

Batch Layer with Apache Spark



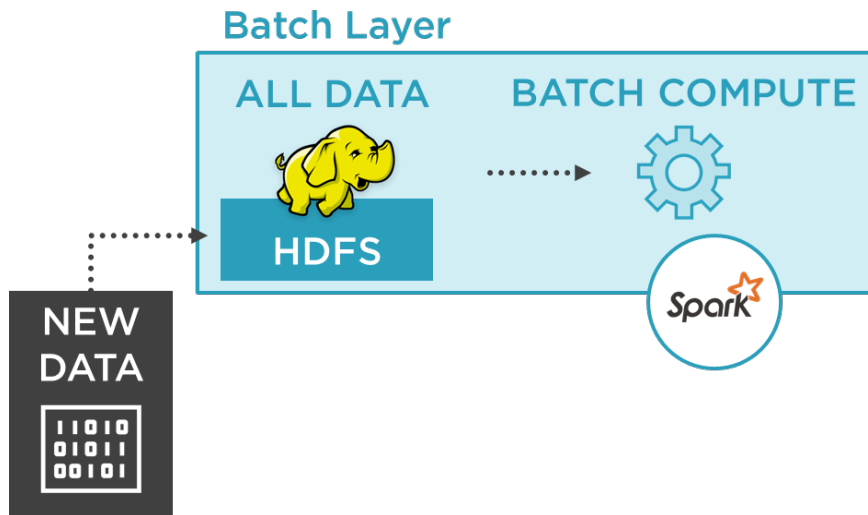
Ahmad Alkilani

DATA ARCHITECT

@akizl



Building the Batch Layer



Fundamentals of batch layer

Aggregations in Spark

- RDD API
- DataFrame and DataSet API
- Caching

Introduction to Spark

Spark components and scheduling



Spark is a general purpose cluster computing platform designed with components for scheduling and executing against large datasets

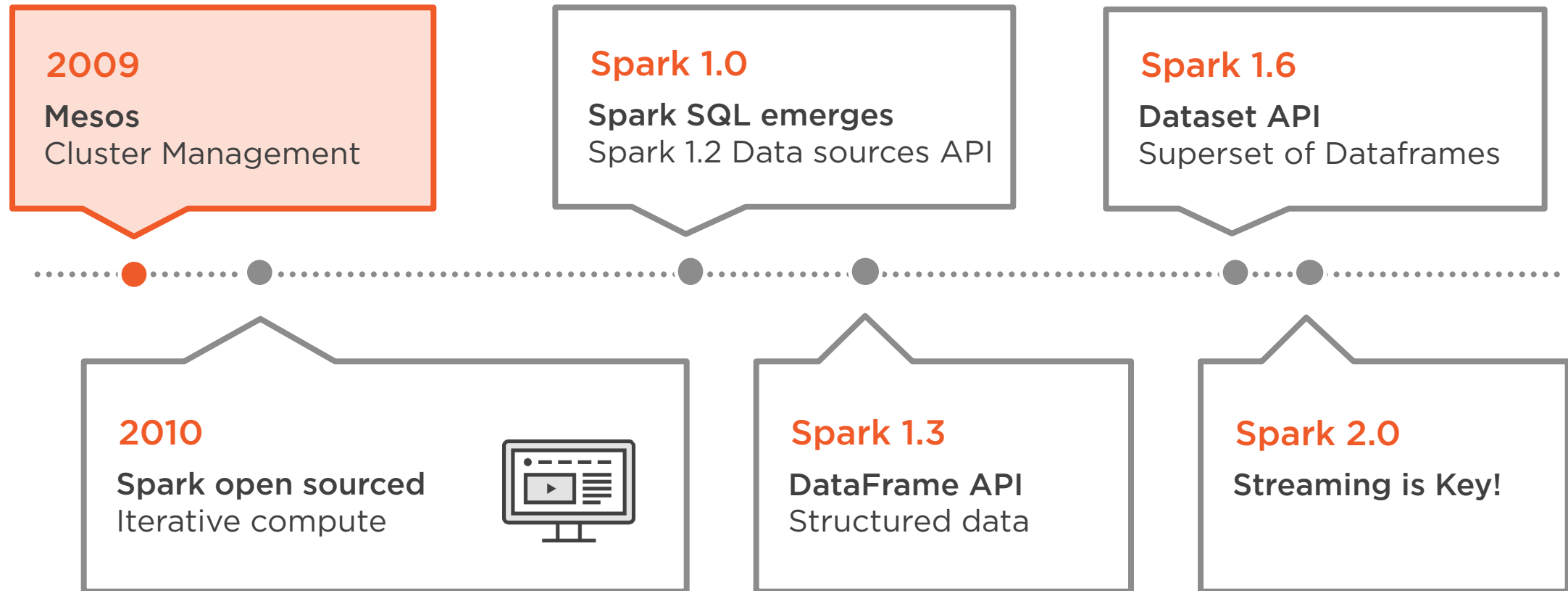




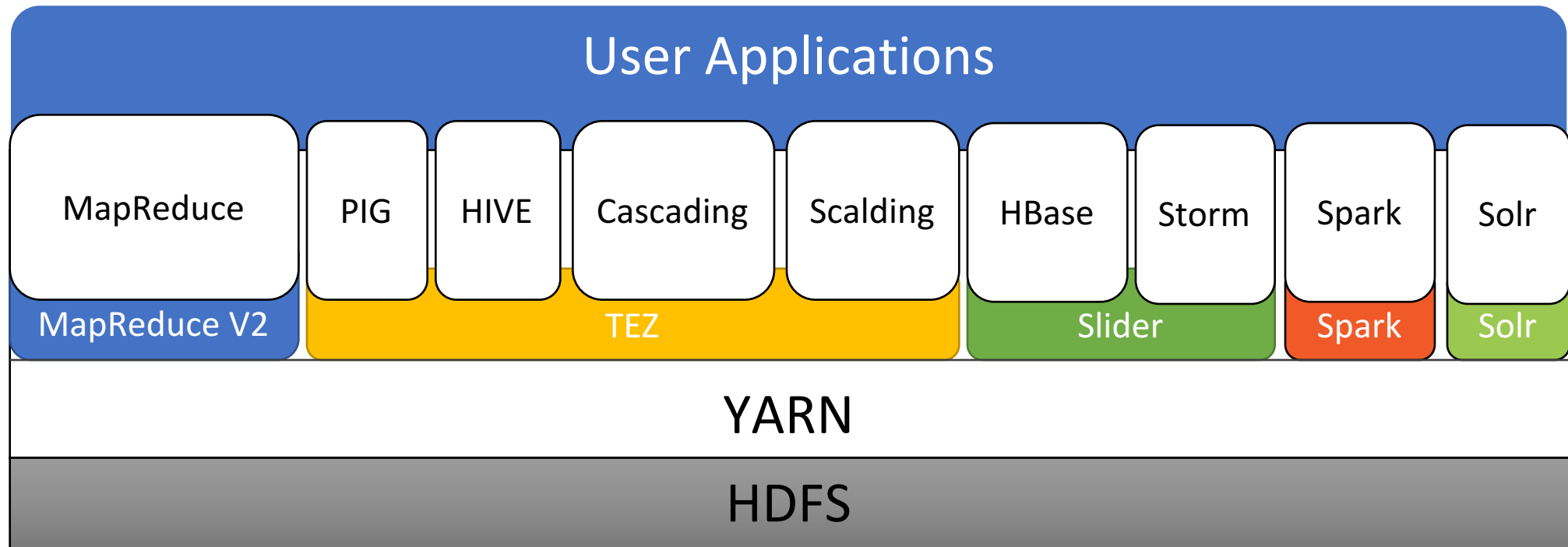
Batch workloads, Iterative algorithms, interactive queries, streaming applications, machine learning workloads, graph processing



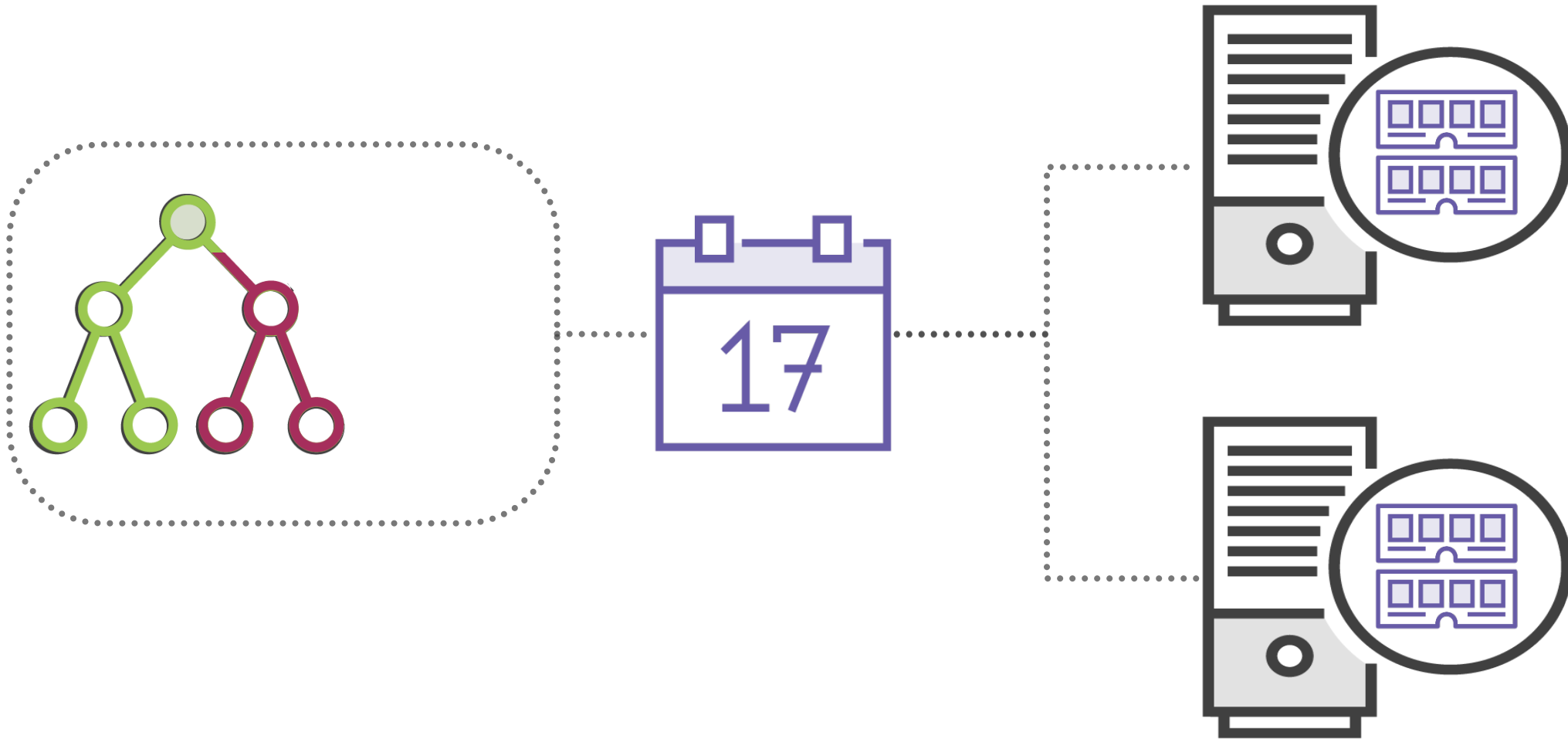
Timeline of Events



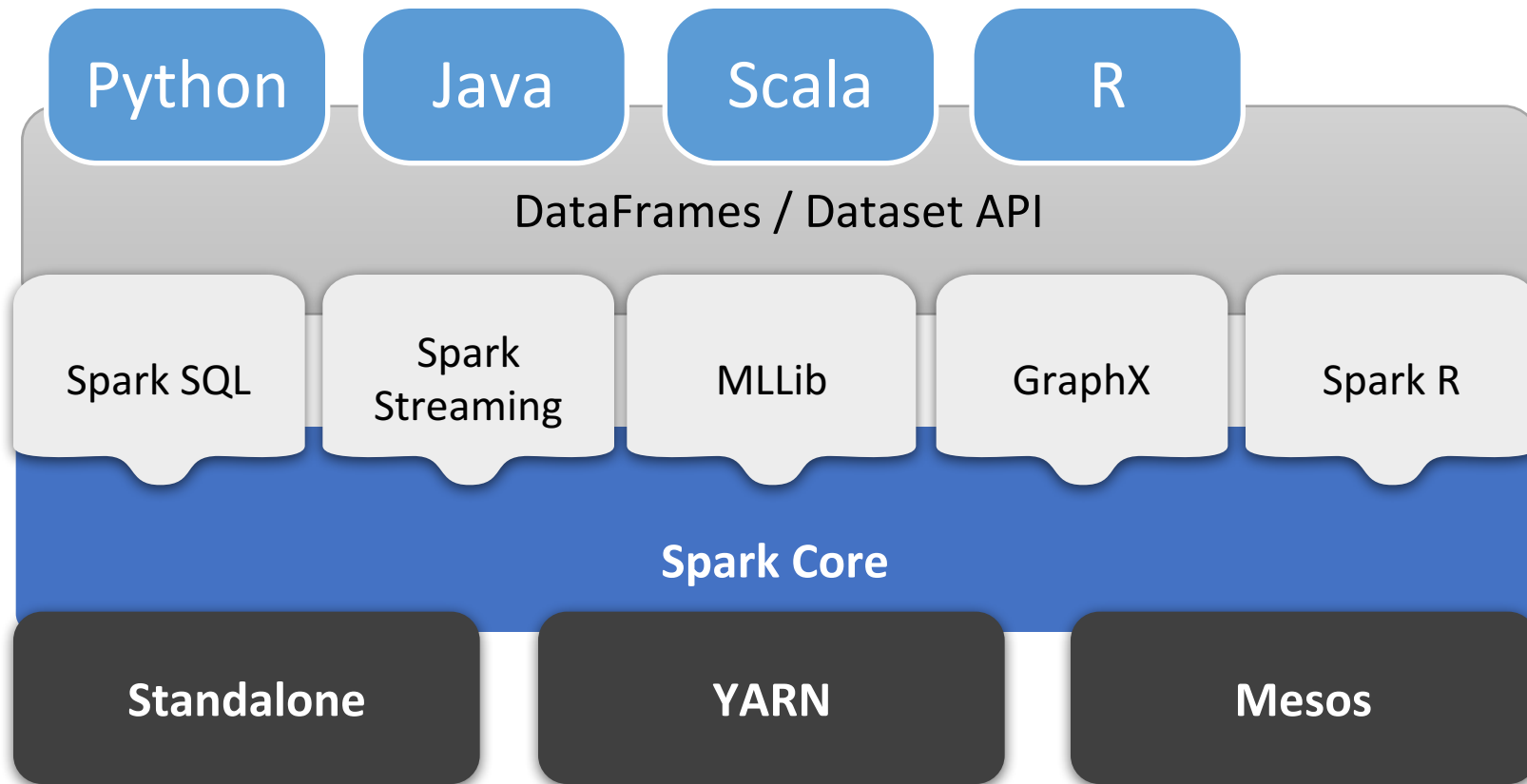
YARN



Spark Fundamentals



Spark Components



Spark RDD

Fundamental abstraction and building block. RDD represents a Resilient Distributed Dataset



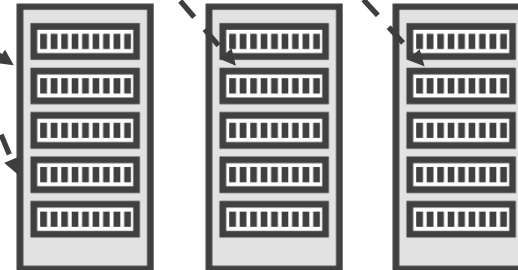
Spark RDD

```
val scalaCollection : Array[  ] = (     )
```

```
scalaCollection.map(item => turnOrange(item))
```

```
val sparkRDD      : RDD [  ] = (     )
```

```
sparkRDD.map(item => turnOrange(item))
```





Generates execution DAG

Optimizes DAG execution

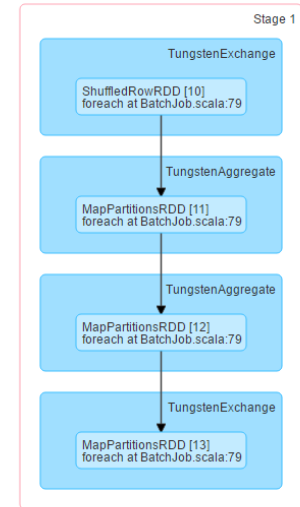
Cache datasets

Resiliency by lineage

Spark modules

- Core
- Streaming
- Spark SQL
- MLlib
- Spark R

▼ DAG Visualization



Spark Components & Scheduling



Spark's Relationship to Hadoop

HDFS

Read/Write

Native integration with
Hadoop APIs and
various file formats

Hive

Tight integration with
Hive and Hive's
metastore. Spark's
HiveContext uses
HiveQL

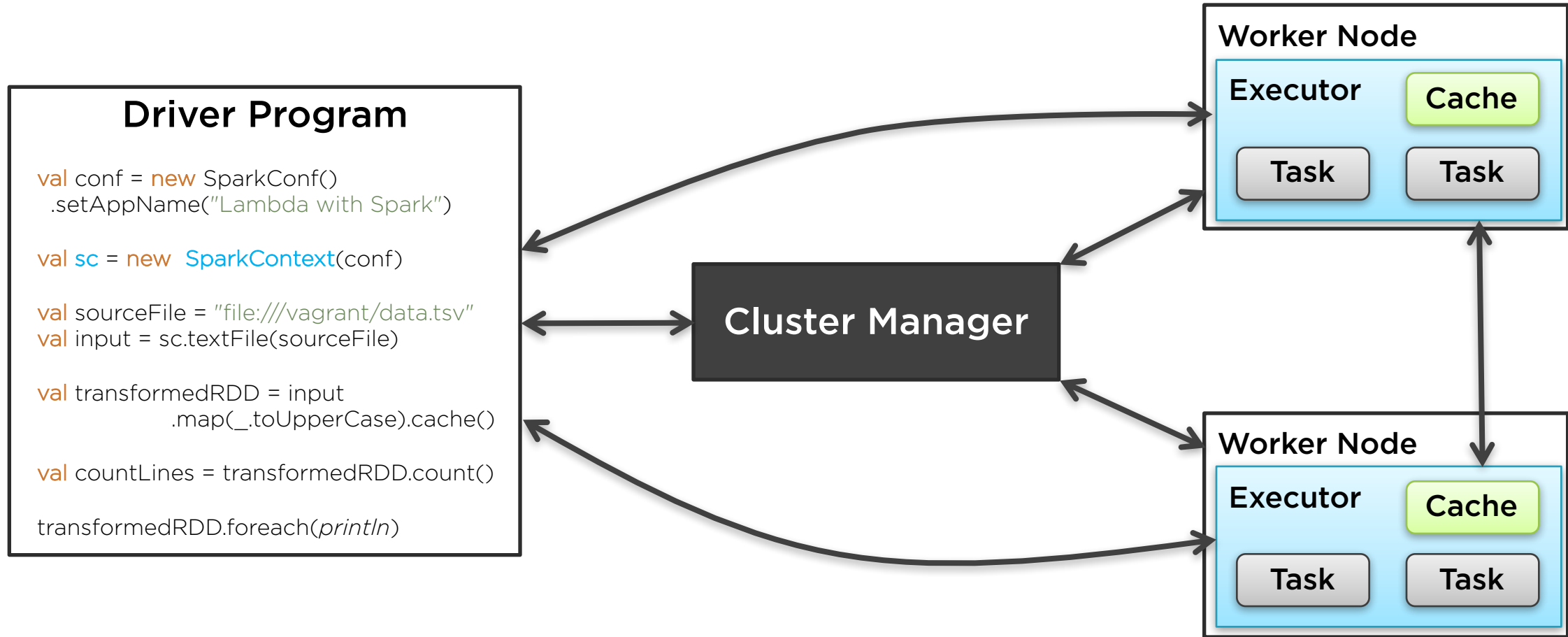
YARN

Cluster Manager

Uses YARN/Hadoop as
a cluster manager to
schedule the execution
of Spark Executors

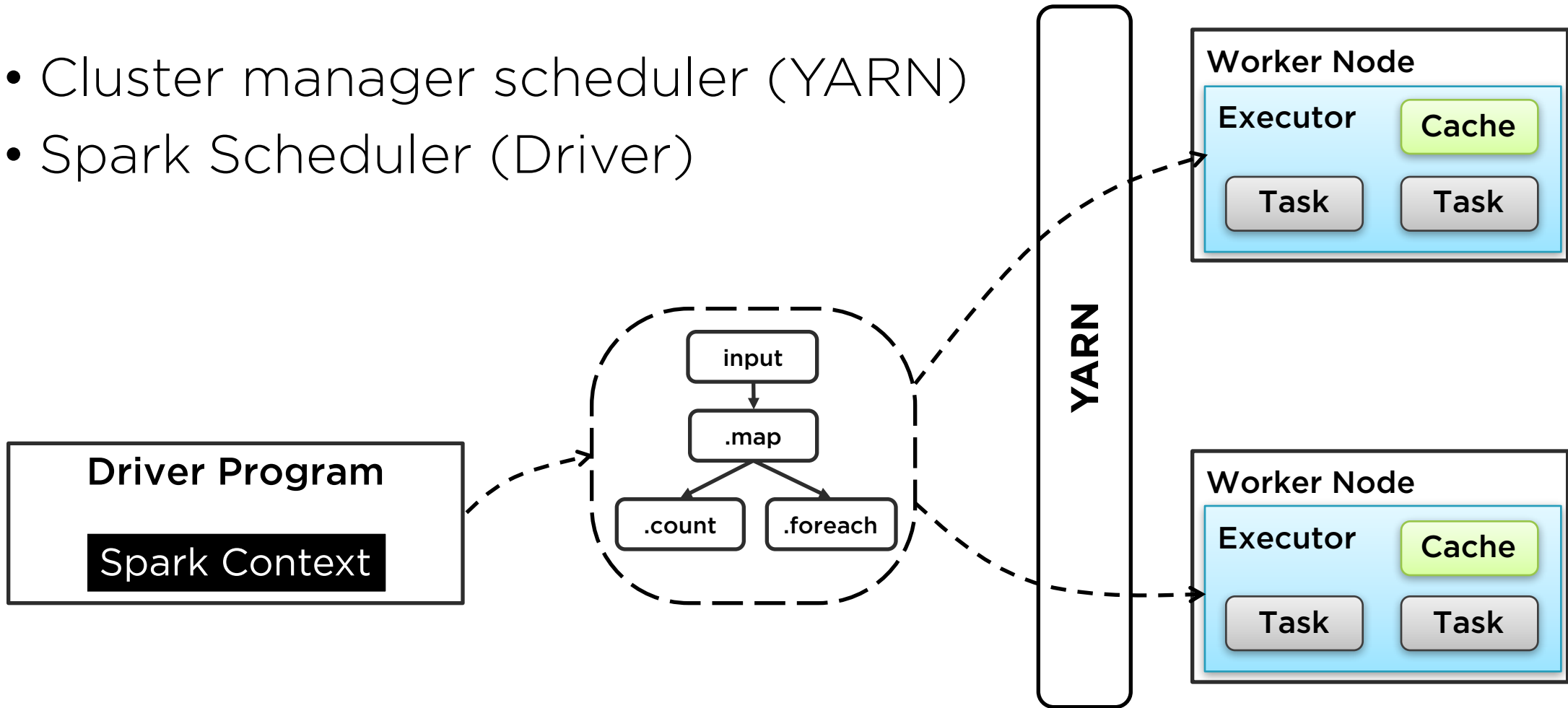


Spark Execution Components



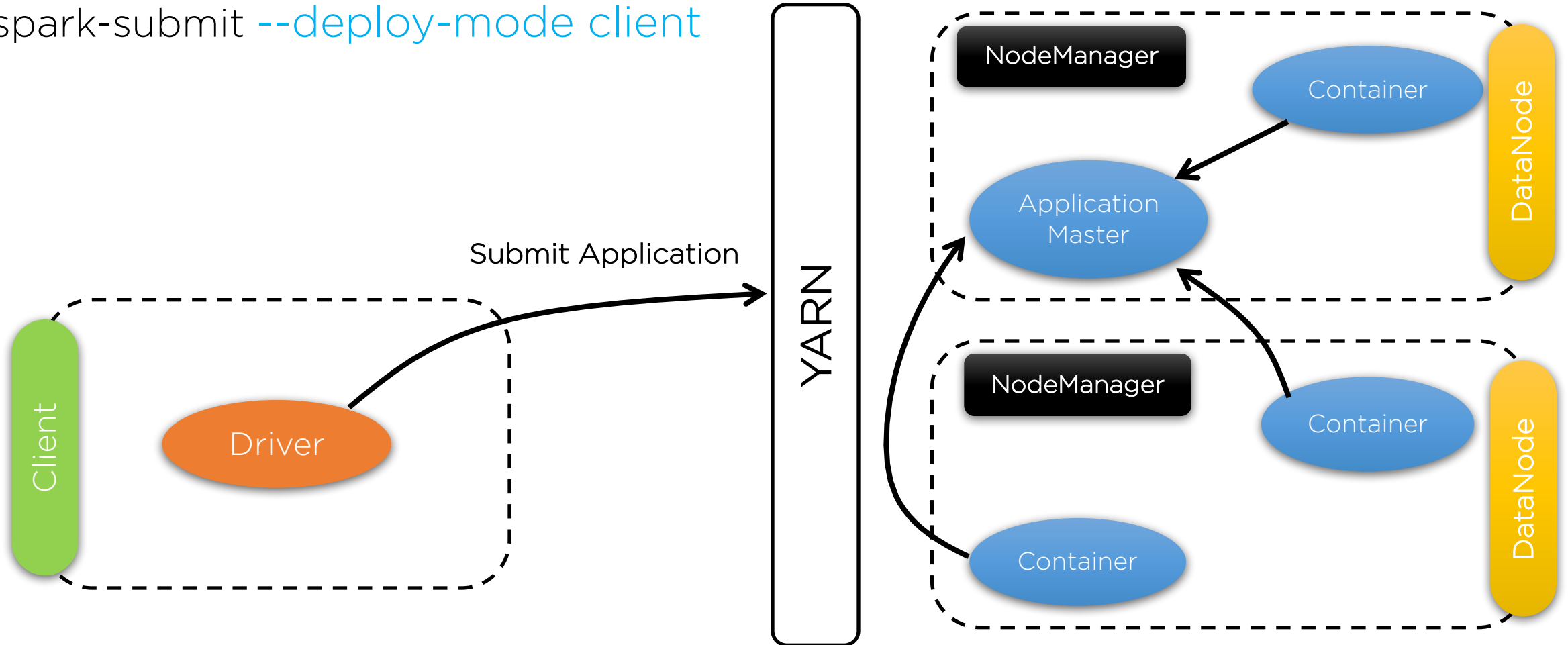
Scheduling

- Cluster manager scheduler (YARN)
- Spark Scheduler (Driver)



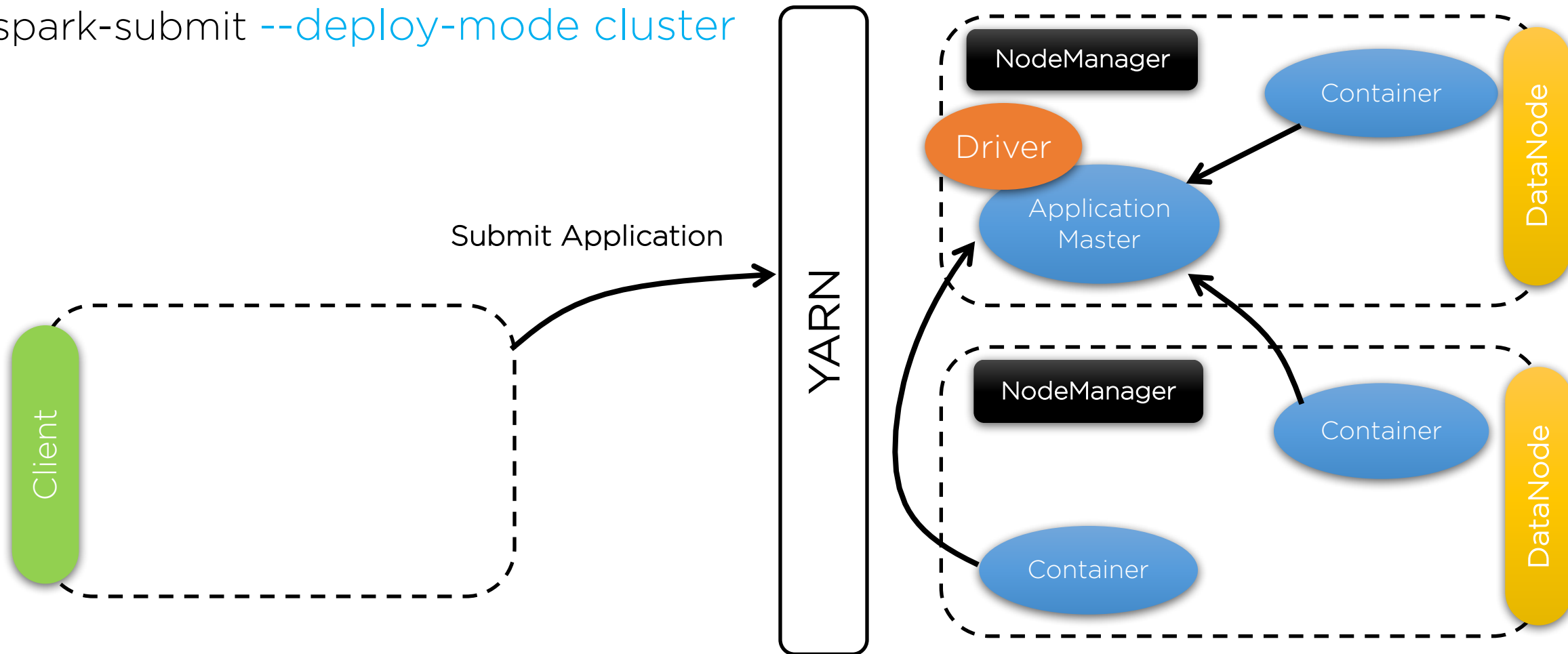
Spark Scheduling - Client

`./spark-submit --deploy-mode client`



Spark Scheduling - Cluster

`./spark-submit --deploy-mode cluster`



Client vs. Cluster Mode

Client mode

- Driver managed by client host
- Availability dependent on client
- Interactive (Spark Shell/REPL)

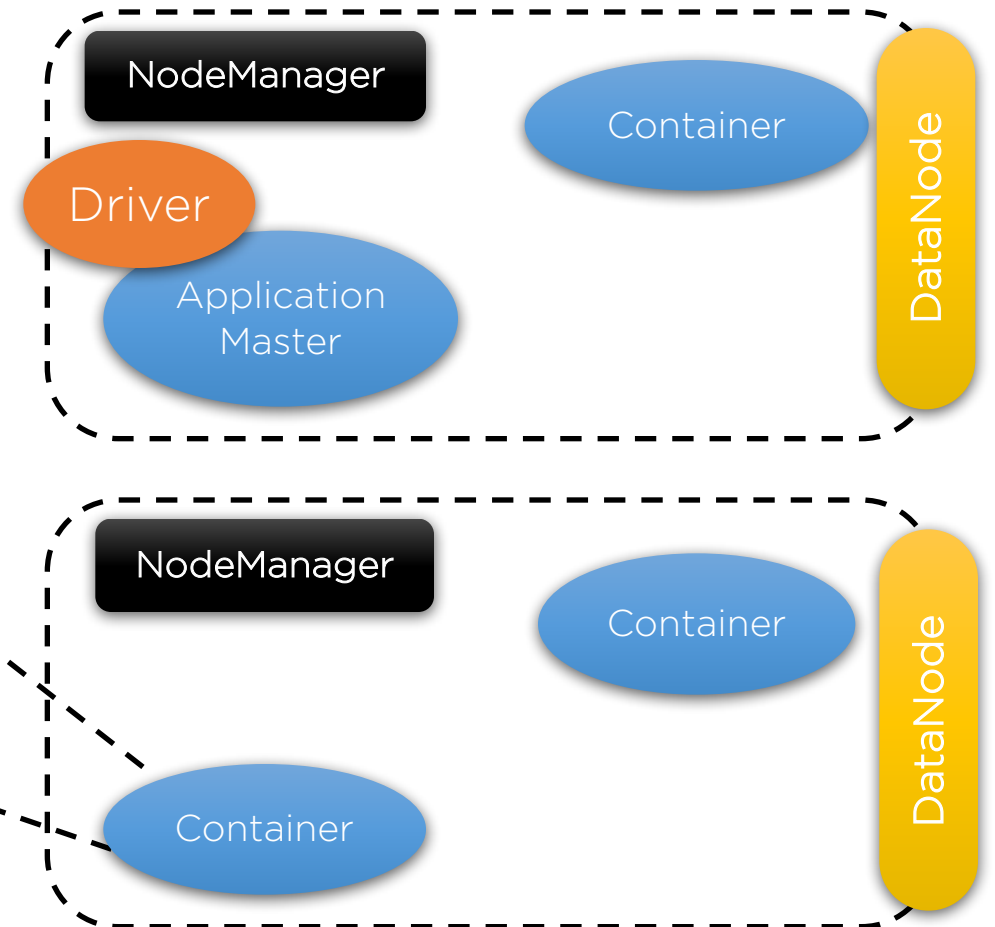
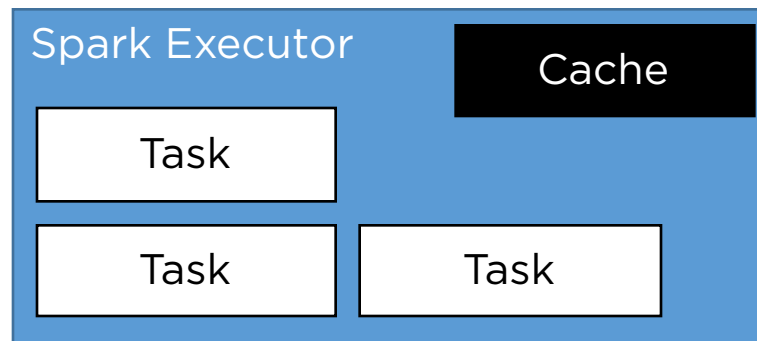
Cluster mode

- Driver managed by cluster manager
- Better availability guarantees
- Non-interactive

Inside the Container

Task

Performs actions on partitions (HDFS InputSplit)



Actions and Transformations

```
val inputRDD = sc
    .textFile("hdfs://...")
```

```
val wordsRDD = inputRDD
    .flatMap(_.split(" "))
    .map(word => (word, 1))
```

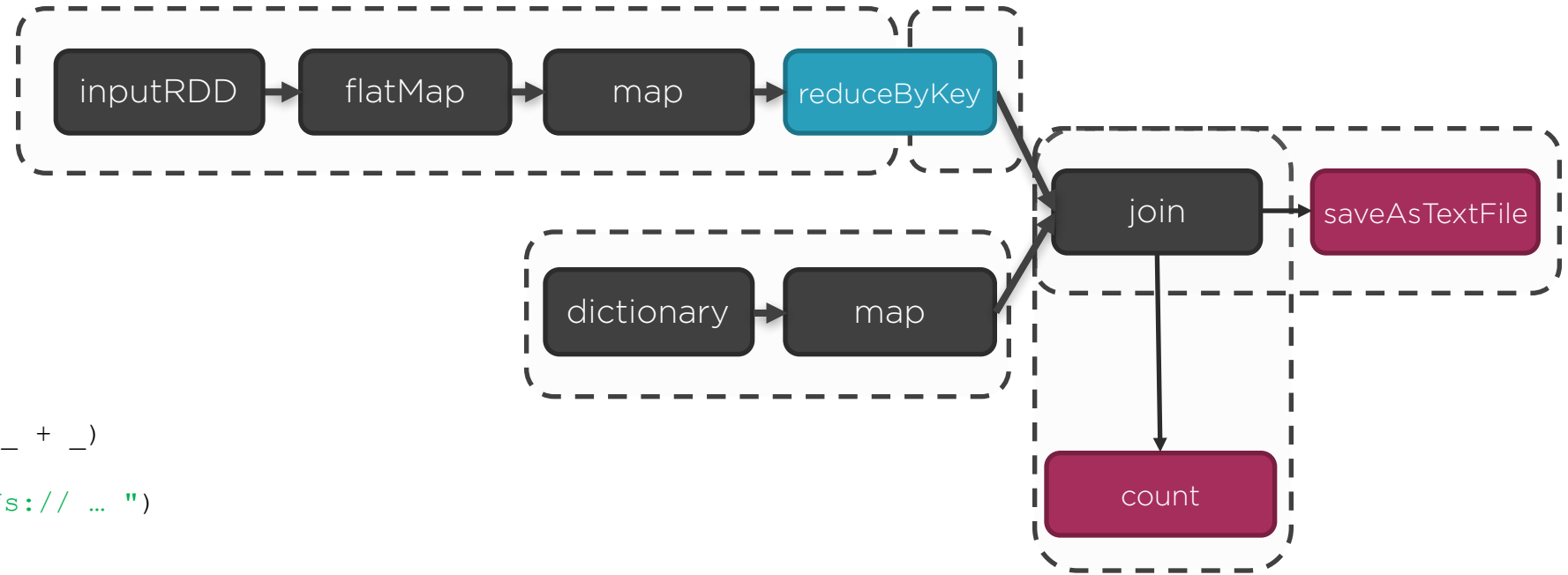
```
val wordCountsRDD = wordsRDD
    .reduceByKey(_ + _)
```

```
val dictionary = sc.textFile("hdfs:// ... ")
    .map { s =>
        val split = s.split(",")
        (split(0), split(1))
    }
```

```
val joined = wordCountsRDD.join(dictionary)
```

```
joined.cache()
```

```
joined.saveAsTextFile("hdfs:// ... ")
joined.count()
```



Demo



Demo



Log Producer – Clickstream Dataset

Typesafe Config - HOCON

- Layers of overrides
- Easier deployments

Stages of Log Producer

- Single file output
- Multiple files simulating stream
- Output to Kafka
- Kafka in Avro format



Log Producer

adjustedTimestamp	referrer	action	prevPage	visitor	page	product
1464253347272	Facebook	add_to_cart		Visitor-548011	Page-6	Chemin Du Soleil,Chardonnay
1464253347272	Other	page_view		Visitor-18032	Page-6	Purell,Advanced Hand Sanitizer -travel size
1464253347272	Twitter	page_view		Visitor-69498	Page-1	David's Tea,Chocolate Chili Chai Black Tea
1464253347272	Google	page_view		Visitor-697731	Page-3	Barilla,Rotini
1464253347272	Twitter	page_view		Visitor-992221	Page-0	Alessi,Garlic Breadsticks
1464253347272	Bing	page_view		Visitor-289265	Page-3	Garnier Fructis Style,Pure Clean Finishing Paste
1464253479272	Facebook	page_view		Visitor-313453	Page-0	Glass Plus,Glass Cleaner
1464253479272	Yahoo	page_view		Visitor-740395	Page-14	Clover Stornetta,Vitamin D Milk
1464253479272	Facebook	page_view		Visitor-576247	Page-9	Andersen's,Creamy Soup Split Pea
1464253479272	Other	purchase		Visitor-957650	Page-8	Dust Destroyer,Compressed-Gas Duster
1464253479272	Facebook	page_view		Visitor-141141	Page-11	Nestle,Leche Condensada
1464254239592	Google	page_view		Visitor-387911	Page-3	Caliber,Security Envelopes
1464254239592	Twitter	page_view		Visitor-410847	Page-12	Glad,Odor Shield Tall Kitchen Drawstring Bags
1464254239592	Yahoo	page_view		Visitor-684141	Page-2	Knorr,Salsa Lista Pizza
1464254239592	Google	page_view		Visitor-930963	Page-0	Gatorade,Fierce Grape
1464254239592	Twitter	page_view		Visitor-547631	Page-3	Tovolo,King Cube Tray

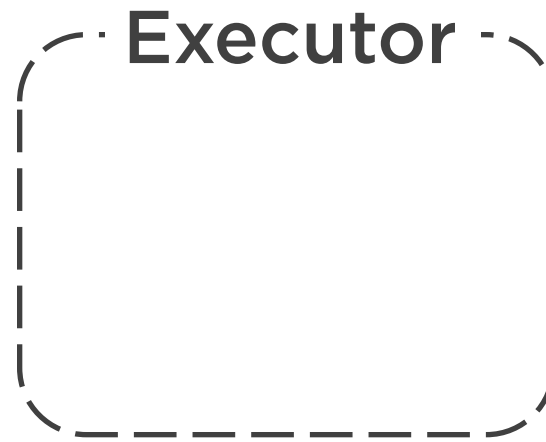
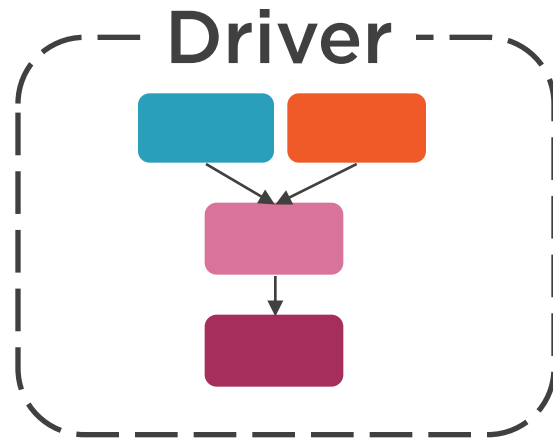
referrers



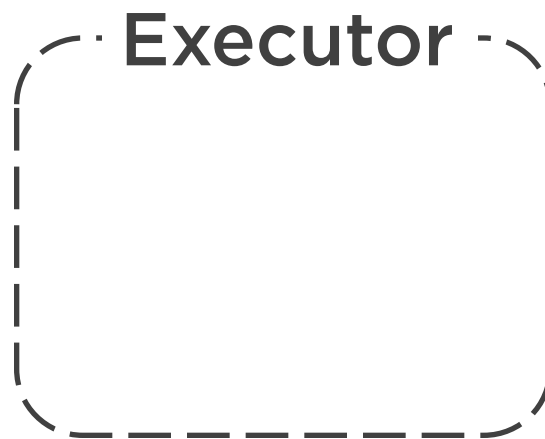
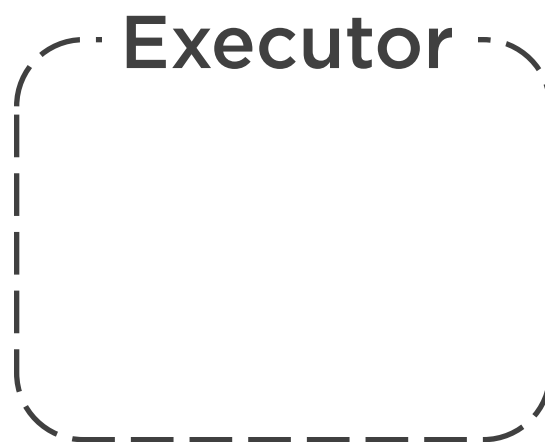
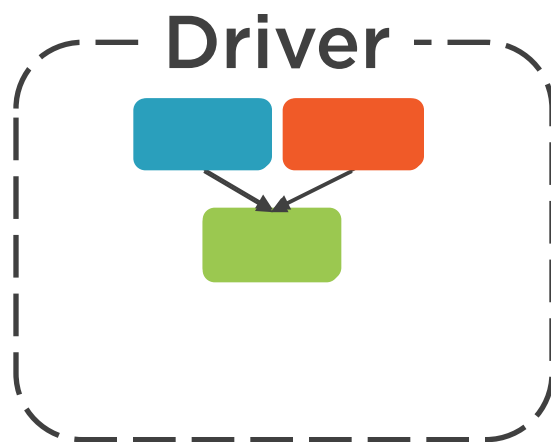
products



Spark Laziness



Spark Laziness



Summary

