

Justify: SPARK is faster than Map reduce.

In Previous Years Questions

Spark is faster than MapReduce for several reasons

1. In-memory processing

Spark primarily processes data in memory (RAM), while MapReduce primarily processes data on disk. This is because Spark uses Resilient Distributed Datasets (RDDs), which are in-memory collections of data that can be partitioned across multiple nodes in a cluster. This allows for much faster processing compared to MapReduce, which needs to read and write data from disk for each operation.

2. Iterative algorithms

Spark is well-suited for iterative algorithms, which involve multiple passes over the data. This is because Spark can cache the intermediate results of each iteration in memory, which can be accessed much faster than reading from disk. MapReduce, on the other hand, is less efficient for iterative algorithms as it needs to write the intermediate results to disk after each iteration, which can be a bottleneck.

3. Pipelines and DAGs

Spark allows you to build complex data processing pipelines using Directed Acyclic Graphs (DAGs). This allows you to optimize the execution of your jobs by running multiple tasks in parallel. MapReduce, on the other hand, has a more rigid two-stage execution model, which can be less efficient for complex workflows.

Justify: SPARK is faster than Map reduce.

4. General-purpose engine

Spark is a general-purpose engine that can be used for various data processing tasks, including machine learning, graph processing, and real-time analytics. MapReduce, on the other hand, is primarily designed for batch processing tasks.

Here are some benchmarks that demonstrate Spark's performance advantage over MapReduce:

- Sort: Spark can sort 100 TB of data in 23 minutes, while MapReduce takes 3 hours.
- Word count: Spark can count the words in 100 TB of text data in 3 minutes, while MapReduce takes 2 hours.
- Machine learning: Spark can train a machine learning model on 100 GB of data in 10 minutes, while MapReduce takes 1 hour.

Overall, Spark is significantly faster than MapReduce for a variety of data processing tasks. This is due to its in-memory processing, efficient iterative algorithms, support for pipelines and DAGs, and general-purpose nature.

However, it is important to note that MapReduce still has its place. It is more mature and stable than Spark, and it can be a good choice for simple batch processing tasks where speed is not a critical factor.