# Apache Spark Interview Guide : Core Concepts

☑ Jobs, Stage and Task

Deepa Vasanthkumar

# Job, Stage and Tasks

The concepts of job, stage, and tasks are fundamental to understanding how Spark manages distributed computing. These elements are part of the execution model that Spark uses to distribute and manage computation across multiple nodes efficiently. Let's break down what each of these terms means and how they relate to each other in the context of executing a PySpark program.

## Job

A job in PySpark is the entire process needed to complete a computation initiated by an action. When you execute an action (like `collect()`, `count()`, `save()`, etc.) on a RDD (Resilient Distributed Dataset) or a DataFrame, Spark submits a job. The job represents the entire piece of computation or query from start to finish. For example, if you have a PySpark DataFrame and you want to count the number of rows that satisfy a particular condition, you might use the `count()` action after a `filter()` transformation. The execution of this `count()` action triggers a job.

Logical Unit of Work: A job in Spark represents a complete computation triggered by an action on a RDD or DataFrame. Jobs are composed of one or more stages, each containing multiple tasks. They provide a high-level abstraction for users to define and manage computations on distributed datasets.

Dependency Management: Jobs manage the dependencies between transformations and actions in a Spark application. Spark constructs a directed acyclic graph (DAG) of operations, where each node represents a transformation, and each

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

edge represents a dependency. Jobs orchestrate the execution of these operations in the correct order, ensuring consistency and correctness of the computation.


Optimization and Monitoring: Jobs enable Spark to optimize the execution plan by applying various optimizations, such as task scheduling, data locality, and memory management. Users can monitor job execution through the Spark UI and control job submission and execution parameters programmatically or through configuration settings.


## Stage


A stage is a set of tasks in a job that can be executed in parallel but must be completed before the next set of tasks (the next stage) can begin. Stages in a Spark job are created based on the transformations applied to the RDD or DataFrame. Specifically, stages are divided by transformations that cause shuffling of data across partitions, such as `groupBy()` or `reduceByKey()`.


Shuffle boundaries (where data must be redistributed across different nodes) are a common reason for new stages. Each stage consists of tasks that execute the same code on different partitions of the data, and all tasks in a stage must be completed before any tasks in subsequent stages start.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Execution Pipelining: Stages are logical units of work in Spark that represent a series of transformations that can be executed together without shuffling data between tasks. Spark breaks down a computation into stages based on data dependencies and transformations. By organizing tasks into stages, Spark can optimize the execution plan by pipelining transformations and minimizing data shuffling.

Shuffle Boundary: Stages are often delineated by operations that require data shuffling, such as `groupByKey()` or `join()`. When a shuffle occurs, Spark creates a new stage, allowing it to manage data redistribution efficiently. Stages help in managing the complexity of distributed computations and enable Spark to optimize performance by minimizing shuffle operations.

## Task

A task is the smallest unit of work that gets sent to an executor. Essentially, a task applies the job's logic to a single partition of the data. Tasks within a stage perform the same operations on different pieces of the data and can be executed in parallel across different executors and nodes in the Spark cluster. Tasks are what actually execute the code on the cluster.

Parallel Execution: Tasks are the basic units of computation in Spark. They allow for parallel execution of operations on distributed data partitions across multiple executor nodes in a cluster. By breaking down computations into smaller tasks, Spark can utilize the processing power of the entire cluster efficiently.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Data Partitioning: Tasks operate on individual partitions of the distributed datasets. This partitioning enables Spark to process large datasets that may not fit into the memory of a single machine. Each task processes a subset of the data independently, facilitating parallelism and scalability.

Fault Tolerance: Spark achieves fault tolerance at the task level. If a task fails during execution (due to hardware failures, network issues, etc.), Spark can rerun that specific task on another node using the data lineage information. By isolating failures to individual tasks, Spark minimizes the impact of failures on the overall computation.

Demonstrating with an example:

Here's a simple example to illustrate how jobs, stages, and tasks might be structured in a PySpark program:

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

```
file_path = "file in s3 or local "



data = spark.read.text(file_path)



words = data.rdd.flatMap(lambda line: line.value.split())



words_starting_with_a = words.filter(lambda word: word.startswith('a'))



count = words_starting_with_a.count()



print("Number of words starting with 'a':", count)
```

**Spark Core Concepts    -  7 -**

1. RDD Creation: Suppose you read in a large dataset into an RDD.

2. Transformation: You apply a `map()` transformation to modify the data.

3. Another Transformation: You then apply a `filter()` transformation to remove some rows.

4. Action: Finally, you call `count()` to get the number of elements in the transformed RDD.

In this example:

- The call to `count()` triggers a job.

- This job may consist of one stage if there's no shuffle needed (if the transformations only involve operations like `map()` and `filter()` that do not necessitate a shuffle).

- The stage will have as many tasks as there are partitions in the RDD. Each task processes one partition of the data independently on different executors.

If there was a shuffle transformation involved (like `reduceByKey()`), then there would be at least two stages:

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

- One stage to process the transformations before the shuffle.


- Another stage to process the shuffle and any transformations after the shuffle.


Understanding these components helps in optimizing and debugging PySpark applications, as you can identify where bottlenecks might occur (like excessive shuffling or uneven task distribution) and tune your Spark application accordingly.


In PySpark, jobs that involve a shuffle operation are more complex in terms of how they are divided into stages and tasks. This complexity arises because shuffle operations require redistributing data across different nodes in the cluster to group it differently according to the keys involved in the shuffle. Let's delve into how stages and tasks are structured when a shuffle operation is part of a PySpark job.


Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Shuffle in Spark

A shuffle occurs when data needs to be redistributed across different executor nodes in a Spark cluster. This redistribution is often necessary when performing "wide" transformations such as `groupByKey()`, `reduceByKey()`, `join()`, and others that involve grouping or aggregating data across partitions. The shuffle operation involves writing data to disk and transferring it over the network, making it one of the most resource-intensive operations in Spark.

Stages in Jobs with Shuffle

In jobs that involve a shuffle, stages are generally divided around the shuffle points:

1. Pre-Shuffle Stage: This is the stage before the shuffle occurs. In this stage, tasks perform transformations on the data within their respective partitions. These transformations are "narrow," meaning they do not require data from other partitions. Examples include `map()`, `filter()`, and `flatMap()`.

2. Shuffle Write: At the end of the pre-shuffle stage, tasks write their results to disk. This output is organized such that data destined for the same partition in the next stage is grouped together.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

3. Shuffle Read: The next stage begins with a shuffle read. Executors fetch the relevant parts of data from across all nodes. This is where the heavy lifting of shuffle happens: data needs to be pulled from potentially all other nodes in the cluster, depending on the keys involved.

4. Post-Shuffle Stage: After the shuffle read, the next set of tasks can begin. These tasks operate on the newly formed partitions, which may now contain data from multiple original partitions. The operations in this stage are also "wide," meaning they generally involve data that needs to be aggregated or joined across the original partitions.

Tasks in Jobs with Shuffle

Tasks are the units of execution that carry out the computation:

1. Pre-Shuffle Tasks: Each task in the pre-shuffle stages processes a partition of the data and performs narrow transformations. The results are then written to local disks in a format ready for shuffling.

2. Shuffle Write Tasks: As part of the end of a pre-shuffle task, the data is organized and written out to disk to prepare for shuffling.

3. Shuffle Read Tasks: During the shuffle read phase, each task is responsible for fetching the appropriate data segments from the shuffle write outputs of tasks in the previous stage.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

4. Post-Shuffle Tasks: These tasks process the gathered data. For example, if the transformation is `reduceByKey()`, each task now has data from all partitions related to specific keys and can perform the reduction operation.