# Recap

- Sqoop

- HBase

# Agenda for today

# Hadoop

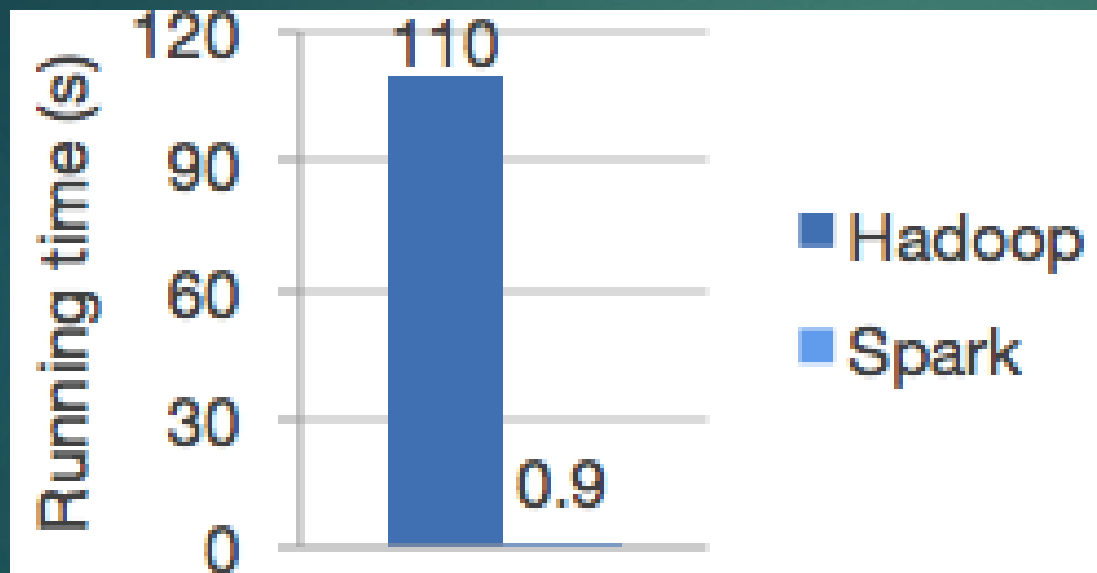| Processing nature | Tool |
| --- | --- |
| Batch processing on structured data | Hive |
| Ad-hoc analysis on structured data | Impala |
| Machine Learning | Apache Mahout |
| Graph processing | Apache Giraph |
| Stream processing | Apache Storm |

# Spark

| Processing nature | Spark |
| --- | --- |
| Batch processing on structured data | SparkSQL |
| Ad-hoc analysis on structured data | SparkSQL |
| Machine Learning | MlLib |
| Graph processing | Graphax |
| Stream processing | Spark Streaming |

# Why Spark?

▶ Faster



Logistic regression in Hadoop and Spark

# Why Spark: cont…

- Ease of use
- Support for multiple languages
- REPL for development and ad-hoc analysis
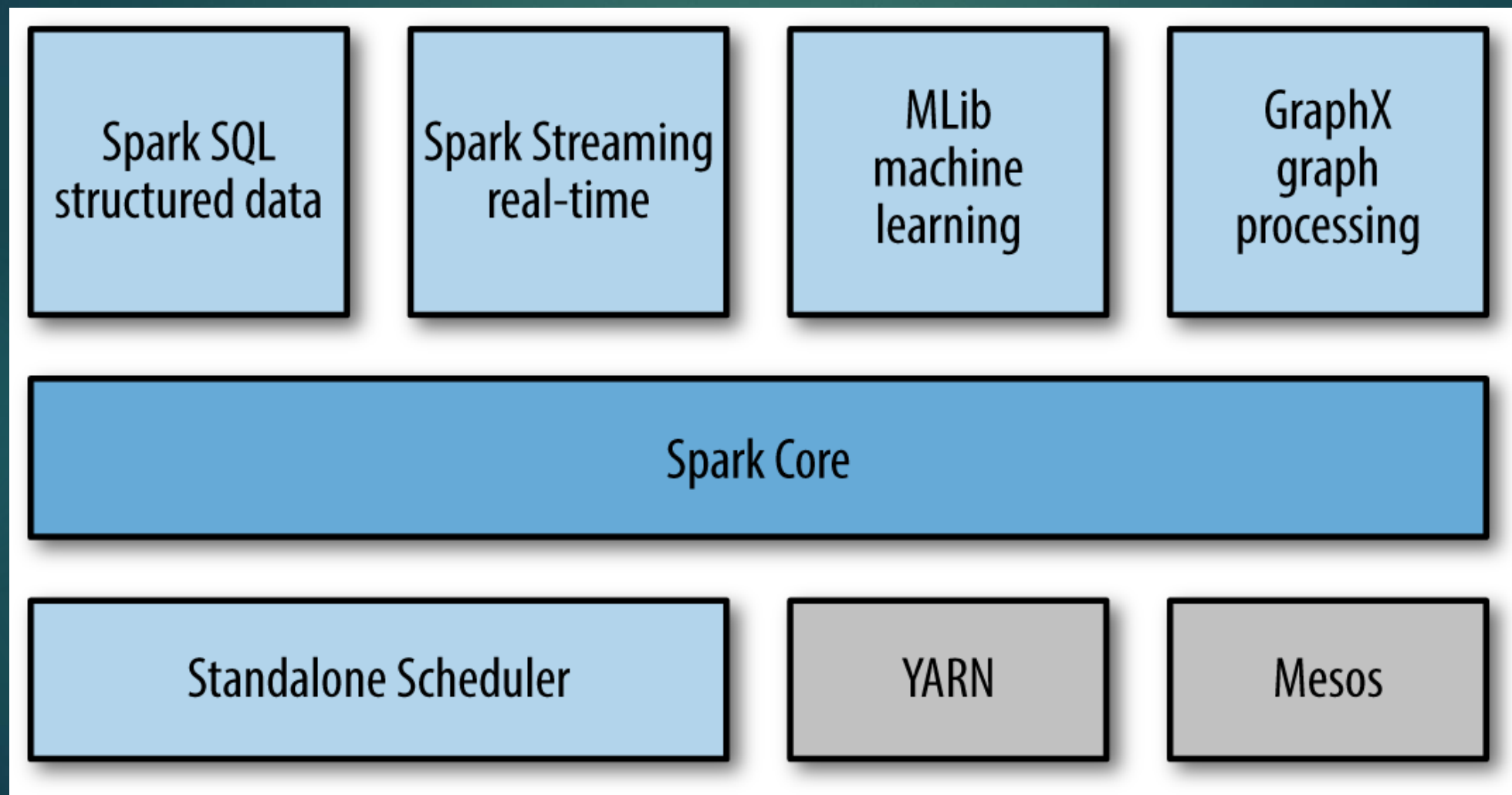- Fever lines of code

# Why Spark: cont…

- ▶ Inter operability with other platforms
- ❑ Hadoop
- ❑ Mesos
- ❑ Hbase
- ❑ Cassandra

# Components

# RDD

- **R**esilient **D**istributed **D**ataset

  Fault tolerant

  Distributed across multiple processes

  Source could be a file or program generated

- Immutable collection of elements, partitioned across multiple processes to operate in parallel

# RDD Operation: Transformation

- ▶ Evaluated lazily

- ▶ Can be applied on any RDD

- ▶ Generates another RDD as result

- ▶ Example: map, flatMap, filter, reduceByKey…

# RDD Operation: Action

- ▶ Call for evaluation of complete DAG

- ▶ Can be applied on any RDD

- ▶ Generates result on driver program
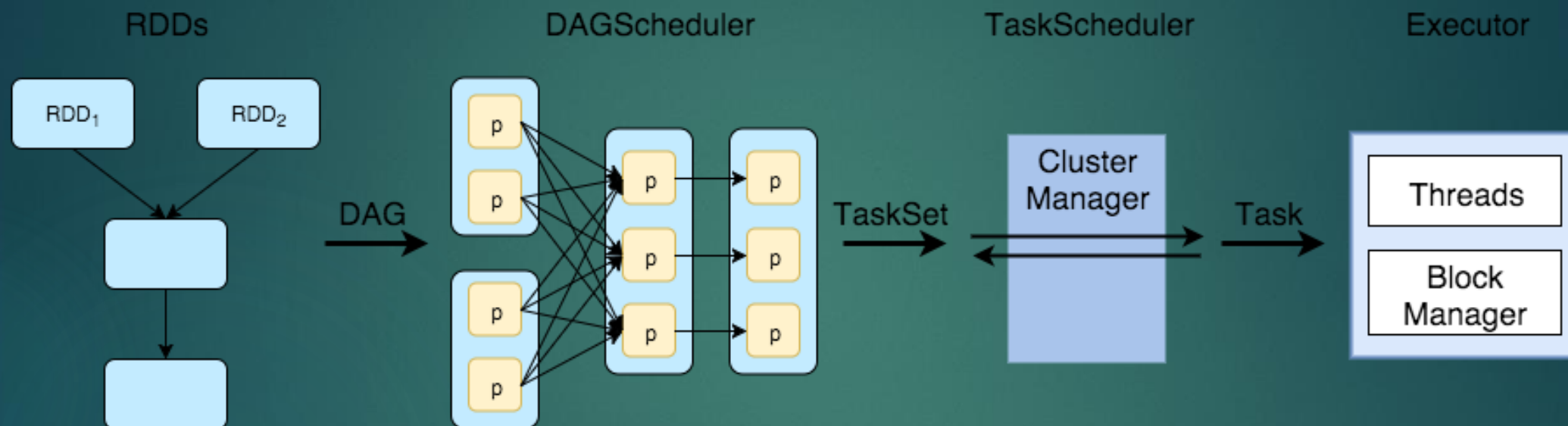
- ▶ Example: count, take, saveAsTextFile, collect…

# DAG

- ▶ Directed Acyclic Graph prepared on RDD

- ▶ DAG scheduler prepares stages and tasks

- ▶ Tasks within a stage will be executed when stage is ready to execute

- ▶ Shuffle operation is the stage boundary

# Execution and Coordination

# Launch spark shell

- Standalone:

./bin/spark-shell --master spark://IP:PORT

- YARN

./bin/spark-shell --master yarn

- Mesos

./bin/spark-shell --master mesos://host:5050

# Submit Application

▶ Scala/Java

spark-submit --class SparkWordCount --master local --deploy-mode client --executor-memory 1g --name wordcount --conf "spark.app.id=wordcount" sparkwordcount.jar <other parameters to JAR file>

▶ Python

spark-submit --master yarn --deploy-mode client --executor-memory 1g --name wordcount --conf "spark.app.id=wordcount" wordcount.py <Other parameters>

# Language comparison matrix

| Metrics | Scala | Java | Python | R |
|---|---|---|---|---|
| Type | Compiled | Compiled | Interpreted | Interpreted |
| JVM based | Yes | Yes | No | No |
| Verbosity | Less | More | Less | Less |
| Code Length | Less | More | Less | Less |
| Productivity | High | Less | High | High |
| Scalability | High | High | Less | Less |
| OOPS Support | Yes | Yes | Yes | Yes |

# REPL

- Scala

spark-shell

- Python

pyspark

# Spark SQL: Why?

▶ Integrated

*context = HiveContext(sc)*
*results = context.sql(*
  *"SELECT * FROM people")*
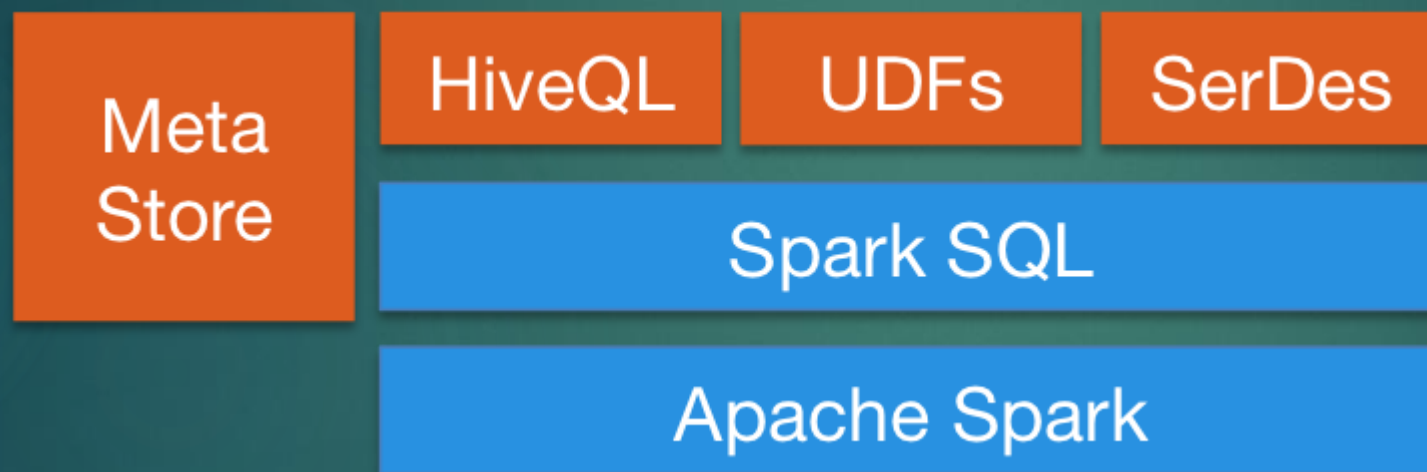*names = results.map(lambda p: p.name)*

# Spark SQL: Why?

▶ Uniform Data access

```
context.jsonFile("filename.json")
  .registerTempTable("json")
results = context.sql(
  """SELECT *
    FROM people
    JOIN json ...""")
```

# Spark SQL: Why?

- ▶ Hive Integration

# Spark SQL: Why?

▶ Standard Connectivity

# Introduction

► Relationship with Shark project

- ❑ Shark had limited integration with Spark

- ❑ Hive optimizer was not best fit for Spark

- ❑ Spark SQL Reused Hive data loading as well as in-memory column storage features of Shark

- ❑ Additionally, introduced RDD-aware optimizer and rich language interface

# Data holders

- RDD

- Dataframe: RDD with schema

- Dataset:
  introduced in 1.6 version
  provides strong type over RDD

# Data source

- ▶ Hive existing table

- ▶ Structured files. Json file for example

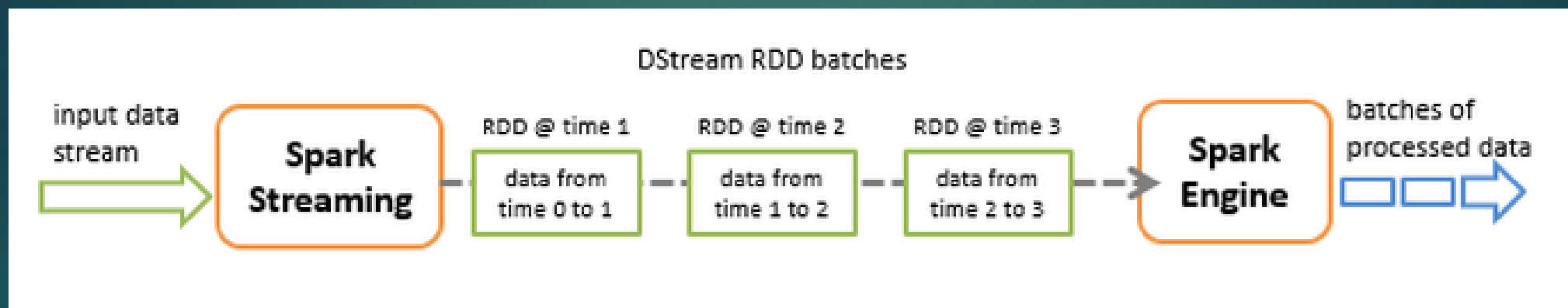- ▶ RDD

# Hive integration

- ▶ Can use most of the SQL features available in Hive

- ▶ Insert data through Spark and read in Hive

- ▶ Executes DDL statements

- ▶ Refers Hive metastore for metadata

# DStream

# Source

- Kafka
- Flume
- HDFS/S3
- Kinesis
- Twitter

# Storage/Target

- HDFS

- Databases

- Dashboard

# Program flow

- ▶ Set streaming context

- ▶ Define source for the streaming context

- ▶ Apply all transformations of Dstream

- ▶ Start the streaming context

# Examples

- https://github.com/apache/spark/tree/master/examples/src/main/scala/org/apache/spark/examples

# Persist/Cache data

▶ Helpful to reuse the same dataset

▶ Multiple storage levels:
https://spark.apache.org/docs/latest/rdd-programming-guide.html

▶ How to check current storage level:
 <Object name>.getStorageLevel

# Further studies

- https://www.cloudera.com/documentation/enterprise/5-6-x/PDF/cloudera-spark.pdf

- https://databricks.com/product/getting-started-guide/quick-start

- Cloud hosted community spark setup
  https://community.cloud.databricks.com/

# References

- http://spark.apache.org/docs/1.3.0/cluster-overview.html

- Hadoop: the definitive guide 4th edition