducing the amount of data that needs to be shuffled across the network.

For example, consider a large dataframe that needs to be joined with another dataframe. If the data in the first dataframe is not sorted, Spark would need to shuffle a large amount of data across the network to perform the join. However, if the data in the first dataframe is sorted, Spark can perform the join more efficiently by merging the sorted data.

Here's an example of how you can sort a dataframe using the `sortWithinPartitions` method in Spark:

```
from pyspark.sql import SparkSession

# create two dataframes
df1 = spark.createDataFrame([(1, "John", 25), (2, "Jane", 30), (3, "Jim", 35)],
                             ["id", "name", "age"])
df2 = spark.createDataFrame([(1, "NYC"), (2, "LA")],
                             ["id", "city"])

# sort the dataframes on the join key
df1 = df1.sortWithinPartitions("id")
df2 = df2.sortWithinPartitions("id")

# perform the join
join_df = df1.join(df2, df1.id == df2.id)

# display the joined dataframe
join_df.show()
```

In this example, we create a dataframe `df` and then sort the data within each partition using the `sortWithinPartitions` method with the `id` column as the sort key. Finally, we display the sorted dataframe.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Note that the `sortWithinPartitions` method sorts the data within each partition, not the entire dataframe. The data in each partition is sorted in ascending order by default, but you can specify a different sort order by passing the `ascending` argument.

Q16. What is bucketing and how it can help in join optimization

Bucketing is a feature in Apache Spark that allows you to distribute the data across multiple files or directories based on the values in a specific column, called the bucketing column. Bucketing can help optimize joins in Spark by reducing the amount of data that needs to be shuffled over the network.

When joining two dataframes, Spark needs to shuffle data across the network to ensure that the data is correctly grouped for the join. By bucketing the dataframes on the join key, Spark can ensure that the data is already grouped on the join key and can avoid shuffling the data across the network.

Here's an example of how you can bucket a dataframe in Spark:

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

# create a Spark session
spark = SparkSession.builder.appName("BucketingExample").getOrCreate()

# create a dataframe
df = spark.createDataFrame([(1, "John", 25), (2, "Jane", 30), (3, "Jim", 35)],
                           ["id", "name", "age"])

# bucket the dataframe
#df = df.repartition(1, expr("id")).write.bucketBy(2, "id").sortBy("id").mode("overwrite").parquet("/tmp/bucketed_data").saveAsTable("bucketExample")

df.write.format("parquet").bucketBy(2, "id").sortBy("id").option("path", "/tmp/bucketed_data").saveAsTable("bucketExample")

# read the bucketed data
bucketed_df = spark.read.parquet("/tmp/bucketed_data")

# display the bucketed data
bucketed_df.show()
```

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Q17. how to do incremental data load in databricks

Incremental data load is the process of loading only new or updated data into a data storage system, rather than loading the entire dataset every time. In Databricks, there are several approaches to perform incremental data loads, including using Spark Structured Streaming, Delta Lake, and Databricks SQL.

1. Using Spark Structured Streaming: Spark Structured Streaming is a high-level API for processing real-time data streams in Apache Spark. It can be used to perform incremental data loads by processing new data as it arrives and updating the data in the storage system.
2. Using Delta Lake: Delta Lake is a storage layer that sits on top of your data lake and provides ACID transactions, scalable metadata handling, and data versioning. With Delta Lake, you can perform incremental data loads by appending new data to an existing Delta Lake table.

```
# Load new data into a dataframe
new_data_df = spark.read.parquet("/data/new_data")

# Append the new data to an existing Delta Lake table
new_data_df.write.format("delta").mode("append").save("/delta/events")
|
```

3. Using Databricks SQL: Databricks SQL is a Spark-based SQL engine that allows you to perform SQL-based data processing in Databricks. With Databricks SQL, you can perform incremental data loads by using the **MERGE** statement to update the data in the storage system.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

```
# Load new data into a dataframe
new_data_df = spark.read.parquet("/data/new_data")

# Create a temporary table for the new data
new_data_df.createOrReplaceTempView("new_data")

# Perform a MERGE statement to update the data in the storage system
spark.sql("""
    MERGE INTO events
    USING new_data
    ON events.id = new_data.id
    WHEN MATCHED THEN
        UPDATE SET *
    WHEN NOT MATCHED THEN
        INSERT *
""")
```

Q 18. What is map join in spark

Map join is a join optimization technique in Apache Spark where the data of one of the joining tables is broadcasted to all the executors, while the data of the other table is partitioned and co-located with the corresponding partition of the broadcasted table. In a map join, the join operation is performed in-memory on each executor, reducing the amount of data that needs to be shuffled over the network. This makes map joins more efficient than traditional reduce-side joins for large datasets.

Map joins are particularly effective when one of the joining tables is small enough to fit in memory on each executor. In this case, the join operation can be performed entirely in-memory, avoiding the overhead of shuffling data over the network.

To perform a map join in Spark, you can use the `broadcast` function from the `pyspark.sql.functions` module. Here's an example:

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Cmd 18

```
from pyspark.sql.functions import broadcast

# Load the large table
large_table = spark.read.parquet("/data/large_table")

# Load the small table
small_table = spark.read.parquet("/data/small_table")

# Broadcast the small table
broadcast_small_table = broadcast(small_table)

# Perform the join
result = large_table.join(broadcast_small_table, "id")
```

Q19. how to handle out of memory issue in spark

Handling out of memory issues in Apache Spark can be done through several methods:

1. **Increase Executor Memory:** The first step to handle out of memory issues is to increase the executor memory. You can do this by setting the `spark.executor.memory` configuration property when starting the Spark application.
2. **Decrease the Number of Executors:** Another way to handle out of memory issues is to decrease the number of executors. By reducing the number of executors, you reduce the memory requirements for each executor, which can help mitigate out of memory issues.
3. **Cache RDDs and DataFrames:** If you have RDDs or DataFrames that are frequently used, you can cache them in memory to reduce the overhead of recomputing them. You can use the `persist()` or `cache()` method on the RDD or DataFrame to cache it in memory.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

4. **Use Serialization:** By default, Spark serializes data using the Java Serialization framework, which can be slow and consume a lot of memory. To reduce the memory overhead, you can use a more efficient serialization framework such as Kryo or Avro.
5. **Reduce the Size of RDDs and DataFrames:** Another way to reduce memory usage is to reduce the size of RDDs and DataFrames by filtering out unnecessary data or aggregating data before it is stored in memory.
6. **Use Spark Configuration Properties:** Spark provides several configuration properties that can help you control memory usage. For example, you can use the `spark.shuffle.spill` property to control when Spark spills data to disk during shuffles, which can help reduce memory pressure.
7. **Monitor Memory Usage:** Regularly monitoring memory usage can help you identify and address out of memory issues. You can use Spark's web UI or the Spark Monitoring API to monitor memory usage and identify the source of memory leaks.

In conclusion, handling out of memory issues in Spark requires a combination of adjusting the memory configuration, caching data in memory, using efficient serialization, reducing the size of data, and monitoring memory usage. You can use these methods to resolve out of memory issues and ensure that your Spark applications run smoothly.

Q20. what are the different types of memory in spark

Apache Spark has a multi-layer memory architecture that allows for efficient and scalable processing of large data sets.

1. **Resilient Distributed Datasets (RDDs):** RDDs are the fundamental data structure in Spark, and they are stored in memory as partitions. RDDs are partitioned across multiple nodes in a cluster, allowing for parallel processing. RDDs can be cached in memory for reuse, reducing the overhead of recomputing them.
2. **Executor Memory:** Each node in a Spark cluster runs one or more executors, which are responsible for executing tasks. Each executor has its own dedicated memory, which is used to store intermediate results, RDDs, and other data structures. The amount of memory available to each executor is specified by the `spark.executor.memory` configuration property.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

3. **Storage Memory:** Storage memory is used to store the cached RDDs and other data structures that are frequently accessed. Storage memory is managed by the Spark Block Manager, which caches data in memory and evicts data to disk when necessary. The amount of storage memory available to each executor is specified by the `spark.memory.storageFraction` configuration property.
4. **Task Result Servers:** When a task is completed, its results are stored in a task result server, which is responsible for serving the results to other nodes in the cluster. The task result server stores the results in memory, allowing for fast access to the results by other tasks.
5. **Shuffle Memory:** Shuffle memory is used to store intermediate results during shuffles, which are operations that redistribute data across nodes in a cluster. Shuffle memory is managed by the Spark Shuffle Manager, which spills data to disk when necessary to prevent out of memory issues. The amount of shuffle memory available to each executor is specified by the `spark.shuffle.memoryFraction` configuration property.

In conclusion, Spark's memory architecture is designed to efficiently manage the memory used by executors, storage, task result servers, and shuffles, allowing for scalable processing of large data sets. The memory architecture can be tuned by adjusting the configuration properties, allowing you to balance memory usage between different components and minimize the risk of out of memory issues.

Q21. what is AQE in spark

AQE (Adaptive Query Execution) is a feature in Apache Spark that allows for dynamically optimizing query execution plans based on the data and the resources available in the cluster. AQE was introduced in Spark 2.3 and it aims to improve the performance of Spark SQL queries by dynamically adapting the execution plan based on the data and the resources available in the cluster.

AQE works by monitoring the execution of queries in real-time, and based on the information gathered, it adapts the execution plan to better utilize the resources available in the cluster. For example, AQE can decide to use a broadcast join instead of a shuffle join, or to change the number of partitions used in a join, based on the size of the data and the resources available.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

AQE is transparent to the user and requires no configuration. It automatically adjusts the execution plan to provide the best performance, without the need for manual tuning. This makes it easier to get good performance from Spark SQL queries, even on large and complex data sets.

Overall, AQE provides an important step forward in the evolution of Spark SQL, and it makes it easier for users to get good performance from their queries, even on large and complex data sets.

Q22. what is broadcast variable in spark with pros , cons and example

Broadcast variables in Apache Spark are read-only variables that are used to cache data on each worker node, instead of sending it over the network for every task. The main purpose of broadcast variables is to optimize the execution of Spark jobs by reducing network traffic and improving the performance of the application.

Pros of using broadcast variables in Spark:

1. **Improved Performance:** By caching the data on each worker node, Spark can avoid the overhead of sending the data over the network for every task, resulting in faster execution times.
2. **Reduced Network Traffic:** Since broadcast variables are cached on each worker node, the network traffic is reduced, which can help avoid network bottlenecks and improve overall performance.
3. **Ease of Use:** Spark provides a convenient API for creating and using broadcast variables, making it simple to add them to your application.

Cons of using broadcast variables in Spark:

1. **Increased Memory Footprint:** Since broadcast variables are cached on each worker node, they can consume a significant amount of memory, which may lead to OutOfMemory errors or other performance issues.
2. **Limited Use Cases:** Broadcast variables are best suited for small amounts of data that will be used by multiple tasks. For larger amounts of data, it may be more efficient to use an external data store such as Hadoop HDFS.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Example:

Here is a simple example of using a broadcast variable in Spark to share a lookup table of countries and their capitals across multiple tasks:

```
val countries = Map("United States" -> "Washington, D.C.", "France" -> "Paris")
val broadcastCountries = sc.broadcast(countries)

val data = Seq("United States", "France", "Japan")
val rdd = sc.parallelize(data)

val result = rdd.map(x => (x, broadcastCountries.value(x)))
```

In this example, the broadcast variable `broadcastCountries` is created from the `countries` map, and then used in a Spark RDD `map` operation to lookup the capital for each country. By using a broadcast variable, the `countries` map is cached on each worker node, reducing the amount of data that needs to be sent over the network for each task.

Q23. what is accumulators in spark with pros , cons and example

Accumulators in Apache Spark are variables that can be manipulated within a parallel execution environment, allowing multiple tasks to contribute to their values. They are designed to provide a mechanism for efficiently aggregating data across a distributed dataset.

Pros:

1. Accumulators provide an efficient way to aggregate values in a distributed computing environment.
2. They allow multiple tasks to increment the same accumulator, making it easy to compute a global aggregate value.
3. Accumulators can be used to implement counters, sums, and other common aggregation operations.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Cons:

1. Accumulators are write-only variables and are only updated through Spark operations.
2. They are not accessible from the driver program, so the final result of an accumulator must be retrieved from an action operation.

Example: Here's a simple example of using an accumulator in Spark to count the number of elements in a RDD (Resilient Distributed Dataset):

```
from pyspark import SparkContext

sc = SparkContext("local", "accumulator example")

# Create an RDD
rdd = sc.parallelize([1, 2, 3, 4, 5])

# Create an accumulator with an initial value of 0
accumulator = sc.accumulator(0)

# Define a function to increment the accumulator
def increment_accumulator(x):
    global accumulator
    accumulator += x

# Use the foreach action to increment the accumulator for each element in the RDD
rdd.foreach(increment_accumulator)

# Get the final result of the accumulator
result = accumulator.value

# Print the result
print("The sum of the elements in the RDD is", result)
```

Q24. spark session vs spark context

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Apache Spark provides two main APIs for interacting with Spark: SparkContext and SparkSession. Here's a comparison between the two:

SparkContext:

1. SparkContext is the entry point for Spark functionality. It represents the connection to a Spark cluster and is the primary means for manipulating RDDs (Resilient Distributed Datasets).
2. SparkContext can be created directly in a Spark application, and is used to create RDDs, broadcast variables, and accumulators.
3. Only one SparkContext can be active in a JVM at the same time.
4. SparkContext provides the low-level API for Spark, which is more complex and requires more code to use compared to SparkSession.

SparkSession:

1. SparkSession is a new entry point introduced in Spark 2.0 that provides a higher-level API for Spark functionality and simplifies the process of creating SparkContext.
2. SparkSession automatically creates a SparkContext for you and is used to create DataFrames, SQLContext, and StreamingContext.
3. Multiple SparkSessions can be active in the same JVM.
4. SparkSession provides a more convenient and less verbose way to use Spark compared to SparkContext.

Q25. Explain udf with example in spark

UDF (User-Defined Functions) in Apache Spark are functions that you can define in Spark SQL or Spark DataFrame API and then use in Spark SQL or Spark DataFrame operations. UDFs are written in Scala, Java, or Python and can be used to extend the functionality of Spark SQL or Spark DataFrame API.

Here's an example of using a UDF in Spark SQL to convert a column of strings to uppercase:

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Q26. how executor memory get split up internally in spark

The memory used by an executor in Apache Spark is divided into several regions for different purposes. The following are the major regions and how they are used:

1. **Executor Heap:** This is the main memory region for Spark executors and it is used to store data structures and metadata related to Spark tasks. The size of the executor heap is controlled by the `spark.executor.memory` configuration property.
2. **User Cache:** This memory region is used to cache data and RDDs that are frequently accessed. The size of the user cache can be controlled using the `spark.storage.memoryFraction` configuration property.
3. **Shuffle Memory:** This memory region is used to store intermediate results during the shuffle process. The size of the shuffle memory can be controlled using the `spark.shuffle.memoryFraction` configuration property.
4. **Off-Heap Memory:** Spark allows users to store data off-heap, which is stored outside the JVM heap and can be accessed directly by Spark tasks. Off-heap memory can be used to store large RDDs that cannot fit in the executor heap. The size of off-heap memory can be controlled using the `spark.memory.offHeap.enabled` and `spark.memory.offHeap.size` configuration properties.

It's important to note that these regions are not isolated and their sizes can overlap, depending on the configuration and the amount of memory available. Spark also has a built-in memory manager that monitors the memory usage of the executor and makes adjustments as necessary to ensure that the executor has enough memory to complete its tasks.

In conclusion, the memory used by a Spark executor is divided into several regions for different purposes, including the executor heap, user cache, shuffle memory, and off-heap memory. The size of each region can be controlled using configuration properties.

Q27. Why technically running the udf in pyspark is not recommended. Give detailed reason.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal



Mission100 Azure Data Engineer Course By Deepak Goyal

<https://adeus.azurelib.com>

Email at: admin@azurelib.com

Ask Queries here: <https://www.linkedin.com/in/deepak-goyal-93805a17/>

Running User-Defined Functions (UDFs) in PySpark can be less performant compared to using built-in functions or using Spark SQL functions. There are several reasons why UDFs in PySpark are considered a less efficient option:

1. Overhead of serialization and deserialization: PySpark serializes and deserializes the Python function and its inputs and outputs between the worker nodes and the driver program. This overhead can add significant latency to the execution of the UDF, particularly for large data sets or complex functions.
2. Performance impact of Python code: Python is a dynamically-typed interpreted language, which can be slower than Spark's built-in functions and SQL functions, which are optimized for performance. Additionally, because Python is an interpreted language, Spark has to launch a Python process for each UDF, which can also add to the overhead of UDF execution.
3. Lack of optimization opportunities: Spark's built-in functions and SQL functions are optimized for performance and scalability, and are optimized for the Spark execution engine. In contrast, PySpark UDFs are executed as regular Python functions and are not optimized for the Spark execution engine. This means that UDFs may not perform as well as built-in or SQL functions.
4. Difficulty in debugging: UDFs can be more difficult to debug compared to built-in functions or SQL functions. This is because the error messages generated by UDFs can be less informative, and it can be harder to isolate and diagnose issues with UDFs.

For these reasons, it's generally recommended to use built-in functions or Spark SQL functions whenever possible, and to use UDFs only when they are necessary to solve a specific problem. When using UDFs, it's also important to carefully consider the performance implications of the UDF and to test the performance of the UDF in a development environment before deploying it to production.

Course Link: <https://adeus.azurelib.com>

Email at: admin@azurelib.com

Azurelib Academy By Deepak Goyal