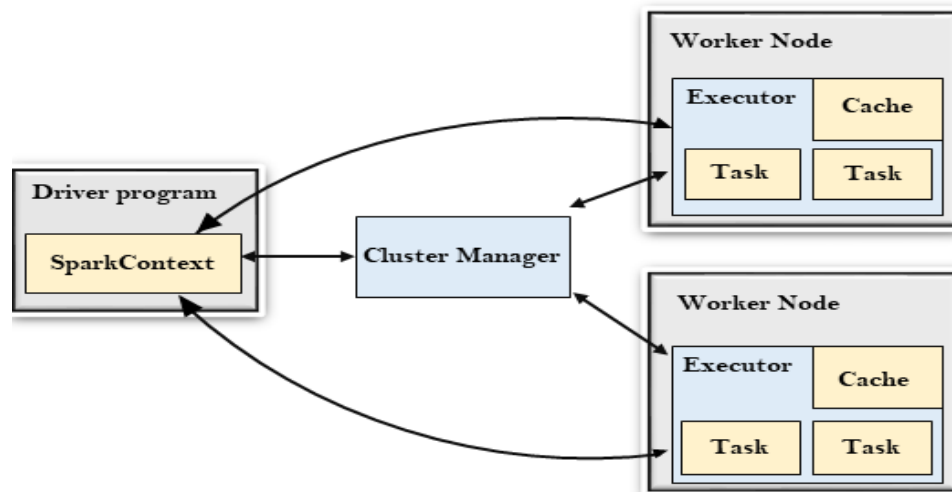


Master Spark Concepts Zero to Big Data Hero:

Detailed Notes on Spark Architecture and Execution Flow

Apache Spark is a powerful open-source distributed computing system that enables fast and efficient data processing. Here's a quick overview of its architecture to help you understand how it works:



Key Components of Spark Architecture

- **Driver Program**

Description: The central coordinator that converts user code into tasks.

Role: Manages the execution of the Spark application and maintains information about the status of tasks.

- **Cluster Manager**

Description: Manages resources across the cluster.

Types: Standalone, YARN, Mesos, Kubernetes.

Role: Allocates resources and schedules tasks on the cluster.

- **Executors**

Description: Workers that run the tasks assigned by the driver program.

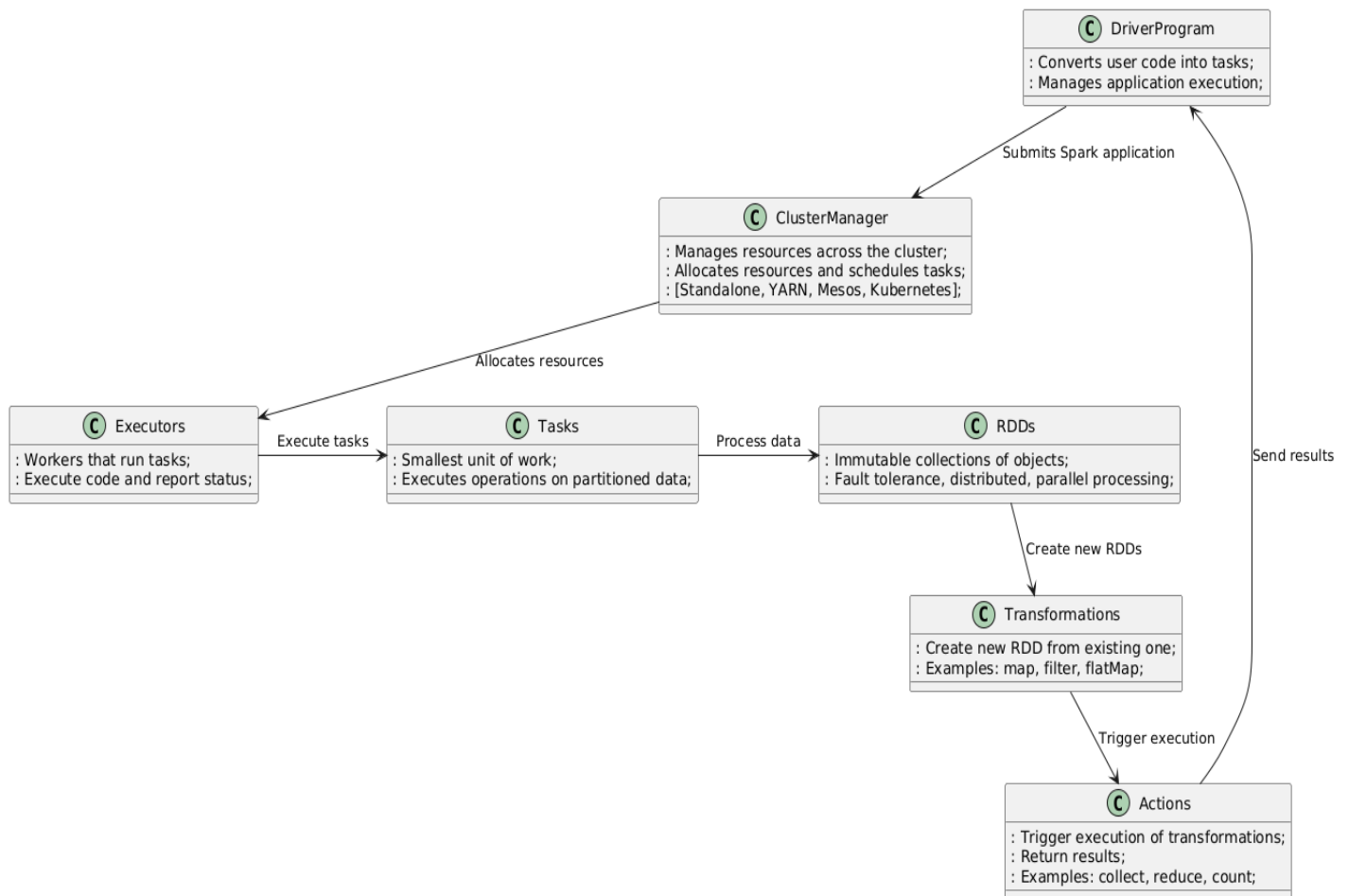
Role: Execute code and report the status of computation and storage.

- **Tasks**

Description: The smallest unit of work in Spark.

Role: Executes individual operations on the partitioned data.

Apache Spark Architecture



Data Processing in Spark

- RDDs (Resilient Distributed Datasets)**

Description: Immutable collections of objects that can be processed in parallel.

Features: Fault tolerance, distributed processing, and parallel execution.

- Transformations**

Description: Operations that create a new RDD from an existing one.

Examples: map, filter, flatMap.

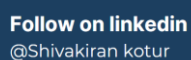
- Actions**

Description: Operations that trigger the execution of transformations and return a result.

Examples: collect, reduce, count.



Follow on linkedin
@Shivakiran kotur



The execution flow in Apache Spark outlines how data processing occurs from the initial job submission to the final result collection. Here's a step-by-step breakdown:

1. SparkContext Creation

- The driver program starts and initializes the SparkContext.
- This context serves as the main entry point for Spark functionality and manages the entire Spark application.

2. Cluster Manager Interaction

- The SparkContext interacts with the Cluster Manager (e.g., YARN, Mesos, or Standalone).
- It requests resources (CPU, memory) needed for the application to run.

3. Job Submission

- The user defines transformations and actions within the Spark application.
- A job is created when an action is called (e.g., collect, count).
- At this point, Spark begins to build a logical execution plan.

4. DAG Scheduler

- The Directed Acyclic Graph (DAG) Scheduler takes the job and breaks it down into stages.
- Each stage corresponds to a set of transformations that can be executed in parallel.

5. Task Scheduler

- The Task Scheduler takes the stages defined by the DAG Scheduler and converts them into tasks.
- It distributes these tasks to the executors based on data locality, optimizing resource usage by minimizing data transfer.

6. Executor Execution

- The tasks are executed on the worker nodes by the executors.
- Executors process the data according to the tasks assigned and handle intermediate data storage as required.

7. Data Shuffling

- If operations like reduceByKey or join are performed, data shuffling may occur.
- This involves redistributing data across the cluster, which can be resource-intensive.

8. Execution and Monitoring

- The driver program continuously monitors the execution of tasks.
- It handles any failures that may occur during processing, re-executing tasks if necessary.

9. Completion

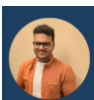
- Once all tasks are completed, the results are sent back to the driver program.
- The SparkContext is now ready to process more jobs, maintaining a seamless workflow.

Conclusion

Apache Spark's architecture is designed to handle large-scale data processing efficiently and effectively. Understanding its components and workflow can help you leverage its full potential for your big data projects.

Important Interview Question from previous post

1. Explain Hadoop Architecture?
2. How MapReduce Works?
3. Difference between MapReduce and Spark?
4. Why Spark is better than MapReduce?
5. What are the components of Spark?
6. What you mean by JVM in Spark and what are its component?
7. What is use of Driver node and Work Node?
8. Explain Spark Architecture?
9. Explain the flow of execution in Spark?



Follow on linkedin
@Shivakiran kotur