# Apache Spark

## Overview

Apache spark is an open source lightning fast unified analytics engine for large scale processing and machine learning. It provides high-level APIs in Java, Scala, Python and R, and an optimized engine that supports general execution graphs.

## Quick Start

### Installation

- Download a docker image : docker pull bitnami/spark

- Once downloaded , run the image : docker run -d --name spark bitnami/spark

- Open the container in interactive mode: docker exec -it spark shell

- Create a text file using command: cd /tmp && echo "apple orrage guava mango pineapple peaches apple mango oragne orange \n apple orrage guava mango pineapple peaches apple mango oragne orange " > words.txt

- Run: `spark-shell` . It should open a scala based spark shell like below.
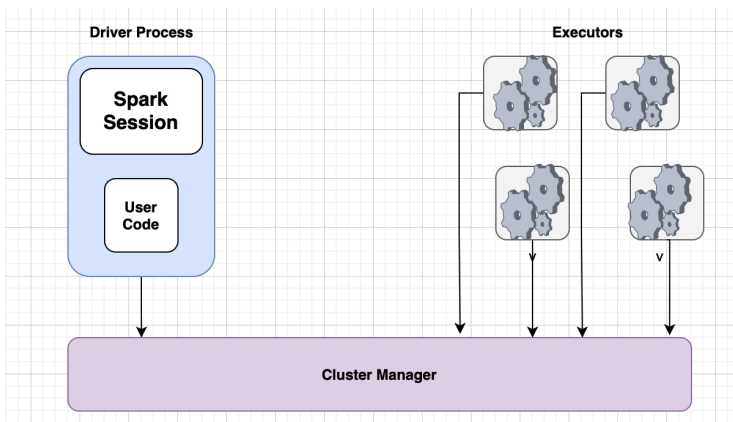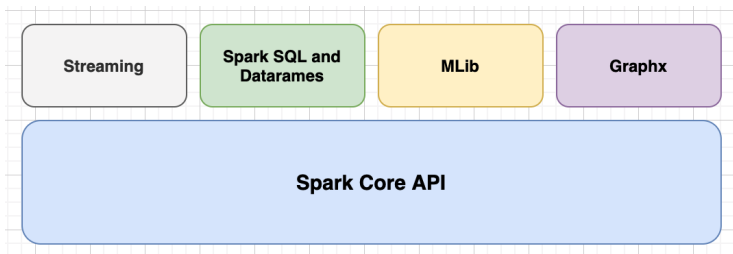


# Write a word-count Program

- Inside spark-shell Run this below code

```
val textFile = spark.read.textFile("/tmp/words.txt")
textFile: org.apache.spark.sql.Dataset[String] = [value: string]

val wordCounts = textFile.flatMap(line => line.split(" ")).groupByKey(identity).count()
wordCounts: org.apache.spark.sql.Dataset[(String, Long)] = [key: string, count(1): bigint

wordCounts.collect()
res0: Array[(String, Long)] = Array((orange,2), (apple,4), (mango,4), (orrage,2), (pineap
```

# Architecture

# Highlights of Spark

- Speed

- Ease of Use

- A Unified Engine

- Spark Web UI (localhost:4040)



# Core Concepts

- DAG & Lazy Evaluation

- Transformations and actions

- RDD & Dataframe

- SparkSession

# Reference

https://spark.apache.org/docs/latest/quick-start.html

https://www.oreilly.com/library/view/spark-the-definitive/9781491912201/

https://alvinalexander.com/scala/collection-scala-flatmap-examples-map-flatten/

https://spark.apache.org/docs/latest/configuration.html

https://spark.apache.org/docs/latest/api/java/index.html