

Системы анализа больших данных
(САБД)

What is Apache Spark?

Санкт-Петербургский политехнический университет Петра Великого
Институт компьютерных наук и технологий
Высшая школа программной инженерии

Ли Ицзя

2023.12



CONTENT

01

History of Spark

02

What is Spark?

03

Spark Features

04

Components of Apache Spark

05

Spark Architecture

06

Applications of Spark

History of Apache Spark

2009

Started as a project at UC Berkley
AMPLab



2013

Spark became an Apache top level
project



Now

Exists as a next generation real-
time and batch processing
framework



2010

Open sourced under a BSD license



2014

Used by Databricks to sort large-
scale datasets and set a new world
record



What is Apache Spark?



Apache Spark™ is an open-source data processing engine across clusters of computers using simple programming constructs.

Support various programming languages



Developers and data scientists incorporate Spark into their applications to rapidly query, analyze and transform data at scale



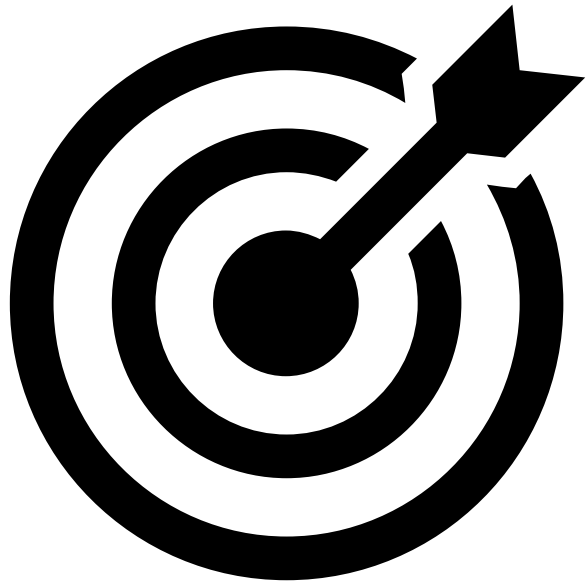
query



analyze



transform



Spark Features

Spark Features



Fast processing



Spark contains [Resilient Distributed Dataset \(RDD\)](#) which saves them time taken in reading, and writing operations and hence, it runs almost ten to hundred times faster than Hadoop

In-memory computing



In Spark, data is stored in the [RAM](#), so it can access the data quickly and accelerate the speed of analytics

Flexible



Spark supports [multiple languages](#) and allow the developers to write applications in Java, Scala, R or Python

Fault tolerance



Spark contains [Resilient Distributed Dataset \(RDD\)](#) that are designed to handle the failure of any worker node in the cluster. Thus, it ensures that the loss of data reduces to zero

Better analytics



Spark has a rich set of [SQL queries, machine learning algorithms, complex analytics](#), etc. With all these functionalities, analytics can be performed better

Hadoop vs Spark



Processing data using MapReduce in Hadoop is slow

Hadoop has more lines of code. Since it is written in Java, it takes more time to execute

Performs batch processing of data

Hadoop relies on data replication and checkpointing to ensure fault tolerance

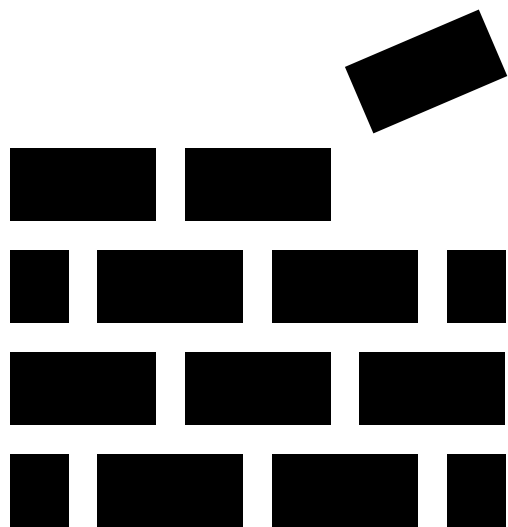


Spark processes data **100 times faster** than MapReduce as it is done in-memory

Spark **has fewer lines of code** as it is implemented in Scala

Performs both batch processing and **real-time processing** of data

Spark provides fault-tolerance through RDD, which **duplicates nothing**



Components of Spark

Components of Apache Spark



Spark Core



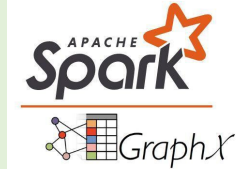
Spark SQL



Spark
Streaming



MLlib



GraphX



Components of Apache Spark –

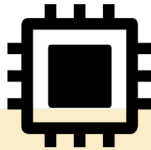
Spark Core



Spark Core

Spark core is the base engine for large-scale parallel and distributed data processing

It is responsible for:



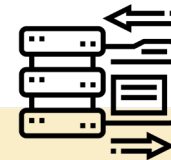
Memory management



Fault recovery



Scheduling, distributing and monitoring jobs on a cluster



Interacting with storage systems



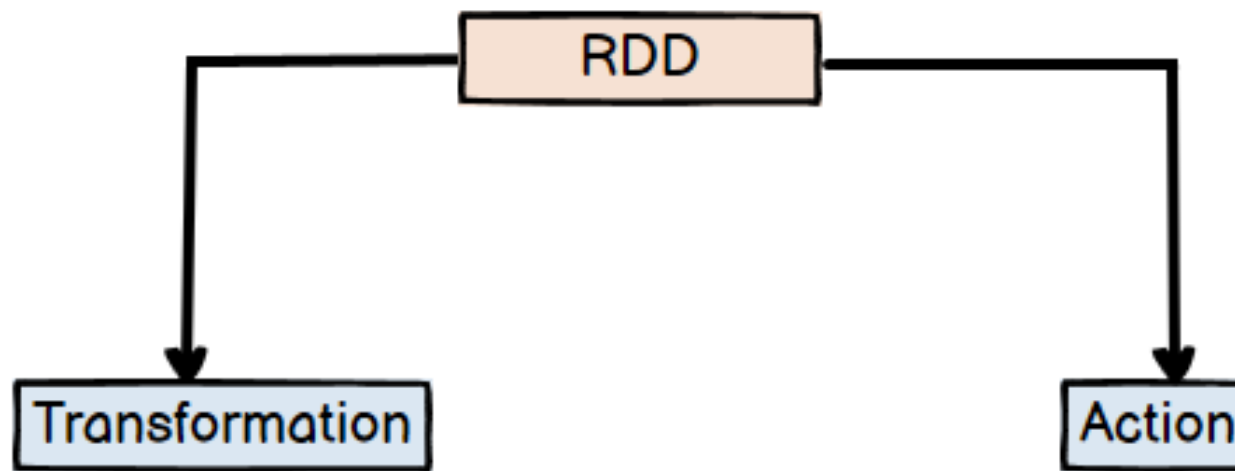
Spark Core

Resilient Distributed Dataset (RDD)

Spark core is embedded with RDDs (Resilient Distributed Dataset), an immutable fault-tolerant, distributed collection of objects that can be operated on in parallel



Spark Core



These are operations (such as *map, filter, union*) that are performed on an RDD that **yields a new RDD** containing the result

These are operations (such as *reduce, first, count*) that **return a value** after running a computation on an RDD

Components of Apache Spark –

Spark SQL



Spark SQL

Spark SQL framework component is used for structured and semi-structured data processing

Spark SQL

Spark SQL

Spark SQL Architecture

DataFrame DSL

Spark SQL & HSQL

DataFrame API

Data Source API

CSV

JSON

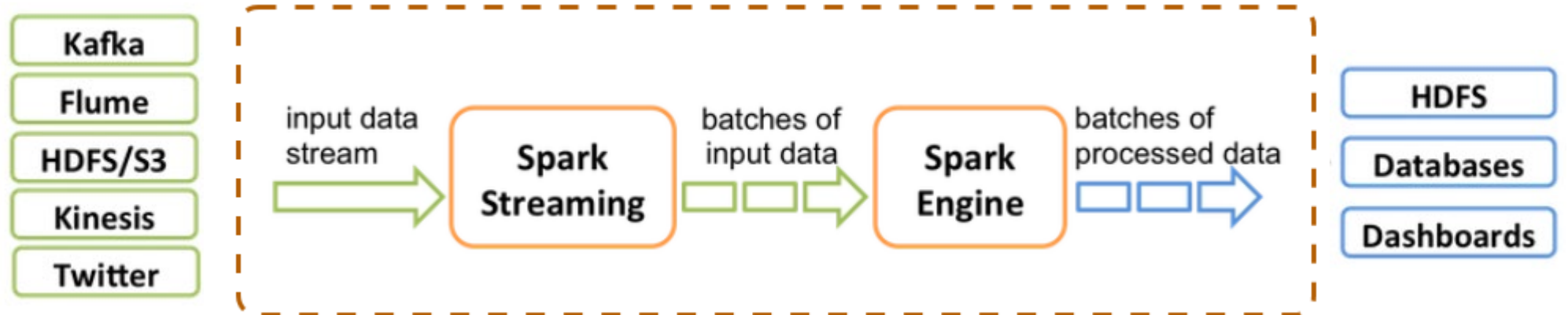
JDBC

Spark Streaming

Spark Streaming is a lightweight API that allows developers to perform batch processing and real-time streaming of data with ease

Provides secure, reliable, and fast processing of live data streams

Spark
Streaming
Spark
Streaming



Spark MLlib

Spark MLlib is low-level machine learning library that is simple to use, is scalable, and compative with various programming languages

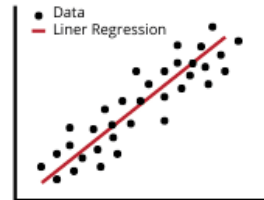


It contains machine learning libraries that have an implementation of various machine learning algorithms

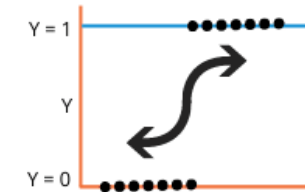


MLlib

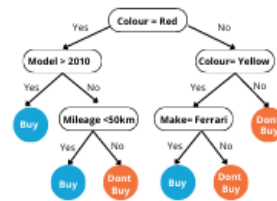
Linear Regression



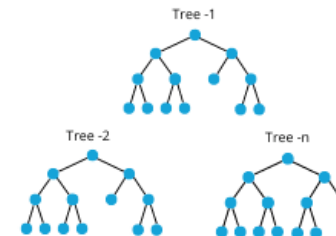
Logistic Regression



Decision Trees



Random Forest

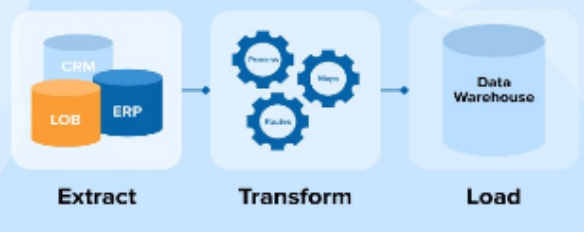


GraphX

GraphX is Spark's own Graph Computation Engine and data store

Provides a uniform tool for ETL

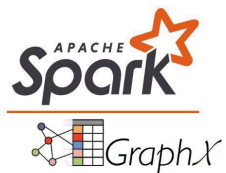
ETL Process (Extract, Transform, Load)



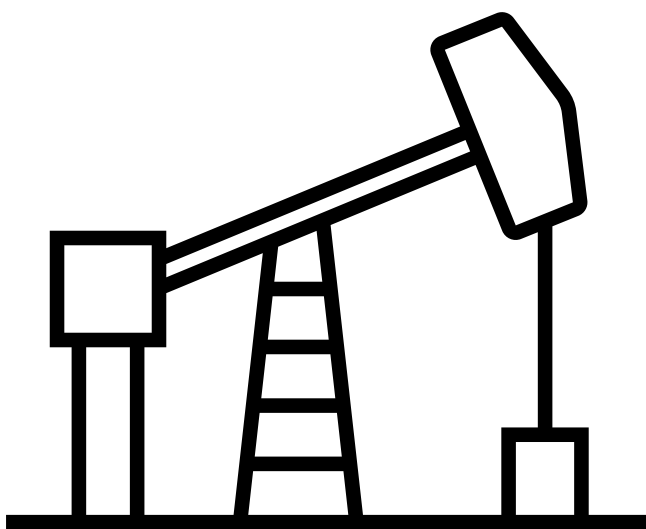
Exploratory data analysis



Interactive graph computations



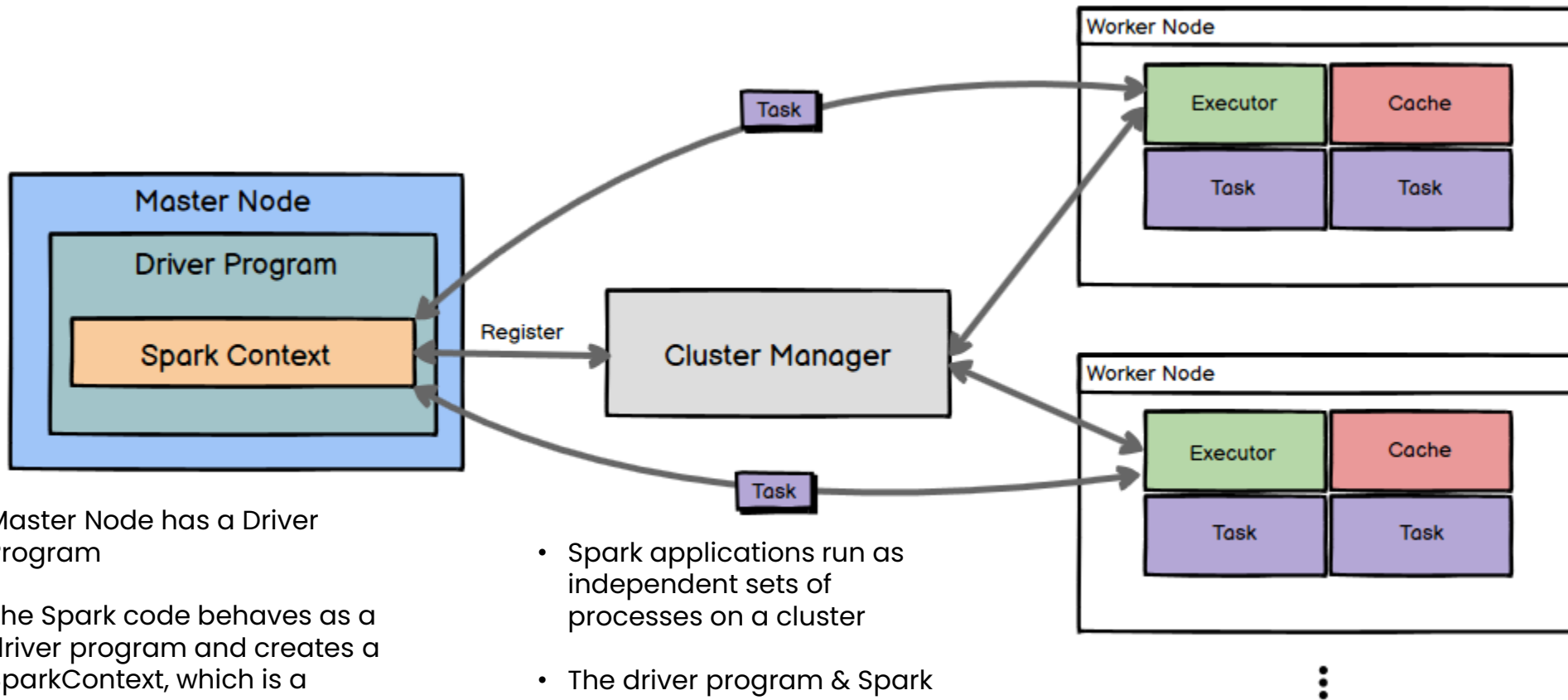
GraphX



Spark Architecture

Spark Architecture

Apache Spark uses a master-slave architecture that consists of a driver, that runs on a master node, and multiple executors which run across the worker nodes in the cluster



- Master Node has a Driver Program
- The Spark code behaves as a driver program and creates a SparkContext, which is a gateway to all the Spark functionalities
- Spark applications run as independent sets of processes on a cluster
- The driver program & Spark context takes care of the job execution within the cluster

Spark Cluster Managers



Standalone mode

1

By default, applications submitted to the cluster will run in FIFO order, and each application will try to use all available nodes.



2

Apache Mesos is an open-source project to manage computer clusters, and can also run Hadoop applications



3

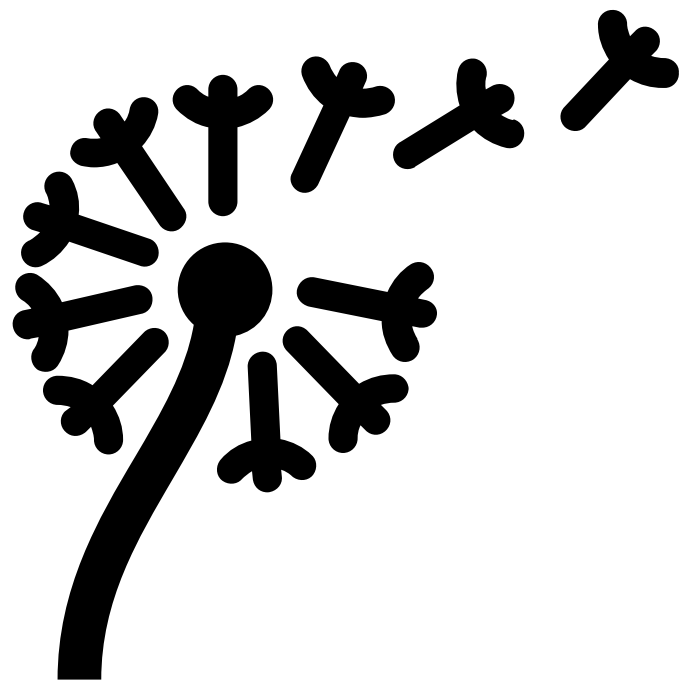
Apache YARN is the cluster resource manager of Hadoop 2. Spark can be run on YARN



kubernetes

4

Kubernetes is an open-source system for automating deployment, scaling, and management of containerized applications



Applications of Spark

Applications of Spark



JPMorgan uses Spark to detect fraudulent transactions, analyze the business spends of an individual to suggest offers, and identify patterns to decide how much to invest and where to invest

JPMORGAN
CHASE & CO.

Banking



Alibaba uses Spark to analyze large sets of data such as real-time transaction details, browsing history, etc. in the form of Spark jobs and provides recommendations to its users.


Alibaba.com™

E-commerce



Entertainment and gaming companies like Netflix and Riot games use Apache Spark to showcase relevant advertisements to their users based on the videos that they watch, share and like

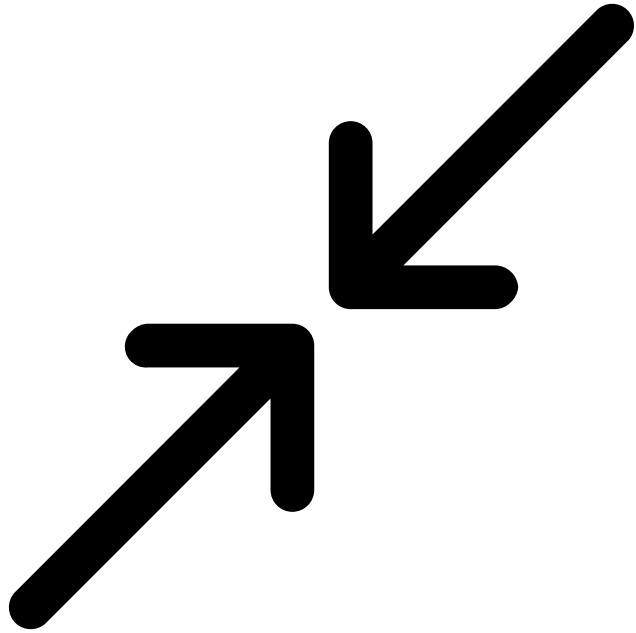


 **RIOT
GAMES™**

Entertainment

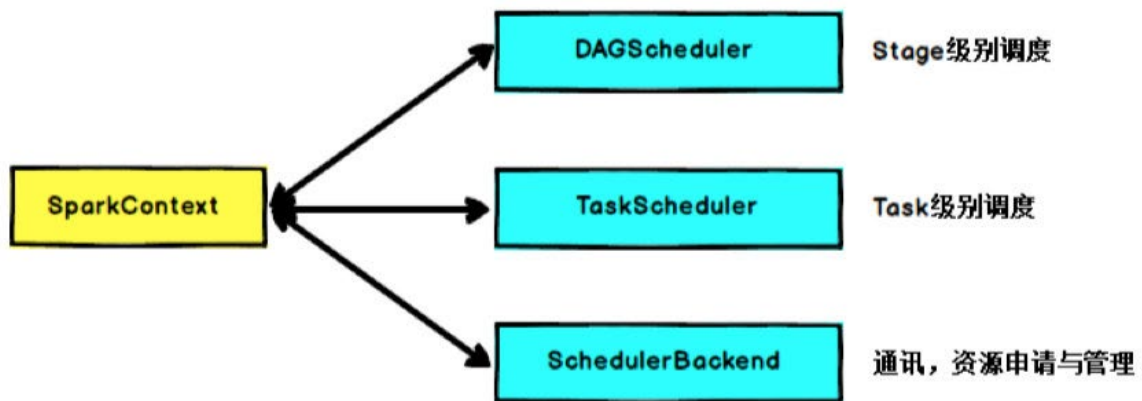


**THANK YOU
for your precious
attention :)**

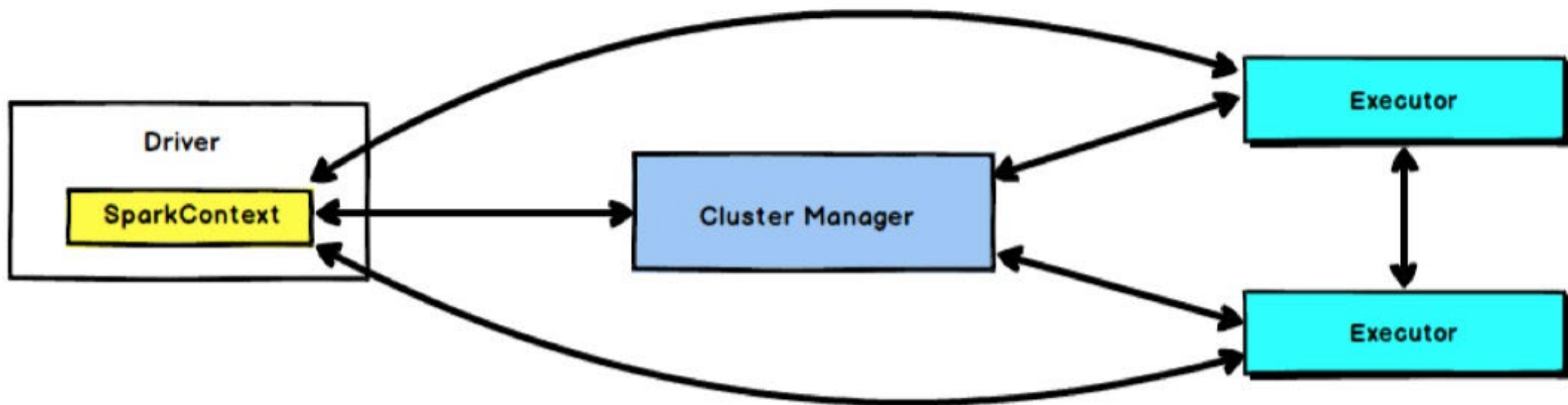


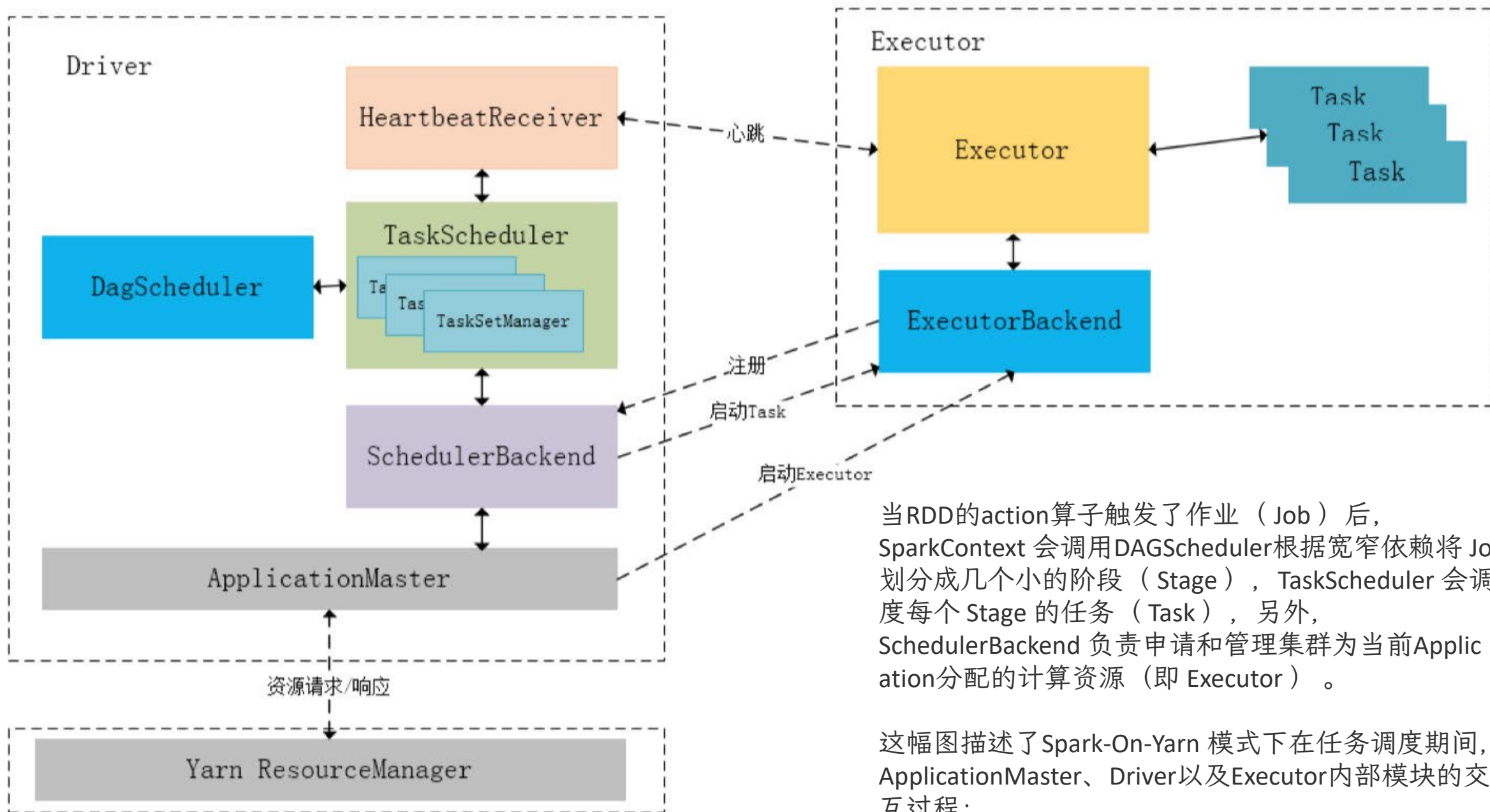
Hadoop vs Spark

SparkContext



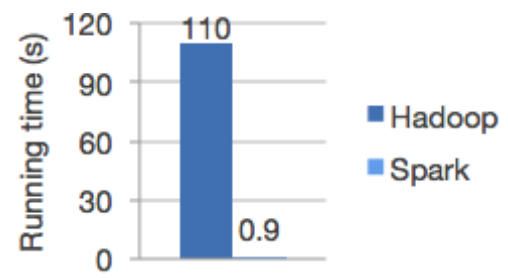
SparkContext的核心作用是初始化 Spark 应用程序运行所需的核心组件，包括高层调度器（DAGScheduler）、底层调度器（TaskScheduler）和调度器的通信终端（SchedulerBackend），同时还会负责Spark程序向ClusterManager的注册等。





当RDD的action算子触发了作业（Job）后，SparkContext 会调用DAGScheduler根据宽窄依赖将 Job 划分成几个小的阶段（Stage），TaskScheduler 会调度每个 Stage 的任务（Task），另外，SchedulerBackend 负责申请和管理集群为当前Application分配的计算资源（即 Executor）。

这幅图描述了Spark-On-Yarn 模式下在任务调度期间，ApplicationMaster、Driver以及Executor内部模块的交互过程：



MapReduce vs Spark (WordCount)

```
1 public class WordCount {
2     public static class TokenizerMapper extends Mapper<Object, Text, Text, IntWritable> {
3         private final static IntWritable one = new IntWritable(1);
4         private Text word = new Text();
5
6         public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
7             StringTokenizer itr = new StringTokenizer(value.toString());
8             while (itr.hasMoreTokens()) {
9                 word.set(itr.nextToken());
10                context.write(word, one);
11            }
12        }
13    }
14
15    public static class IntSumReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
16        private IntWritable result = new IntWritable();
17        public void reduce(Text key, Iterable<IntWritable> values,
18            Context context) throws IOException, InterruptedException {
19            int sum = 0;
20            for (IntWritable val : values) {
21                sum += val.get();
22            }
23            result.set(sum);
24            context.write(key, result);
25        }
26    }
27 }
```

```
1 val textFile = sc.textFile("hdfs://...")
2 val counts = textFile.flatMap(line => line.split(" "))
3     .map(word => (word, 1))
4     .reduceByKey(_ + _)
5 counts.saveAsTextFile("hdfs://...")
```

Spark Use Case

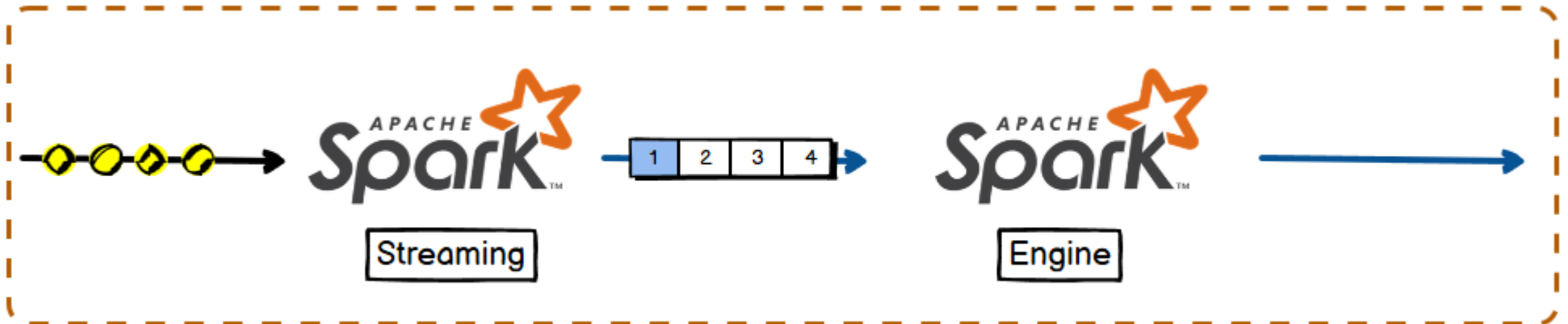
Spark Streaming

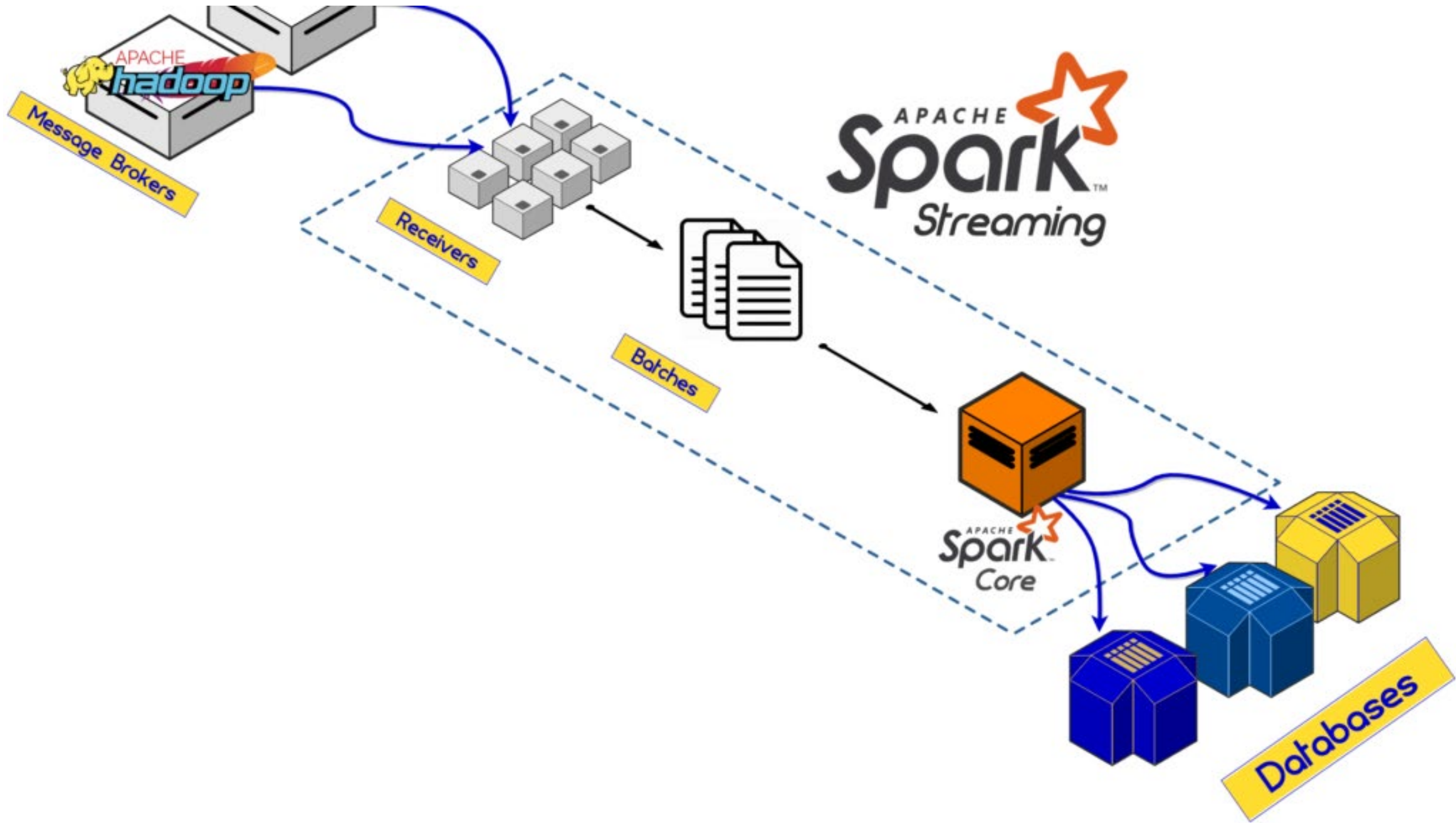
Spark Streaming is a lightweight API that allows developers to perform batch processing and real-time streaming of data with ease

Provides secure, reliable, and fast processing of live data streams

Spark
Streaming

Spark
Streaming





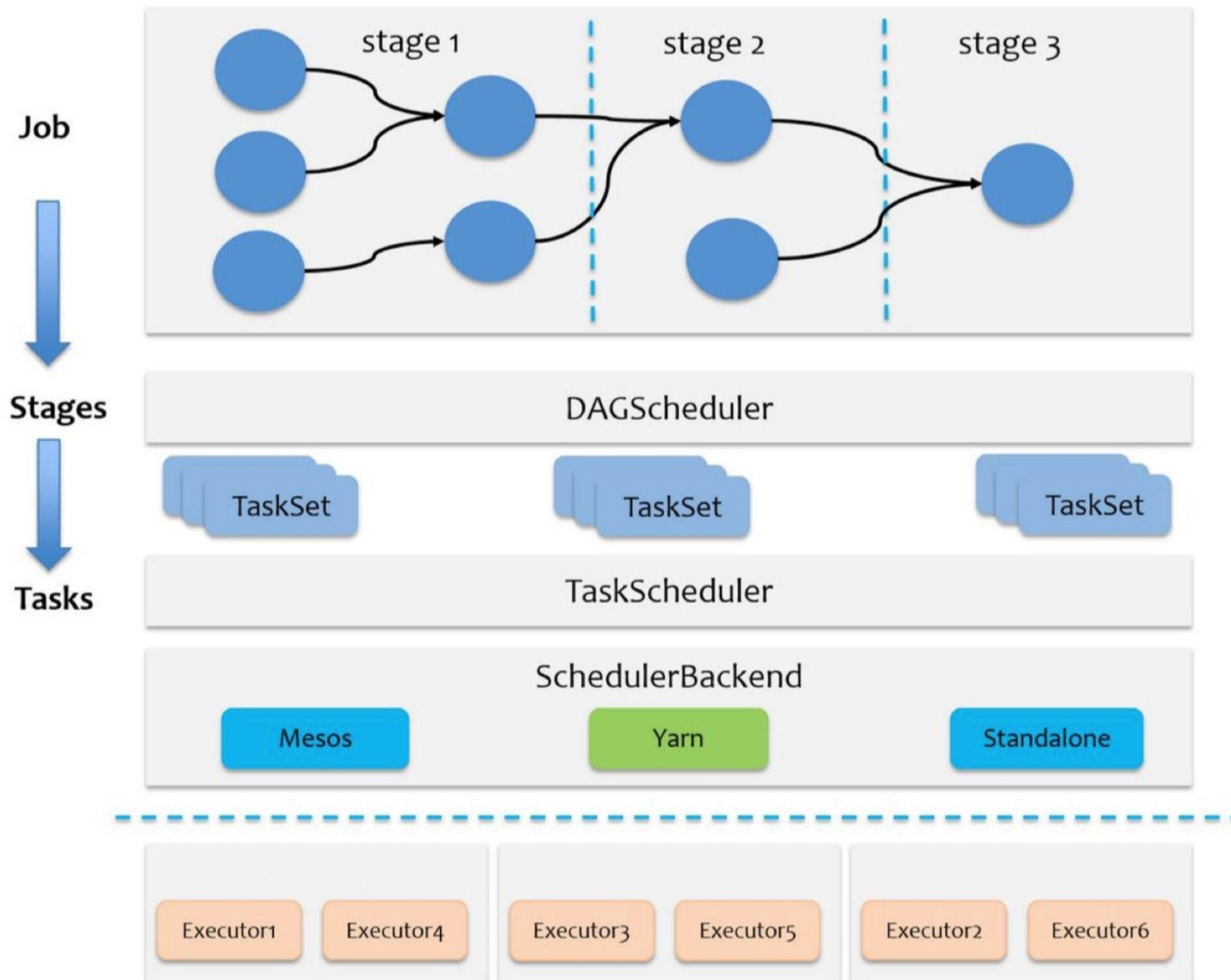
一个Spark程序包括Job、Stage以及Task三个概念：

Job是以Action方法为界，遇到一个Action方法则触发一个Job；

Stage是Job的子集，以RDD宽依赖（即Shuffle）为界，遇到Shuffle做一次划分；

Task是Stage的子集，以并行度（分区数）来衡量，分区数是多少，则有多少个task。

Spark任务的调度总体上分两路进行，一路是Stage级的调度，一路是Task级的调度，总体的调度流程如下：



提交流程

