

# Agenda



- Spark Platform
- Spark Core
- Spark Extensions
- Using Apache Spark

# About me

Vitalii Bondarenko

Data Platform Competency Manager

**Eleks**

[www.eleks.com](http://www.eleks.com)

**20 years in software development**

**9+ years of developing for MS SQL**

**3+ years of architecting Big Data So**

# eleks<sup>®</sup>

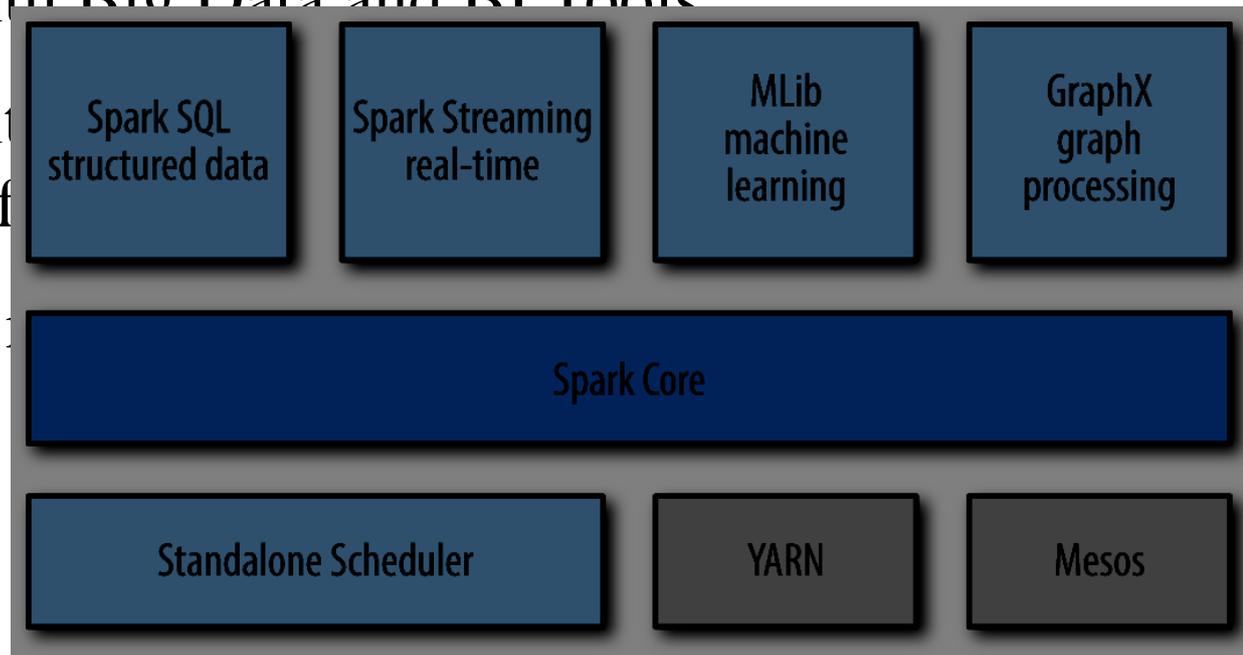


# Spark Stack

- Clustered computing platform
- Designed to be fast and general purpose
- Integrated with distributed systems
- API for Python, Scala, Java, clear and understandable code
- Integrated with Big Data and BI Tools

•Integrated with  
Cassandra, Kaf

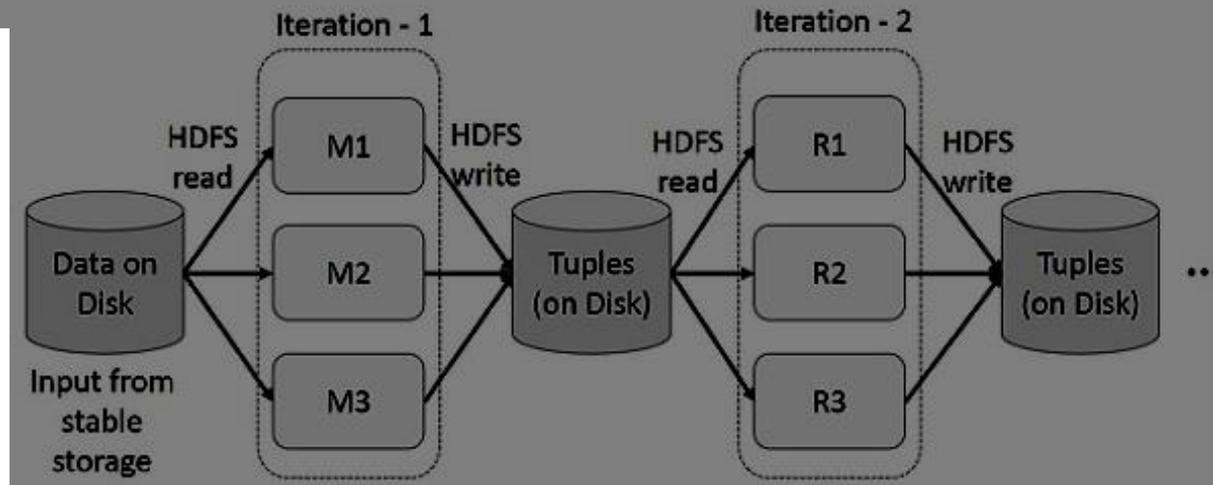
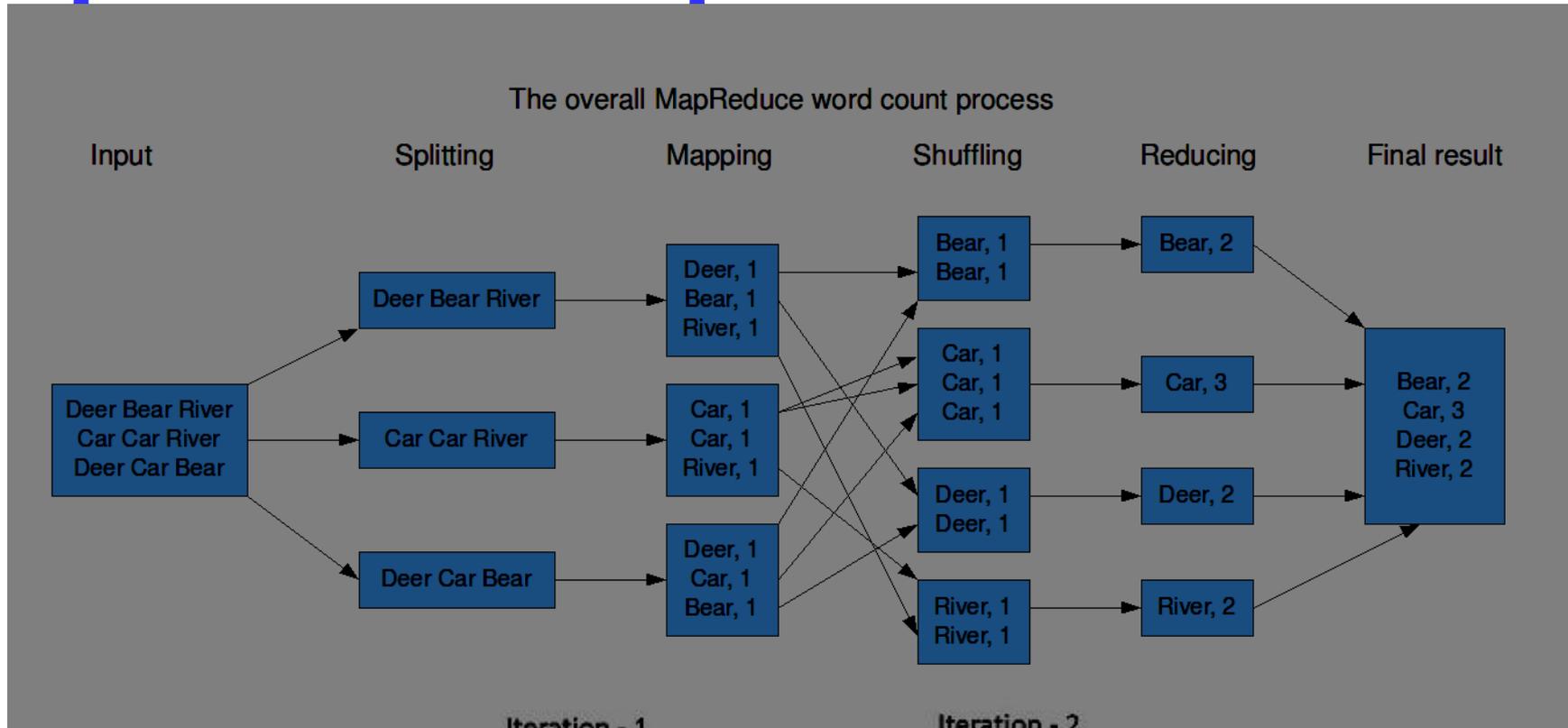
•First Apache



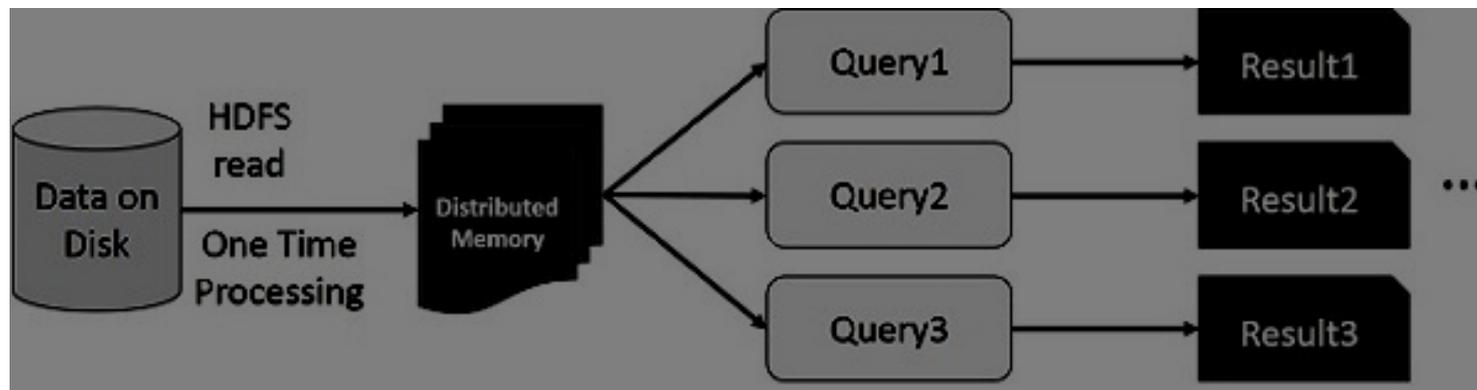
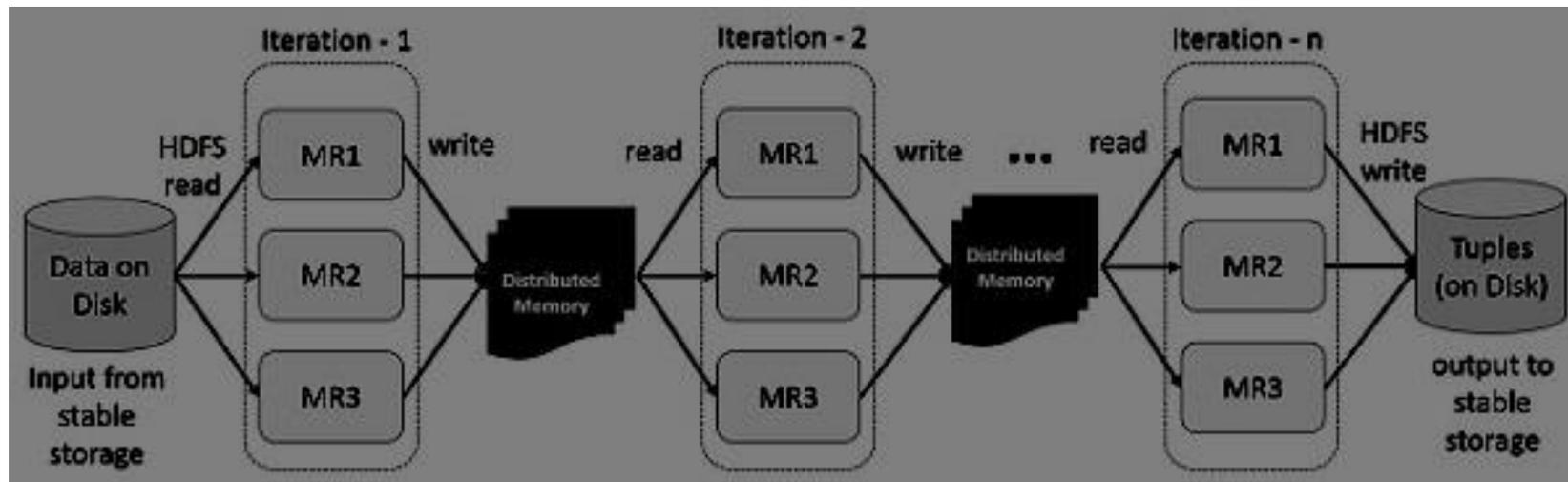
like

ed

# Map-reduce computations



# In-memory map-reduce



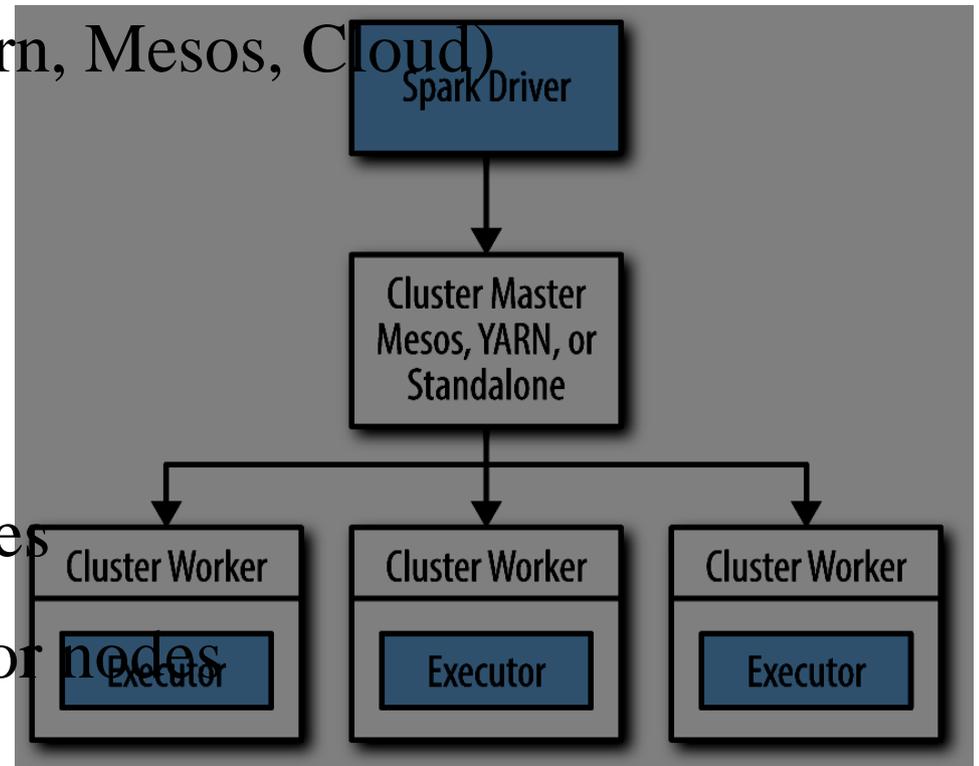
# Execution Model

## Spark Execution

- Shells and Standalone application
- Local and Cluster (Standalone, Yarn, Mesos, Cloud)

## Spark Cluster Architecture

- Master / Cluster manager
- Cluster allocates resources on nodes
- Master sends app code and tasks to nodes
- Executors run tasks and cache data



# RDD: resilient distributed dataset

- Parallelized collections with fault-tolerant (Hadoop datasets)

- **Transformations** set new RDDs (filter, map, distinct, union, subtract, etc)

- **Actions** call to calculations (count, collect, first)

- Transformations are lazy

- Actions trigger transformations computation

```
from pyspark import SparkContext as sc
```

- **Broadcast Variables** send data to executors

```
inputRDD = sc.textFile("log.txt")
```

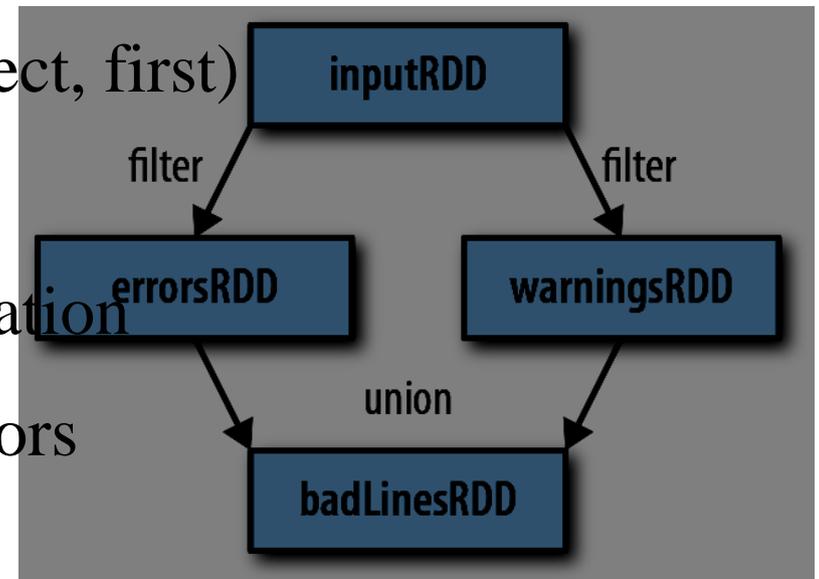
- **Accumulators** collect data on driver

```
errorsRDD = inputRDD.filter(lambda x: "error" in x)
```

```
warningsRDD = inputRDD.filter(lambda x: "warning" in x)
```

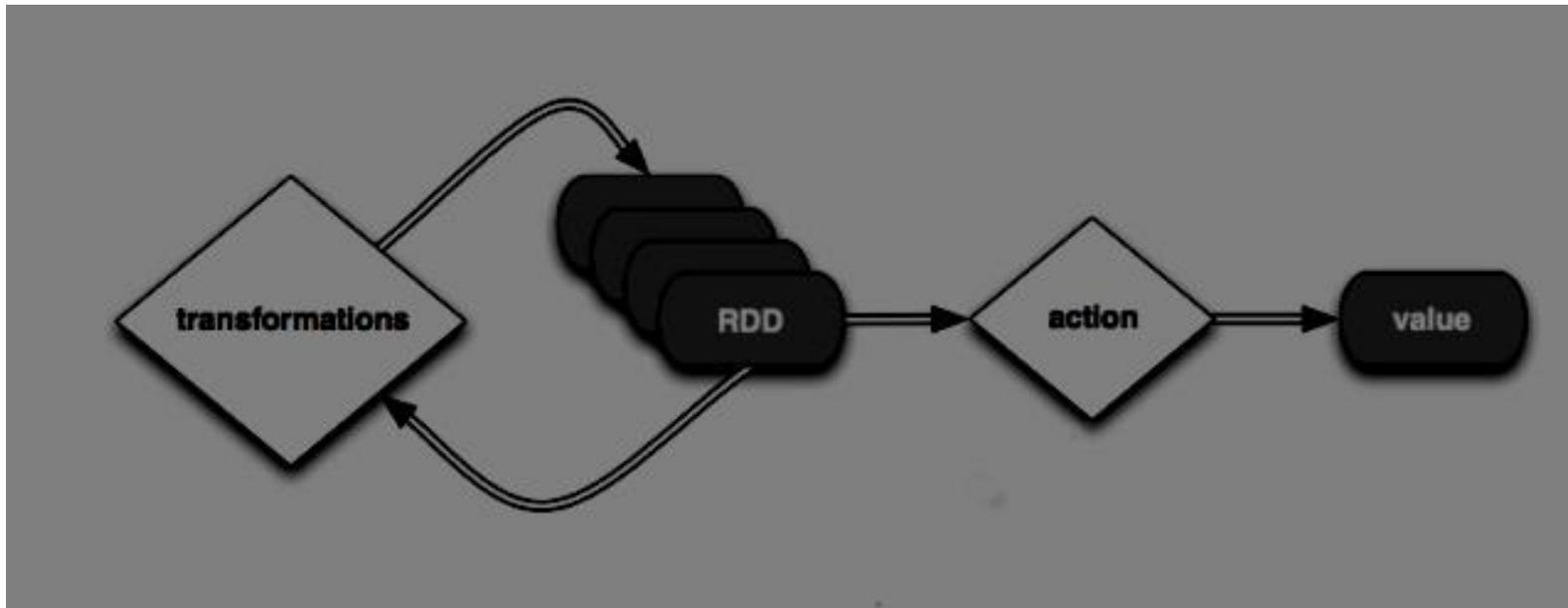
```
badLinesRDD = errorsRDD.union(warningsRDD)
```

```
print "Input had " + badLinesRDD.count() + " concerning lines"
```



# Spark program scenario

- Create RDD (loading external datasets, parallelizing a collection on driver)
- Transform
- Persist intermediate RDDs as results
- Launch actions



# Transformations (1)

<b>Function name</b>	<b>Purpose</b>	<b>Example</b>	<b>Result</b>
<b>map()</b>	Apply a function to each element in the RDD and return an RDD of the result.	<code>rdd.map(x =&gt; x + 1)</code>	{2, 3, 4, 4}
<b>flatMap()</b>	Apply a function to each element in the RDD and return an RDD of the contents of the iterators returned. Often used to extract words.	<code>rdd.flatMap(x =&gt; x.to(3))</code>	{1, 2, 3, 2, 3, 3, 3}
<b>filter()</b>	Return an RDD consisting of only elements that pass the condition passed to filter().	<code>rdd.filter(x =&gt; x != 1)</code>	{2, 3, 3}
<b>distinct()</b>	Remove duplicates.	<code>rdd.distinct()</code>	{1, 2, 3}

# Transformations (2)

<b>Function name</b>	<b>Purpose</b>	<b>Example</b>	<b>Result</b>
<code>union()</code>	Produce an RDD containing elements from both RDDs.	<code>rdd.union(other)</code>	{1, 2, 3, 3, 4, 5}
<code>intersection()</code>	RDD containing only elements found in both RDDs.	<code>rdd.intersection(other)</code>	{3}
<code>subtract()</code>	Remove the contents of one RDD (e.g., remove training data).	<code>rdd.subtract(other)</code>	{1, 2}
<code>cartesian()</code>	Cartesian product with the other RDD.	<code>rdd.cartesian(other)</code>	{(1, 3), (1, 4), ... (3, 5)}

# Actions (1)

Function name	Purpose	Example	Result
<code>collect()</code>	Return all elements from the RDD.	<code>rdd.collect()</code>	{1, 2, 3, 3}
<code>count()</code>	Number of elements in the RDD.	<code>rdd.count()</code>	4
<code>countByValue()</code>	Number of times each element occurs in the RDD.	<code>rdd.countByValue()</code>	{(1, 1), (2, 1), (3, 2)}
<code>take(num)</code>	Return num elements from the RDD.	<code>rdd.take(2)</code>	{1, 2}
<code>top(num)</code>	Return the top num elements the RDD.	<code>rdd.top(2)</code>	{3, 3}

# Actions (2)

<b>takeOrdered(num) (ordering)</b>	Return num elements based on provided ordering.	<b>rdd.takeOrdered(2) (myOrdering)</b>	<b>{3, 3}</b>
<b>reduce(func)</b>	Combine the elements of the RDD together in parallel (e.g., sum).	<b>rdd.reduce((x, y) =&gt; x + y)</b>	<b>9</b>
<b>fold(zero)(func)</b>	Same as <b>reduce()</b> but with the provided zero value.	<b>rdd.fold(0)((x, y) =&gt; x + y)</b>	<b>9</b>
<b>aggregate(zeroValue) (seqOp, combOp)</b>	Similar to <b>reduce()</b> but used to return a different type.	<b>rdd.aggregate((0, 0)) ((x, y) =&gt; (x._1 + y, x._2 + 1), (x, y) =&gt; (x._1 + y._1, x._2 + y._2))</b>	<b>(9, 4)</b>

# Spark Streaming Architecture



- Micro-batch architecture
- SparkStreaming Concept
- Batch interval from 500ms
- Transformation on Spark Engine
- Output operations instead of Actions
- Different sources and outputs



# Spark Streaming Example

```
from pyspark.streaming import StreamingContext

ssc = StreamingContext(sc, 1)

input_stream = ssc.textFileStream("sampleTextDir")

word_pairs = input_stream.flatMap(
    lambda l:l.split(" ")).map(lambda w: (w,1))

counts = word_pairs.reduceByKey(lambda x,y: x + y)

counts.print()

ssc.start()

ssc.awaitTermination()
```

- Process RDDs in batches
- Start after ssc.start()
- Output to console on Driver
- Awaiting termination

# Spark SQL

- SparkSQL interface for working with structured data by SQL
- Works with Hive tables and HiveQL
- Works with files (Json, Parquet etc) with defined schema
- JDBC/ODBC connectors for BI tools
- Integrated with Hive and HiveQL
- DataFrame abstraction



# Spark DataFrames

```
# Import Spark SQL from pyspark.sql
import HiveContext, Row
```

```
# Or if you can't include the hive requirements from pyspark.sql
import SQLContext, Row
```

```
sc = new SparkContext(...)
hiveCtx = HiveContext(sc)
sqlContext = SQLContext(sc)
```

```
input = hiveCtx.jsonFile(inputFile)
```

```
# Register the input schema RDD
input.registerTempTable("tweets")
```

```
# Select tweets based on the retweet
CounttopTweets = hiveCtx.sql("""SELECT text, retweetCount FROM tweets ORDER BY retweetCount LIMIT 10""")
```

• `hiveCtx.cacheTable("tableName")`, in-memory, column-store, while driver is alive

• `df.show()`

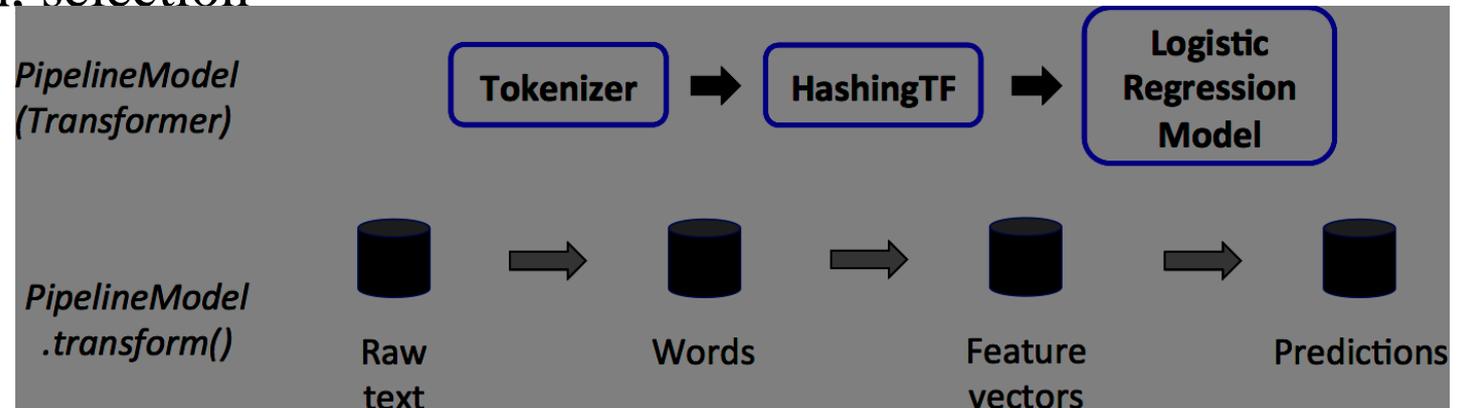
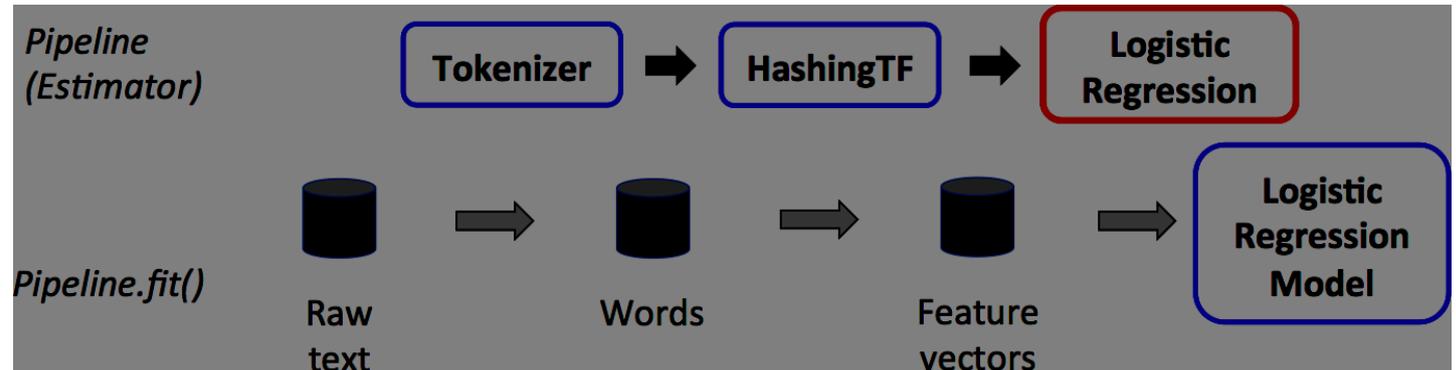
• `df.select("name", df("age")+1)`

• `df.filtr(df("age") > 19)`

# Spark ML

## Spark ML

- Classification
- Regression
- Clustering
- Recommendation
- Feature transformation, selection
- Statistics
- Linear algebra
- Data mining tools



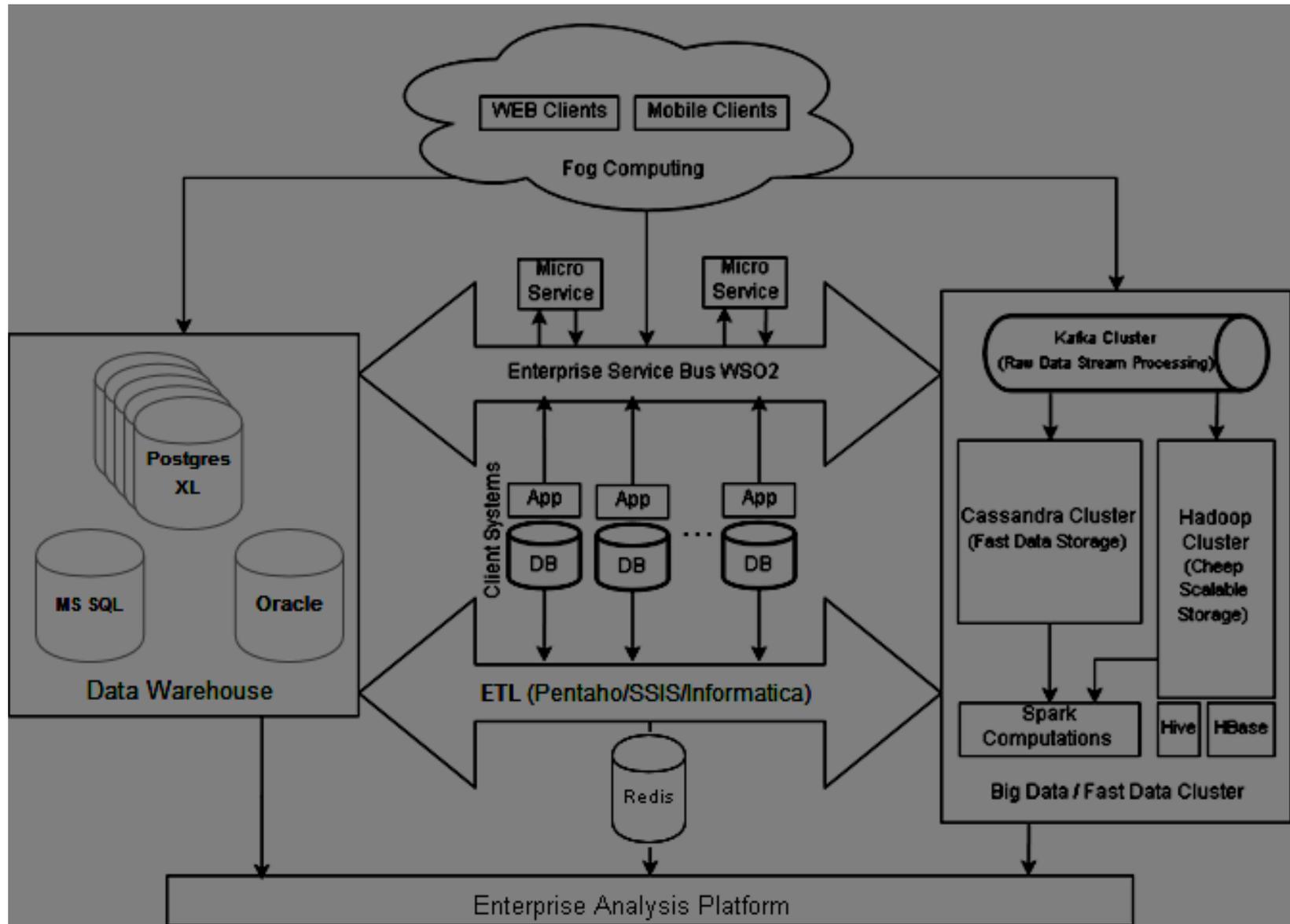
## Pipeline Components

- DataFrame

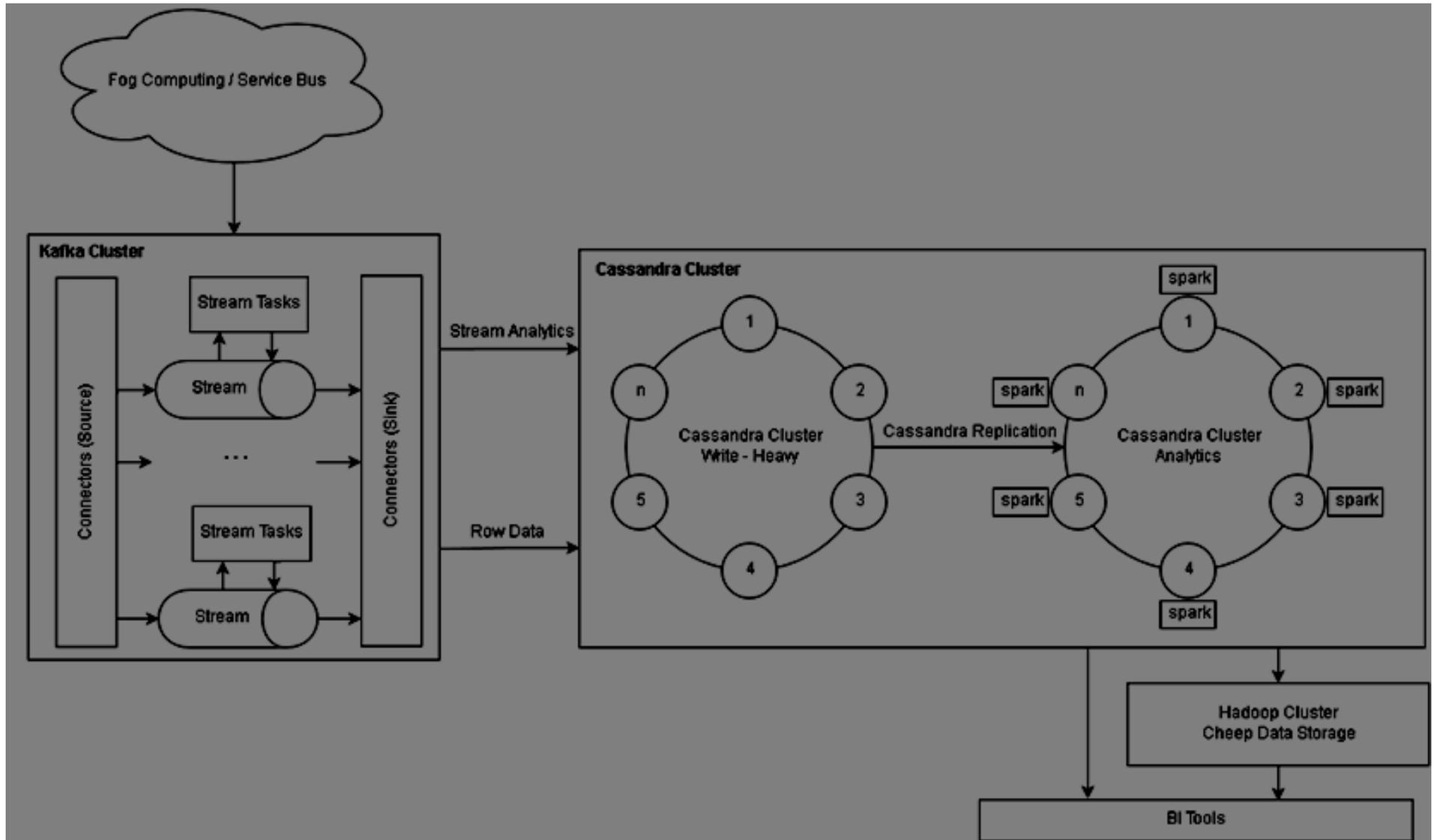
# DEMO: Spark

- Local Spark installation
- Shells and Notebook
- Spark Examples
- HDInsight Spark Cluster
- SSH connection to Spark in Azure
- Jupyter Notebook connected to HDInsight Spark
- Transformations
- ActionsSimple SparkSQL querying
- Data Frames
- Data exploration with SparkSQL
- Connect from BI
- Training a model

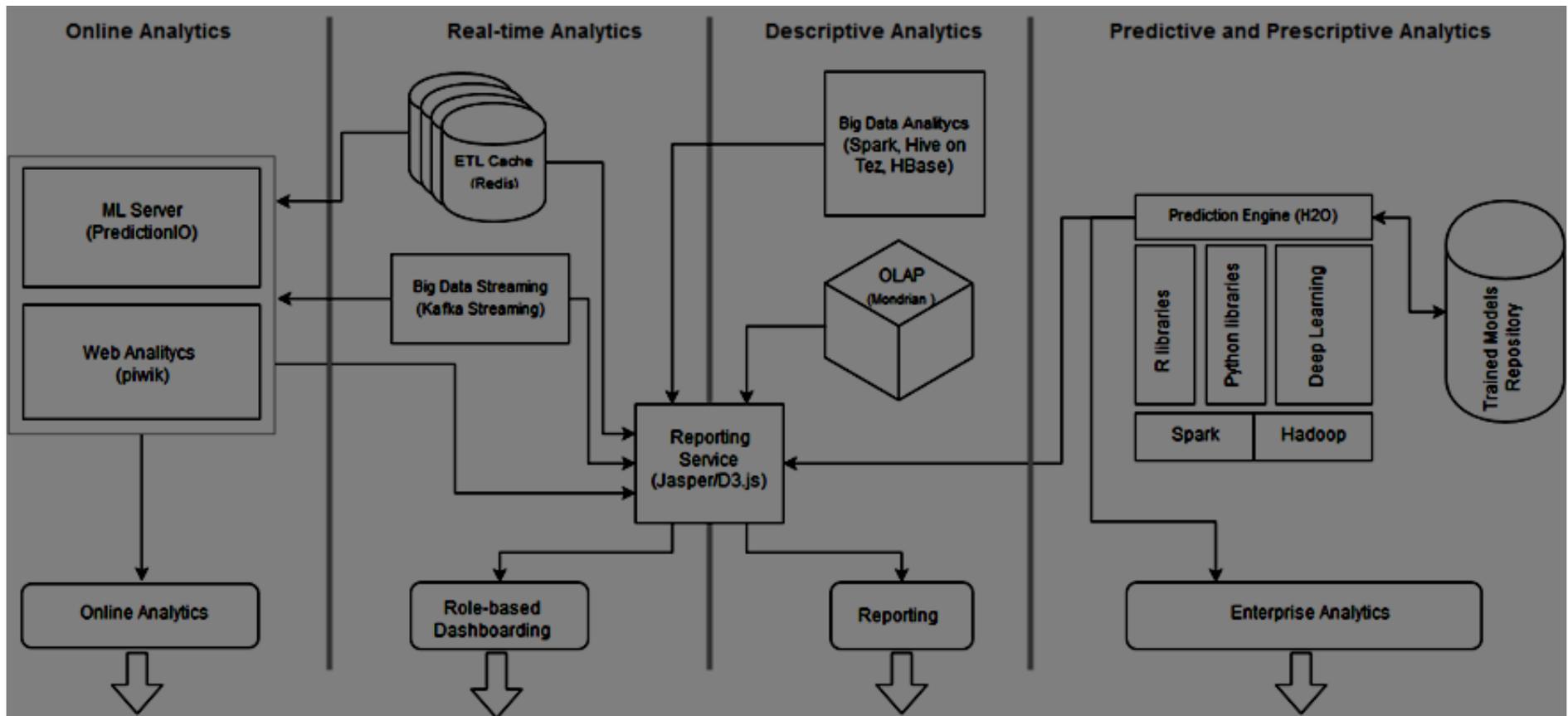
# Eleks Enterprise Platform



# Fast Data Processing



# BI / DS Platform



# Using Spark

1. Visual data exploration and interactive analysis (HDFS)
2. Spark with NoSQL (HBase and Cassandra)
3. Spark with Data Lake
4. Spark with Data Warehouse
5. Machine Learning using R Server, Mlib
6. Putting it all together in a notebook experience
7. Using BI with Spark

## **8. Spark Environmens**

- On-Premis, Cloudera, Hortonworks, DataBricks
- HDInsight, AWS, DataBricks Cloud
- Sparkling Water (H2O), prediction.io

**Q&A**