

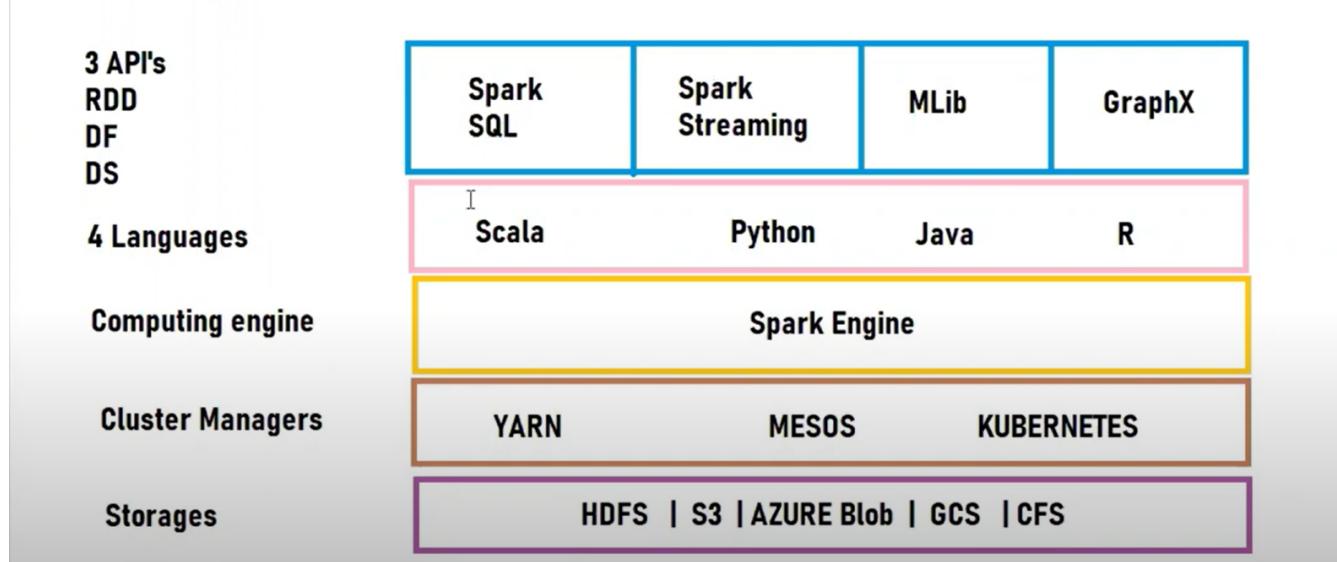
What is Spark?

In memory distributed data processing framework/engine.

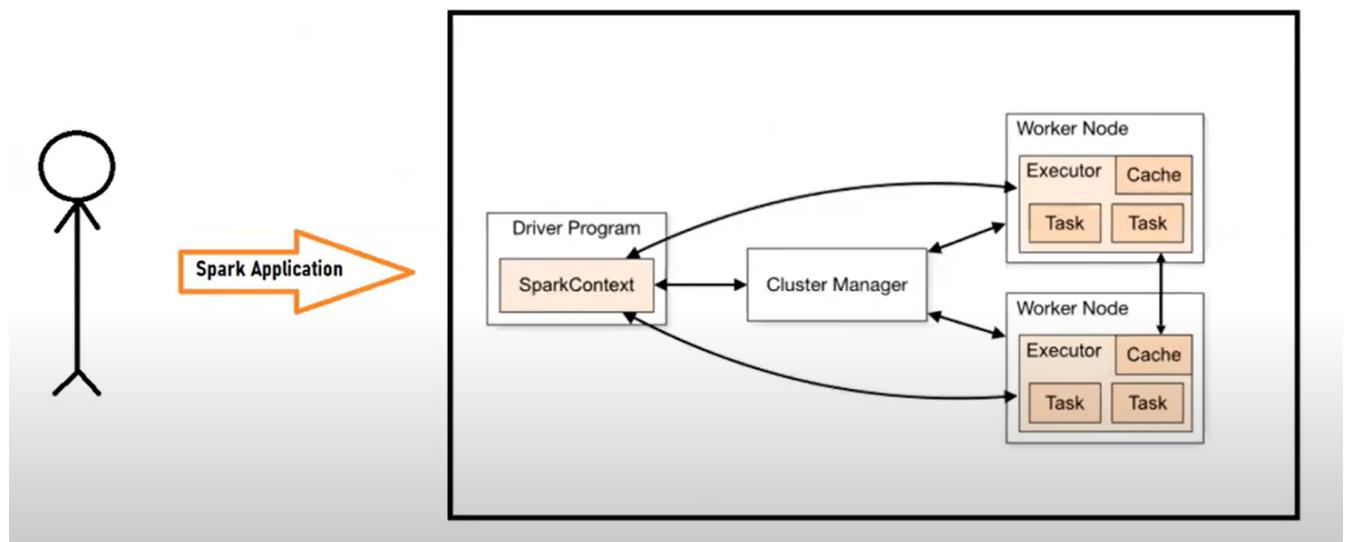
Spark Features:

- In memory computation
- Open source
- Cost effective
- Fault-tolerance
- Supports multiple languages (Python, Scala, Java, and R)
- Lazy Evaluation
- Batch/Near real-time data streaming
- Very rich built-in library support
- As compared to Java (20–25 lines of code), we can replace with single line code in Python/Scala
- 100 times faster in memory than Hadoop MR or other traditional systems
- 10 times faster on disk
- Spark can integrate with Hadoop ecosystem, ETL tools like Talend, Informatica, etc., and cloud platforms (AWS, Azure, GCP)

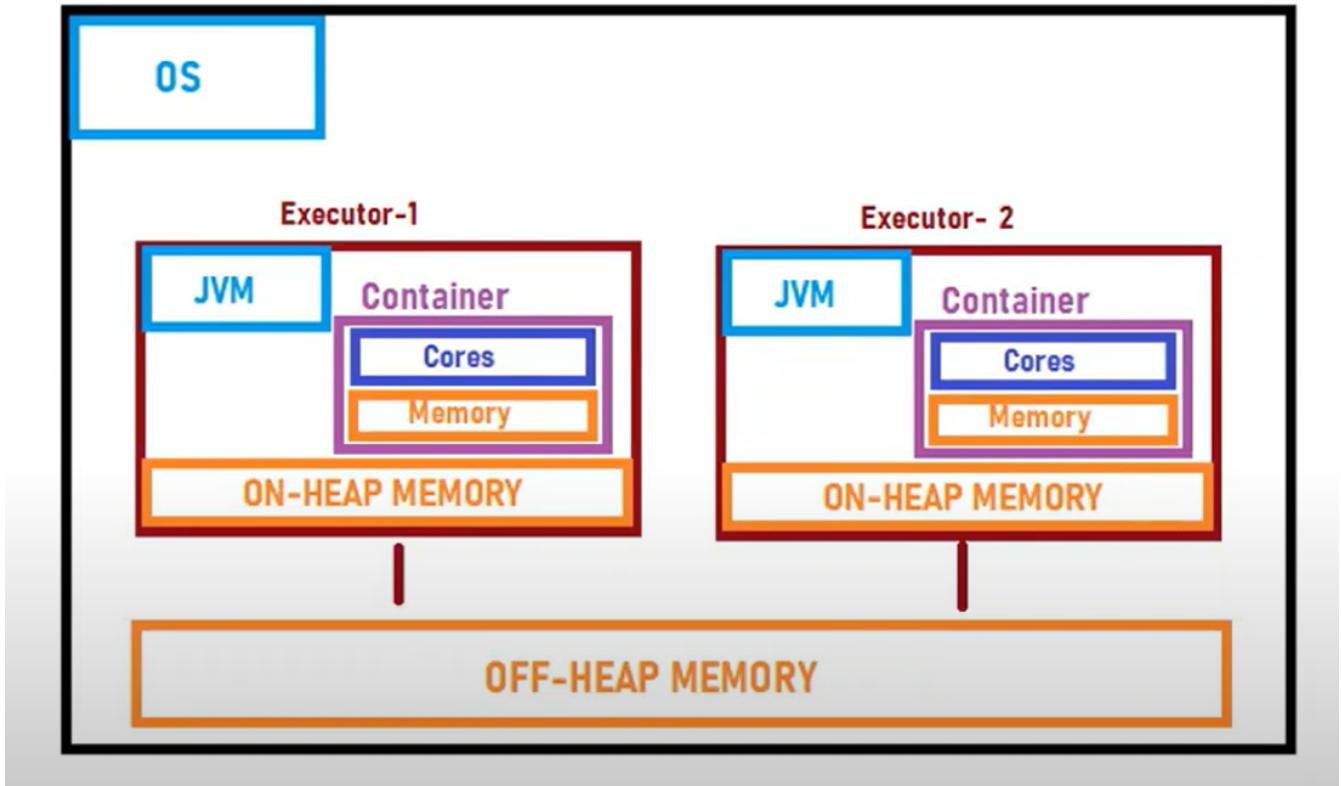
Spark Components :



Spark Application Submission-Cluster Mode



Worker Node / VM

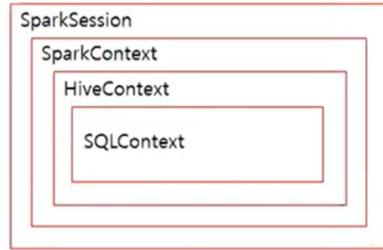


On-Heap vs Off-Heap Memory

Aspect	On-Heap Memory	Off-Heap Memory
Location	Inside JVM Heap	Outside JVM Heap (native memory)
Management	JVM Garbage Collector	Spark (Tungsten engine)
Default	Yes	No (needs config)
Performance	Slower for large data (GC overhead)	Faster for large data (no GC)
Storage Format	Java objects	Serialized binary
Use Cases	Small to medium datasets, default Spark jobs	Large datasets, Spark SQL/DataFrame, caching, joins, shuffles
Risk	GC pauses, OutOfMemoryError	Native memory leaks if misused
Configuration	No extra setup	<code>spark.memory.offHeap.enabled=true</code> and <code>spark.memory.offHeap.size</code>
Example	Default caching <code>.cache()</code>	Optimized Tungsten memory, serialized storage

I. Spark Introduction:

1.7. Spark Session Vs SparkContext



- Prior to 2.0 version of Spark, **SparkContext** is an entry point to Spark programming with **RDD** and to connect to **Spark Cluster**, Since Spark 2.0 **SparkSession** has been introduced and became an entry point to start programming.
- Prior to 2.0 version of Spark, before the creation of SparkContext, SparkConf must be created.
- Prior to 2.0 version of Spark, we have to use a distinct context for each API. We required StreamingContext for Streaming, SQLContext for SQL, and HiveContext for Hive.
- From Spark 2.0, **SparkSession** provides a common entry point for a Spark application. Instead of SparkContext, HiveContext, and SQLContext, everything is now within a SparkSession.

2. Spark RDD's:

2.1.0. Spark RDD – Resilient Distributed Dataset

PySpark RDD (Resilient Distributed Dataset) is a fundamental data structure of PySpark that is fault-tolerant, immutable distributed collections of objects, which means once you create an RDD you cannot change it. Each dataset in RDD is divided into logical partitions, which can be computed on different nodes of the cluster.

RDD Creation

In order to create an RDD, first, you need to create a [SparkSession which is an entry point to the PySpark application](#).

Spark session internally creates a SparkContext

```
# Import SparkSession
from pyspark.sql import SparkSession

# Create SparkSession
spark = SparkSession.builder \
    .master("local[1]") \
    .appName("SparkByExamples.com") \
    .getOrCreate()
```

2. Spark RDD's:

2.1.0. Spark RDD – Ways of Creation

SparkContext has several functions to use with RDDs.

using parallelize()

parallelize() method is used to create an RDD from a list.

```
>>> # Create RDD from parallelize
>>> dataList = [("Surya", 10,"Ind"), ("Butler", 150,"England"), ("Rahul",80,"Ind")]
>>> rdd_player=spark.sparkContext.parallelize(dataList)
>>> rdd_player.collect()
[('Surya', 10, 'Ind'), ('Butler', 150, 'England'), ('Rahul', 80, 'Ind')]
```

using textFile()

RDD can also be created from a text file using textFile() function of the SparkContext..

```
>>> rdd2 = spark.sparkContext.textFile("C:/Users/BNK/Dropbox/Clever Studies/Pyspark/rdd2.csv")
>>> rdd2.collect()
['karuna,ECE,90', 'Varun,EEE,86', 'Spandan,IT,88']
```

2. Spark RDD's:

2.1.1. Spark RDD – Operations

In case of MR



In case of Spark



2. Spark RDD's:

2.1.0. RDD Operations in PySpark

RDD supports two types of operations.

I. Transformations:

Transformations are the operations which are applied to an RDD to create a new RDD.

Transformations follow the principle of Lazy Evaluations.

Lazy Evaluation: RDD's execution will not start until an action is triggered

Few transformations are,

map

flatMap

filter

distinct

reduceByKey

sortBy

o 2. Spark RDD's:

2.1.0. RDD Operations in PySpark

RDD supports two types of operations.

2. Actions:

Actions are the operations which are applied on an RDD, which return a value to the driver program after running a computation on the dataset.

Few actions are,

collect

reduce

countByKey/countByValue

foreach(func)

take

First

saveAsTextFile(path)

Spark Transformations and Actions:

1. Transformations:

Transformations are the operations which are applied to an RDD to create a new RDD. Ex:

map

flatMap

reduceByKey

2. Actions:

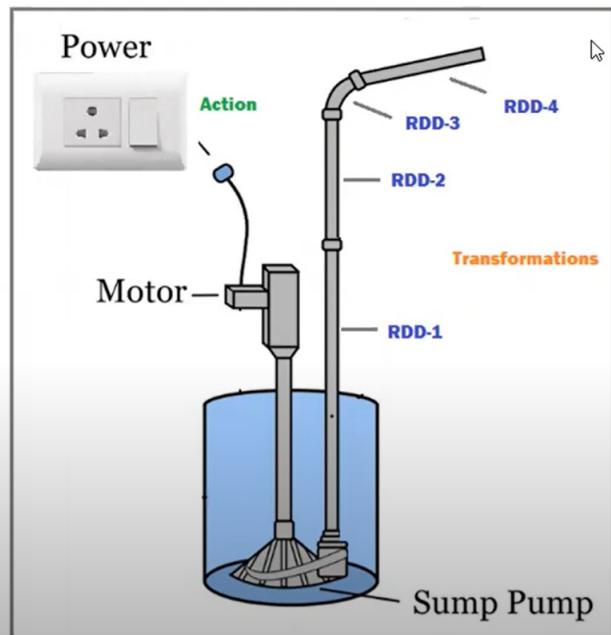
Actions are the operations which are applied on an RDD, which return a value to the driver program after running a computation on the dataset. Ex:

collect

reduce

countByKey/countByValue

foreach(func)

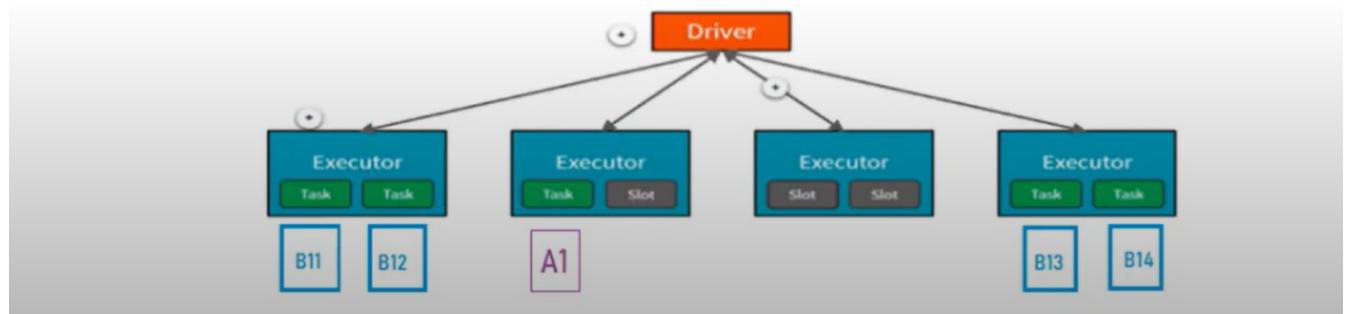
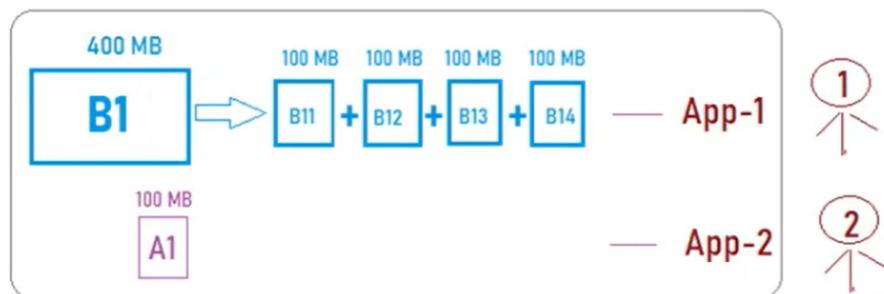


2. Spark Repartition() vs Coalesce()

1. repartition() is used to increase or decrease the RDD/DataFrame partitions
2. coalesce() is used to only decrease the number of partitions in an efficient way.

PySpark repartition() and coalesce() are very expensive operations as they shuffle the data across many partitions hence try to minimize using these as much as possible.

2. Spark partitionBy():



Shuffle and Combiner :

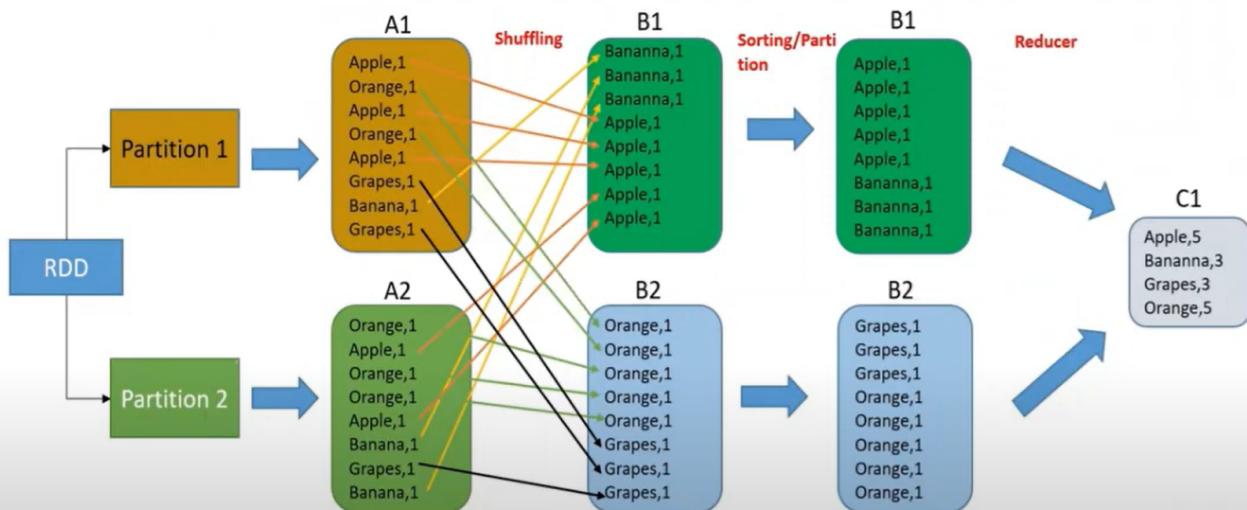
Shuffle:

1. Shuffling is a process of redistributing data across partitions.
2. Shuffling will create a new stage in spark
3. Shuffling is a costly operation as it involves disk I/O , network I/O and data serialization/deserialization.
4. Spark reduceByKey(),groupByKey(),aggregateByKey(),join(), union() etc will involve shuffling.
5. Distinct operation will creates a shuffle.
6. If shuffling is necessary,then use combiner.
7. In combiner,shuffling may happen but it is in lesser side.
8. Since,groupByKey() won't use 'combiner',try to avoid to use this where ever is possible.Instead you can use reduceByKey() and aggregateByKey() if possible.

Combiner :

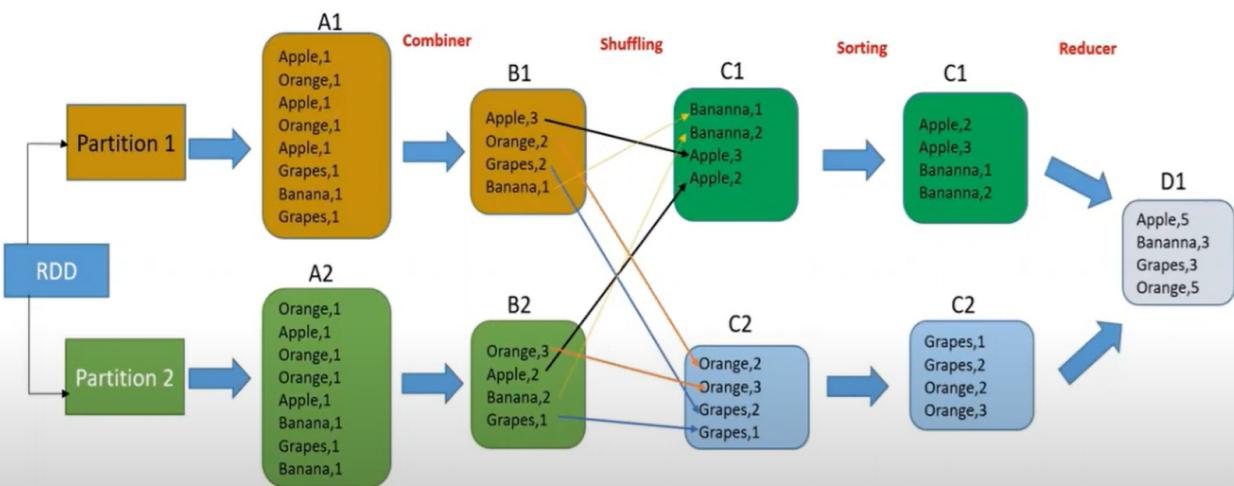
It computes intermediate values for each partition to avoid partial shuffling.

Shuffle without Combiner :



** `groupByKey()` will follow this approach.

Shuffle with Combiner :



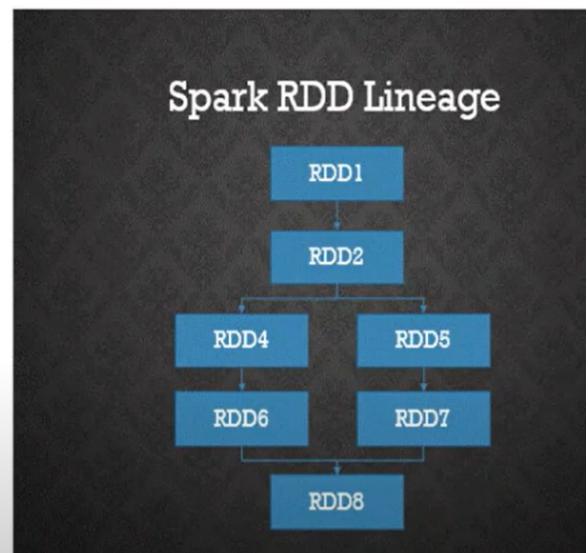
** `reduceByKey()`, `aggregateByKey()` will follow this approach.

Spark Type's of Transformations:

Transformations are the operations which are applied to an RDD. Since RDD are immutable in nature, transformations always create new RDD without updating an existing one hence, this creates an RDD lineage.

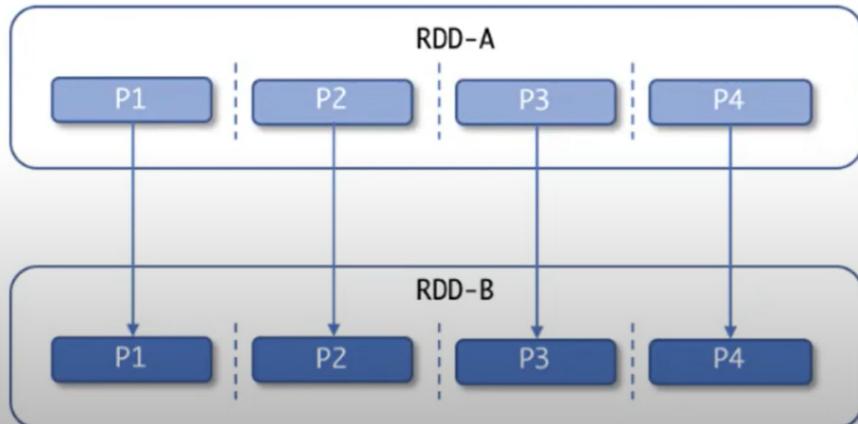
Two types of transformations:

- 1.Wide Transformations
- 2.Narrow Transformations



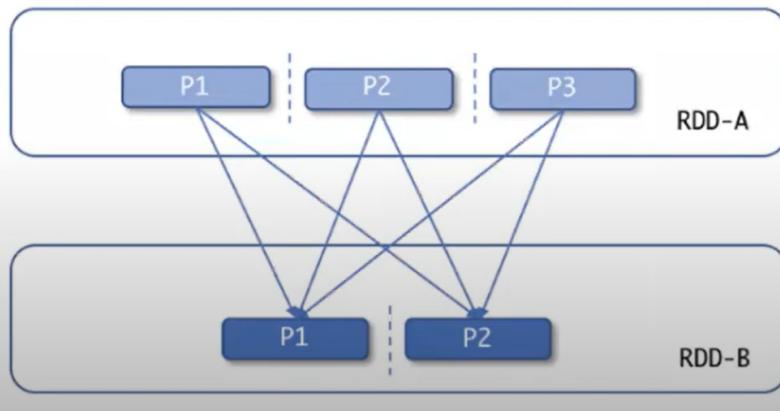
Spark Narrow Transformation:

- ✓ Narrow transformations transform data without any shuffle involved. These transformations transform the data on a per-partition basis; that is to say, each element of the output RDD can be computed without involving any elements from different partitions. This leads to an important point: The new RDD will always have the same number of partitions as its parent RDDs.
- ✓ The map(), filter() are some examples of Narrow transformations.
- ✓ This kind of transformation is basically fast.



Spark Wide Transformation:

- ✓ Wide transformations involve a shuffle of the data between the partitions. In the case of these transformations, the result will be computed using data from multiple partitions and thus requires a shuffle. Wide transformations are similar to the shuffle-and-sort phase of MapReduce.
- ✓ The `groupByKey()`, `reduceByKey()`, `join()`, `distinct()`, and `intersect()` are some examples of wide transformations.
- ✓ Slow as compare to narrow dependencies speed might be significantly affected as it might be required to shuffle data around different nodes when creating new partitions.

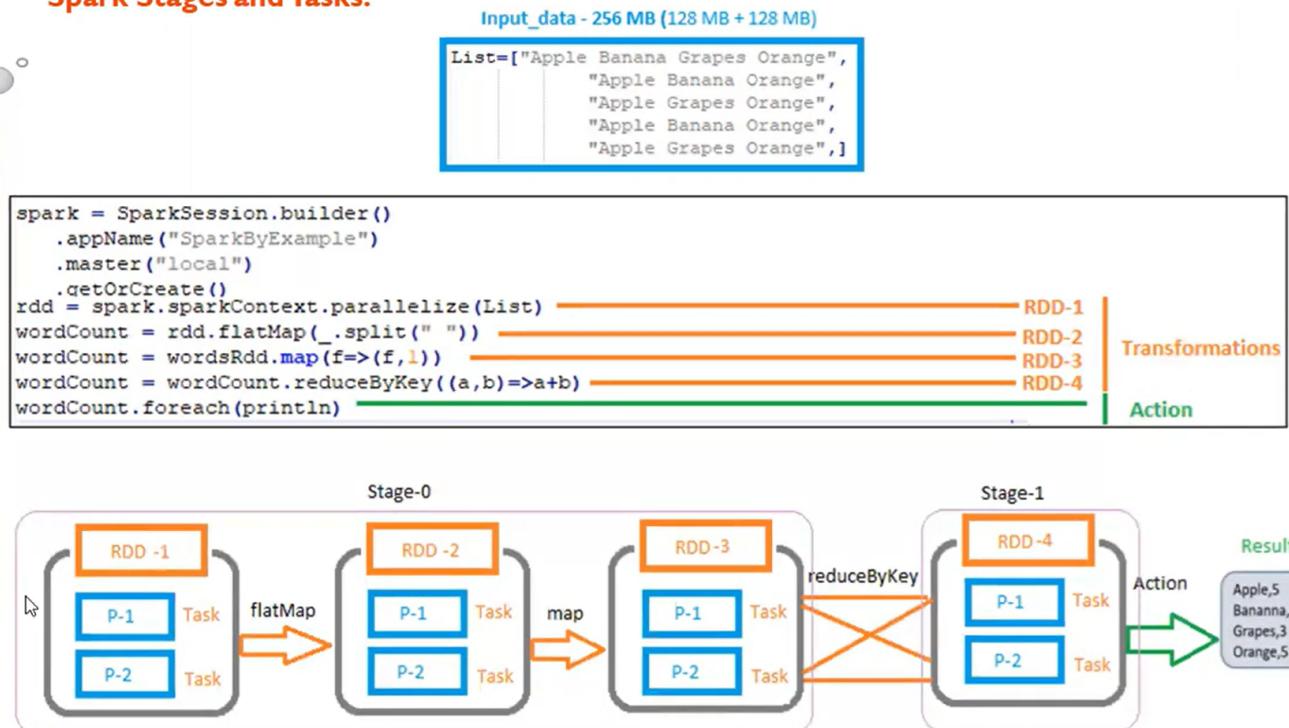


To understand, how spark executes a job internally

We must know the below concepts,

- ◆ RDD - An RDD is an immutable distributed collection of elements of a data, partitioned across nodes in the cluster
- ◆ Partition - A Partition is a single slice (<= 128 MB chunks of data) of a larger Dataset
- ◆ Transformations and Actions
- ◆ Types of Transformations
 - ◆ 1. Narrow Transformations
 - ◆ 2. Wide Transformations
- ◆ Spark Job- It is a single computation action that gets instantiated to complete a Spark Action.
- ◆ Tasks - The single computation unit performed on a single data partition is called a task. It is computed on a single core of the worker node.
- ◆ Stages - A stage is a set of independent tasks all computing the same function that need to run as part of a Spark job

Spark Stages and Tasks:



Note in above code, when we are using reduceByKey it involves shuffling, so it creates a new stage(step).

Spark Stages and Tasks:

- I. A graphical user interface that allows us to view all our Spark Application's jobs, stages, and tasks

The screenshot shows the Apache Spark 3.3.0 Web UI. At the top, there are tabs for Jobs, Stages, Storage, Environment, Executors, and SQL / DataFrame. The 'Jobs' tab is selected. Below it, the 'Spark Jobs (?)' section displays the following information:

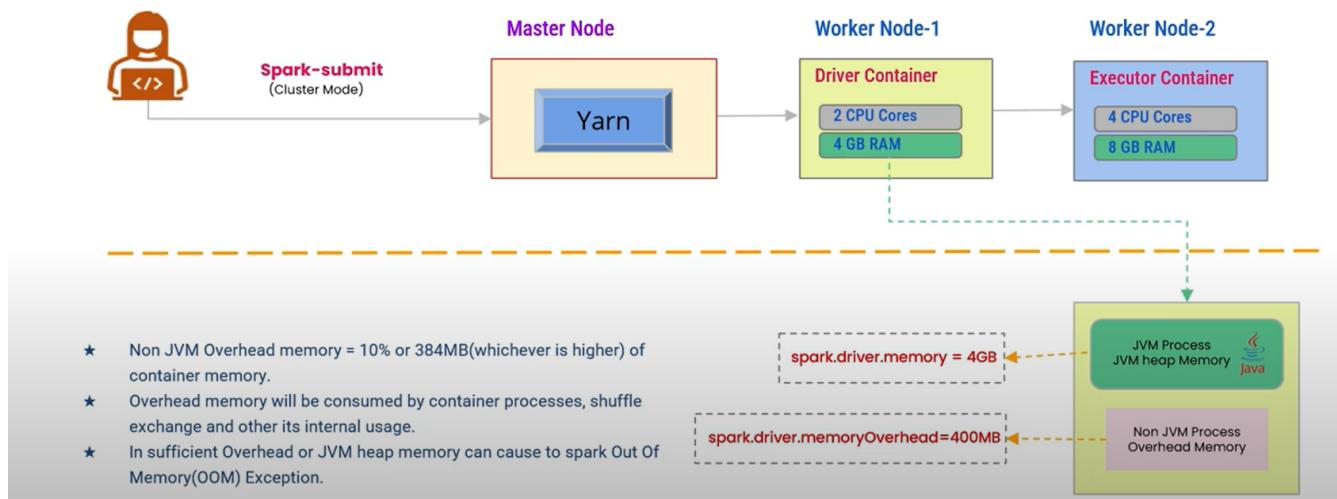
- User: tanmay
- Total Uptime: 10 s
- Scheduling Mode: FIFO
- Completed Jobs: 1

Below this, the 'Completed Jobs (1)' section lists one job:

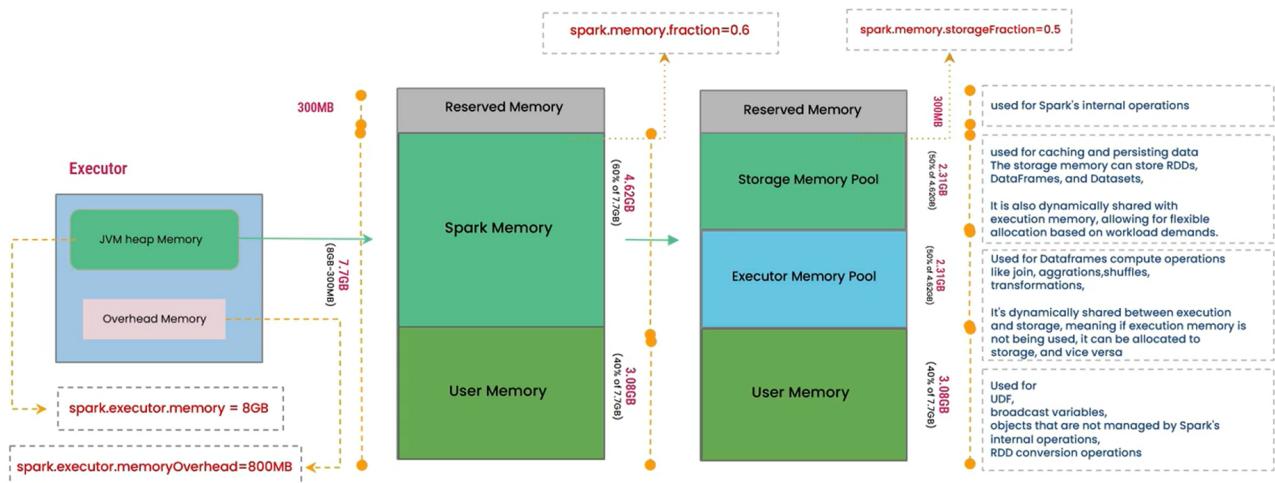
Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0	collect at /tmp/ipykernel_48321/1034515582.py:12 collect at /tmp/ipykernel_48321/1034515582.py:12	2022/09/14 14:52:01	2 s	2/2	5/5

At the bottom of the completed jobs section, there are pagination controls: '1 Pages. Jump to 1', 'Show 100 items in a page', and a 'Go' button.

Driver Memory Allocation



Executor Memory deep dive



To process 25GB data in spark

- How many CPU cores are required?
- How many executors are required?
- How much each executor memory is required?
- What is the total memory required?

How many executor CPU cores are required to process 25 GB data?

$$\begin{aligned} 25\text{GB} &= 25 \times 1024 \text{ MB} = 25600 \text{ MB} \\ \text{Number of Partitions} &= 25600 \text{ MB}/128 \text{ MB} = 200 \\ \text{Number of CPU Cores} &= \text{Number of Partitions} = 200 \end{aligned}$$

How many executors are required to process 25 GB data?

$$\begin{aligned} \text{Avg CPU cores for each executor} &= 4 \\ \text{Total number of executor} &= 200/4 = 50 \end{aligned}$$

Note:

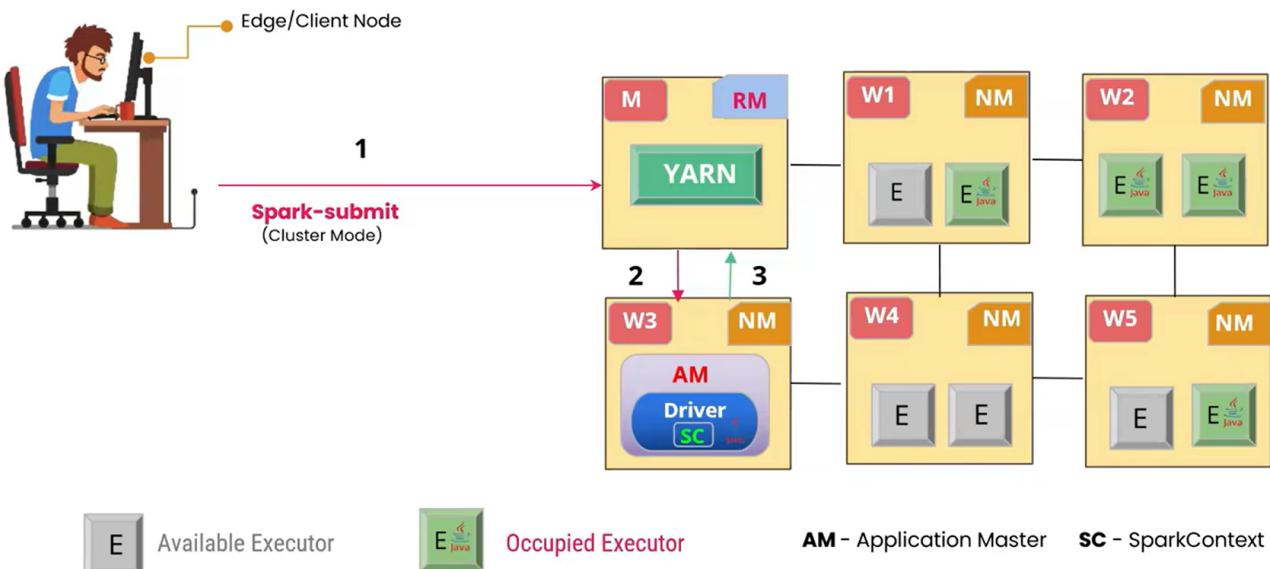
- By default, Spark creates one partition for each block of the file (blocks being 128MB by default in HDFS), but you can also ask for a higher number of partitions by passing a larger value.
- To get the better job performance in spark, researchers have found that we can take 2 to 5 maximum core for each executor
- Expected memory for each core = Minimum 4 x (default partition size)=4*128 MB=512 MB
- Executor memory is not less than 1.5 times of spark reserved memory (Single core executor memory should not be less than 450MB)

How much each executor memory is required to process 25 GB data?

$$\begin{aligned} \text{CPU cores for each executor} &= 4 \\ \text{Memory for each executor} &= 4 \times 512 \text{ MB} = 2 \text{ GB} \\ \text{Total Memory for all the executor} &= 50 \times 2 \text{ GB} = 100 \text{ GB} \end{aligned}$$

What is the total memory required to process 25 GB data?

$$\begin{aligned} \text{Total number of executor} &= 50 \\ \text{Memory for each executor} &= 2 \text{ GB} \\ \text{Total Memory for all the executor} &= 50 \times 2 \text{ GB} = 100 \text{ GB} \end{aligned}$$



In []: