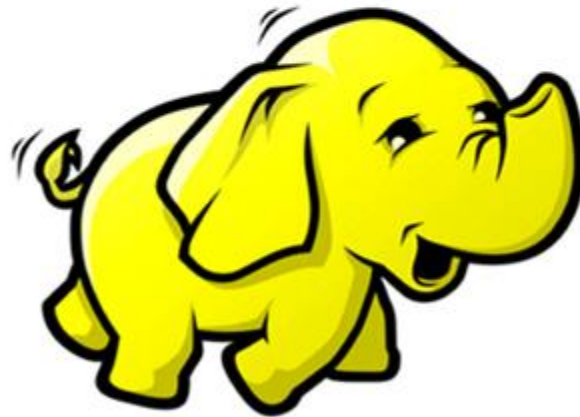




git



Recap

2

- ▶ Hive

- ▶ Impala

Hadoop

3

Processing nature	Tool
Batch processing on unstructured data	Pig
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/Kafka

Lots of tools

4



Image Ref: <http://hunteryoun.com/wp-content/uploads/2017/05/Too-many-Tools-1160x770.jpg>

One solution for all

5



Why Spark?

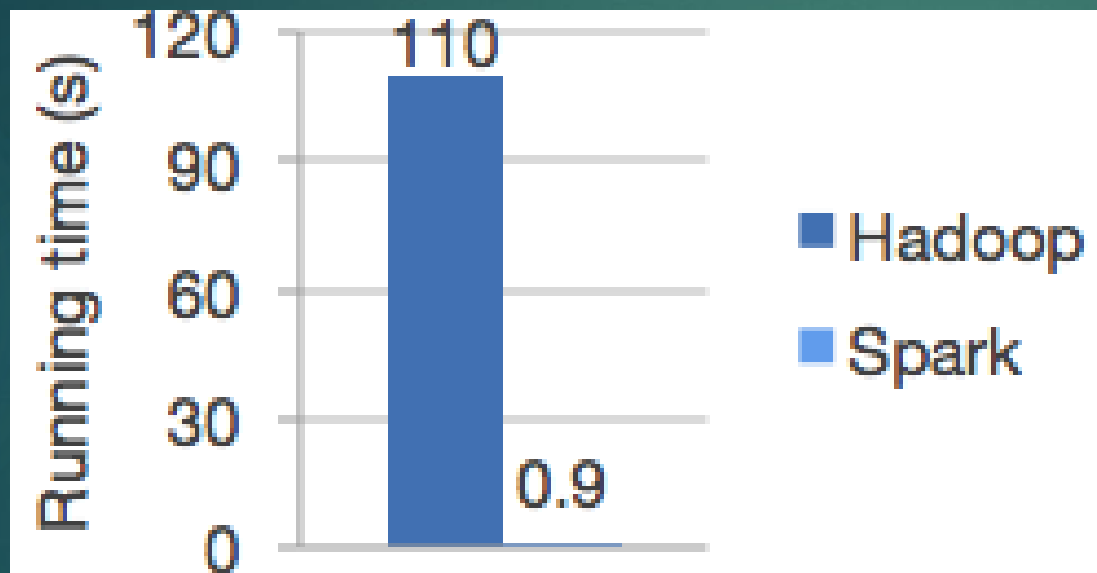
6

Processing nature	Spark
Batch processing on unstructured data	Spark RDD API
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: cont...

7

► Faster



Logistic regression in Hadoop and Spark

Why Spark: cont...

8

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

Why Spark: cont...

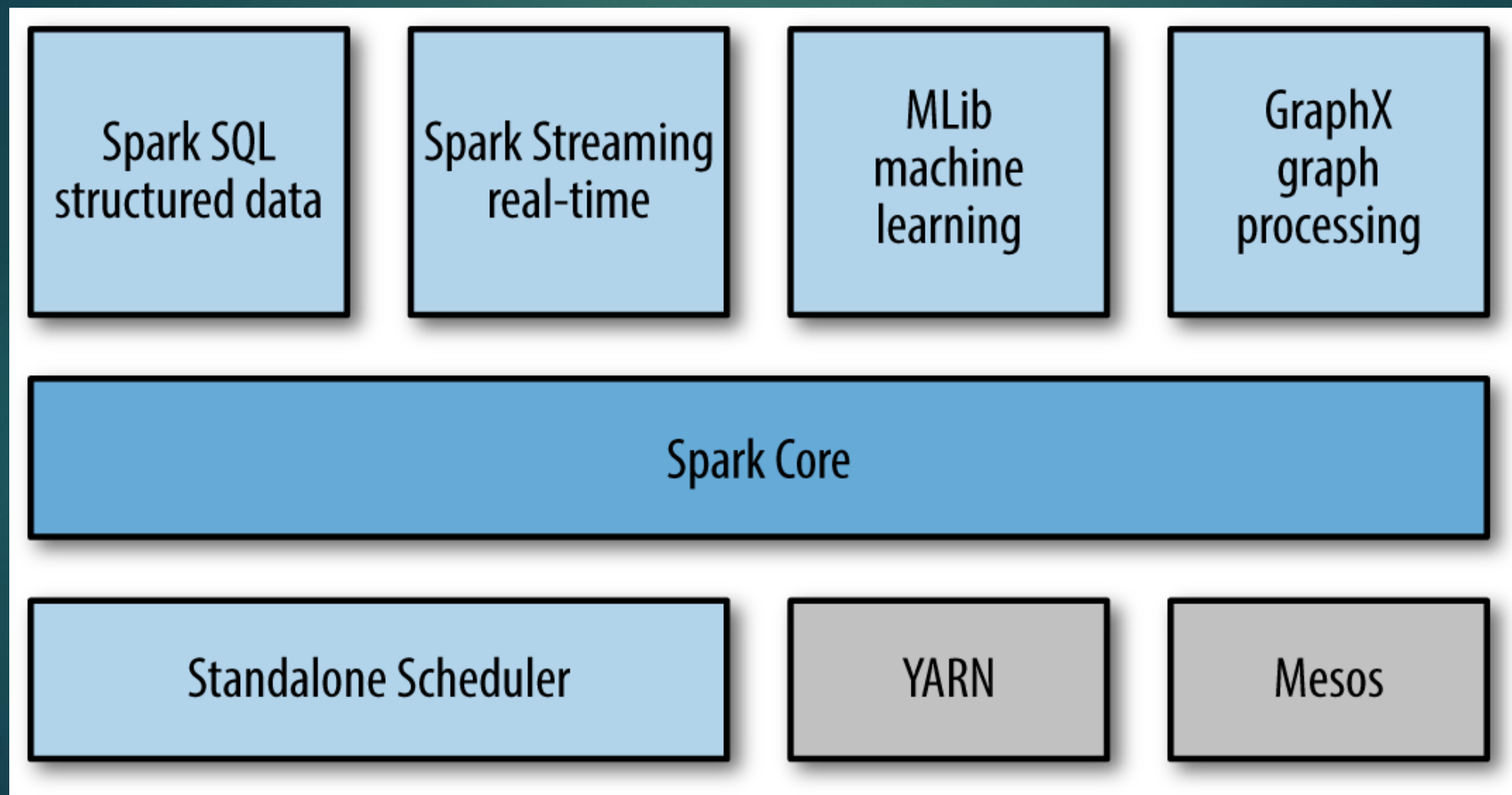
9

- ▶ Inter operability with other platforms

- Hadoop
- Mesos
- Hbase
- Cassandra

Components

10



Core of Spark: RDD

11

- ▶ **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 Operations

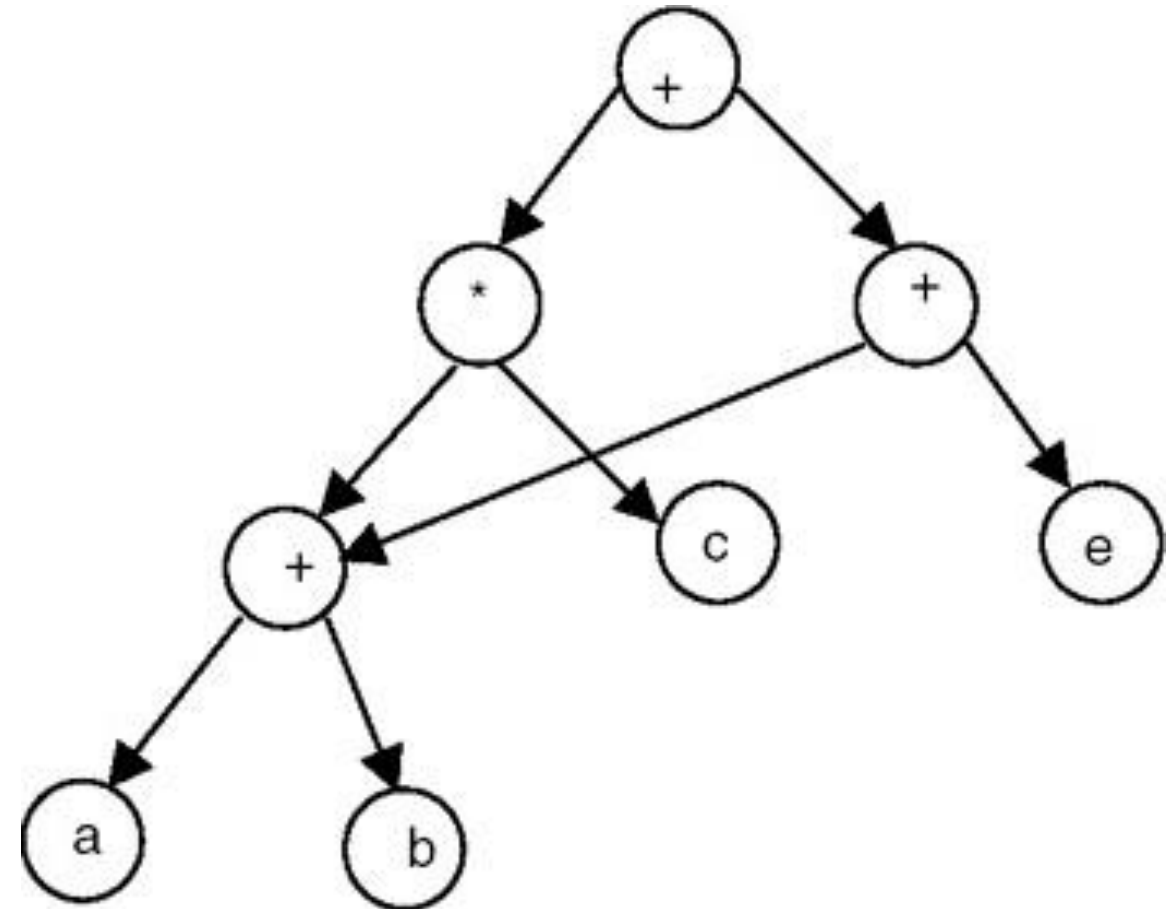
12

- ▶ Transformation

- ▶ Action

DAG

- Directed Acyclic Graph prepared to indicate task dependencies



RDD Operation: Transformation

14

- ▶ Evaluated lazily
- ▶ Can be applied on any RDD
- ▶ Generates another RDD as result
- ▶ Example: map, flatMap, filter, reduceByKey...

RDD Operation: Action

15

- ▶ Call for evaluation of complete DAG
- ▶ Can be applied on any RDD
- ▶ Generates result on driver program
- ▶ Example: count, take, saveAsTextFile, collect...

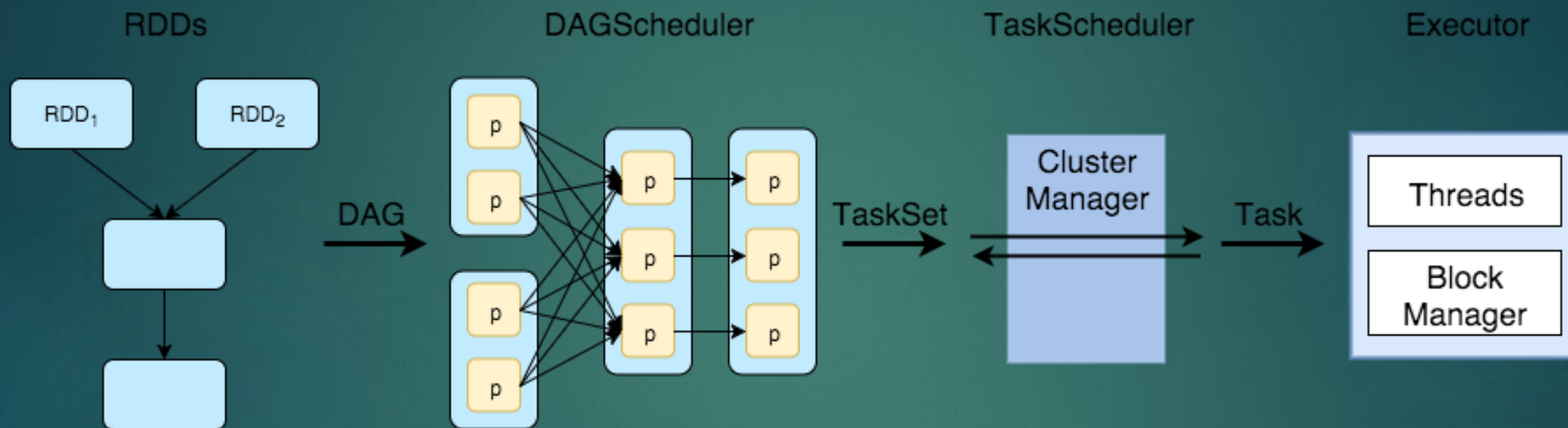
Parallelism

16

- ▶ Tasks within a stage will be executed when stage is ready to execute
- ▶ Shuffle operation is the stage boundary

Execution and Coordination

17



Launch spark shell

18

- ▶ Standalone:

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

- ▶ YARN

```
./bin/spark-shell --master yarn
```

- ▶ Mesos

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

Submit Application

19

► Scala/Java

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

► Python

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

Language comparison matrix

20

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
OOP Support	Yes	Yes	Yes	Yes

REPL

21

- ▶ Scala
spark-shell

- ▶ Python
pyspark

The logo for Spark SQL. It features the word "Spark" in a bold, black, sans-serif font. An orange, stylized five-pointed star is positioned above the letter "k". To the right of "Spark" is the text "SQL" in a black, sans-serif font, with the letters "Q", "L", and "S" being significantly larger than the letter "A".

Spark SQL

Why SparkSQL?

23

► Integrated

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

Why SparkSQL: cont...

24

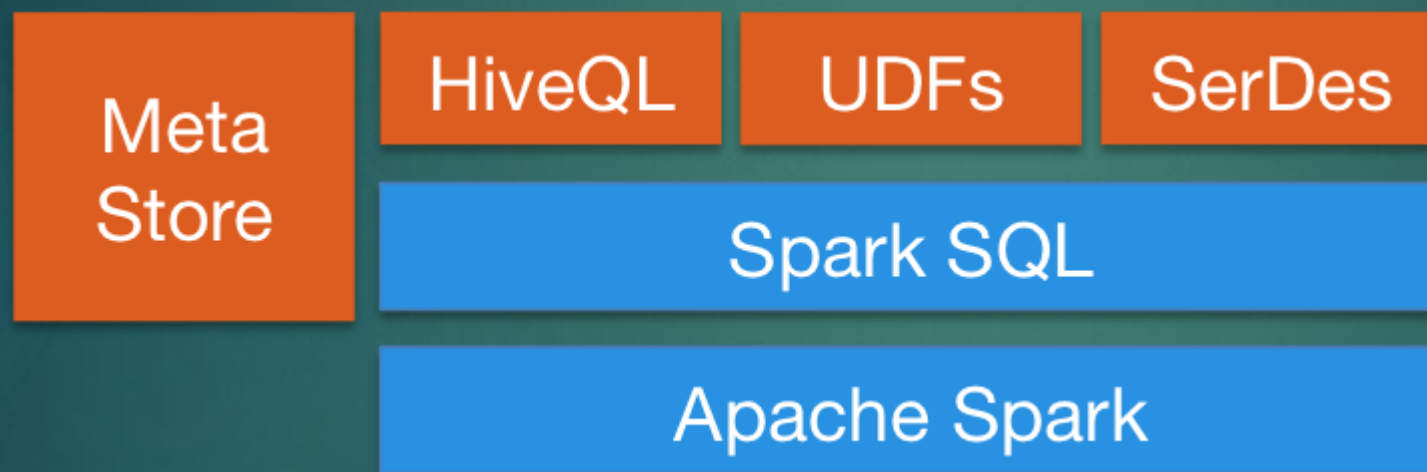
► Uniform Data access

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

Why SparkSQL: cont...

25

► Hive Integration



Why SparkSQL: cont...

26

► Standard Connectivity

BI Tools

...

JDBC / ODBC

Spark SQL

Data holders

27

- ▶ RDD
- ▶ Dataframe: RDD with schema
- ▶ Dataset:
introduced in 1.6 version
provides strong type over RDD

Data source

28

- ▶ Hive existing table
- ▶ Structured files. Json file for example
- ▶ RDD

Hive integration

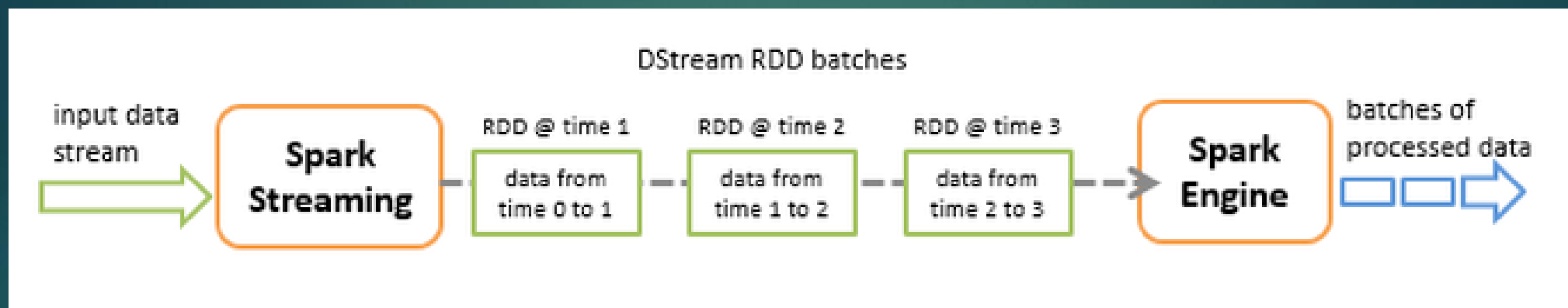
29

- ▶ 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

31



Source

32

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

Storage/Target

33

- ▶ HDFS
- ▶ Databases
- ▶ Dashboard

Program flow

34

- ▶ Set streaming context
- ▶ Define source for the streaming context
- ▶ Apply all transformations of Dstream
- ▶ Start the streaming context

Persist/Cache data

35

- ▶ 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`

Examples

36

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

References

- ▶ <http://spark.apache.org/docs/1.3.0/cluster-overview.html>
- ▶ Hadoop: the definitive guide 4th edition
- ▶ <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/>