



## Intro to Apache Spark

---

## Agenda

- Bio
- IBM
- Intro to Spark
- Jupyter Notebook in Data Science Experience

---

## Speaker

Joseph Kambourakis

- Open Source Analytics Technical Evangelist at IBM
- [joseph.kambourakis@ibm.com](mailto:joseph.kambourakis@ibm.com)

Previously

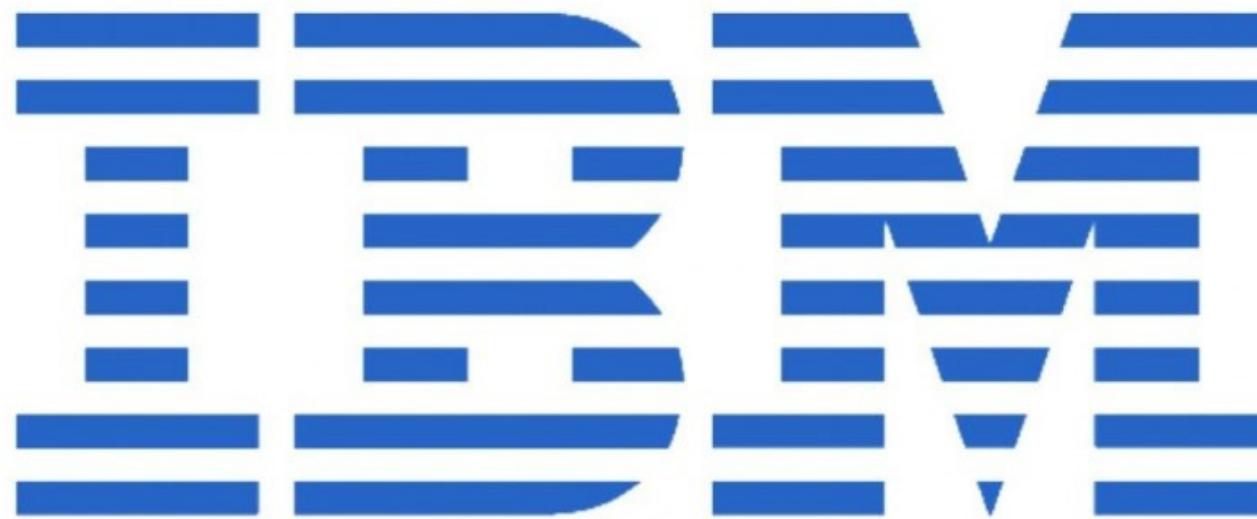
- Cloudera
- EMC Dell

Education

- MBA from Bentley University
- BS from WPI

---

What do you think of?



---

## Founding Member





# Linux

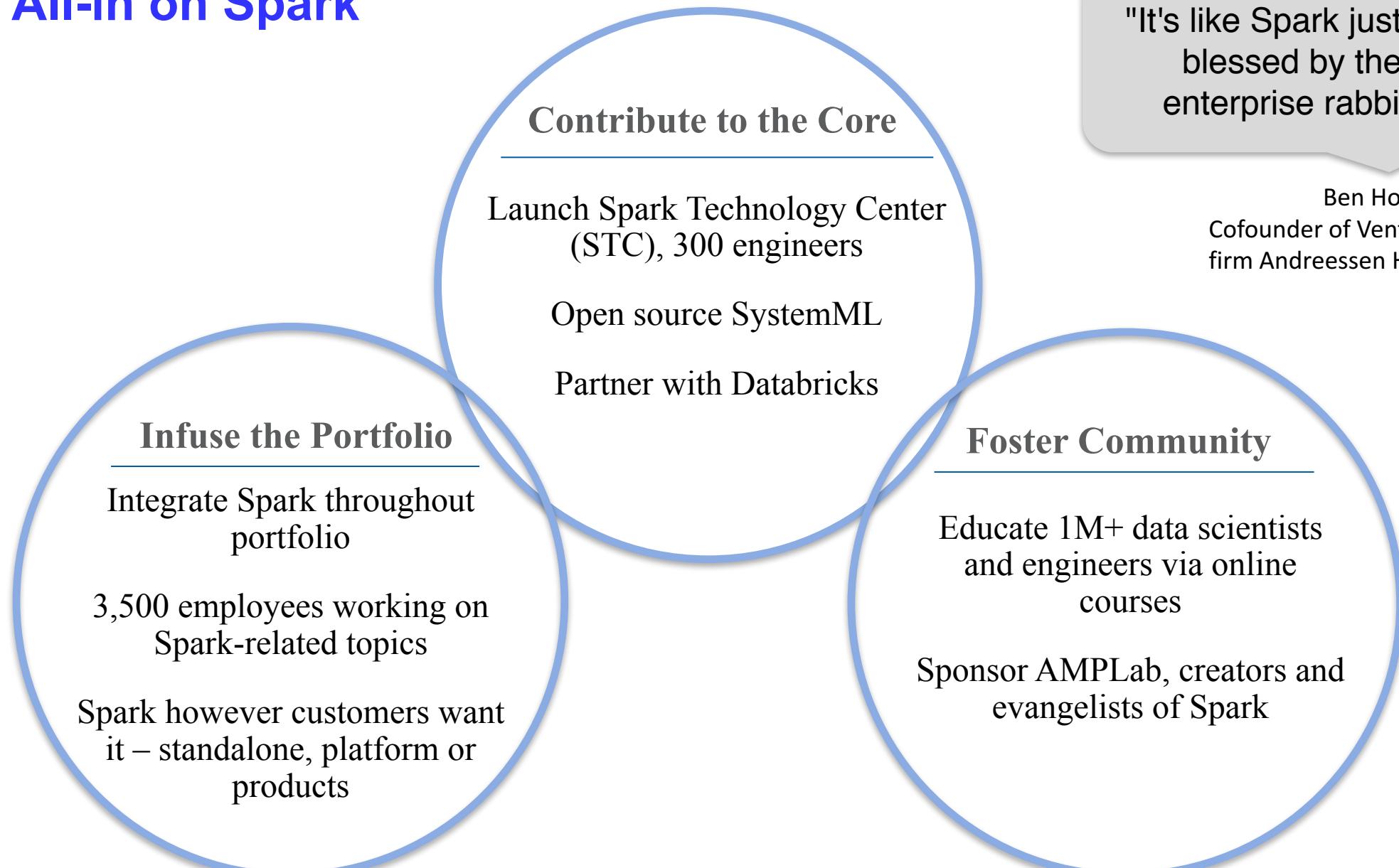
---

## Platinum Member





# IBM is All-in on Spark



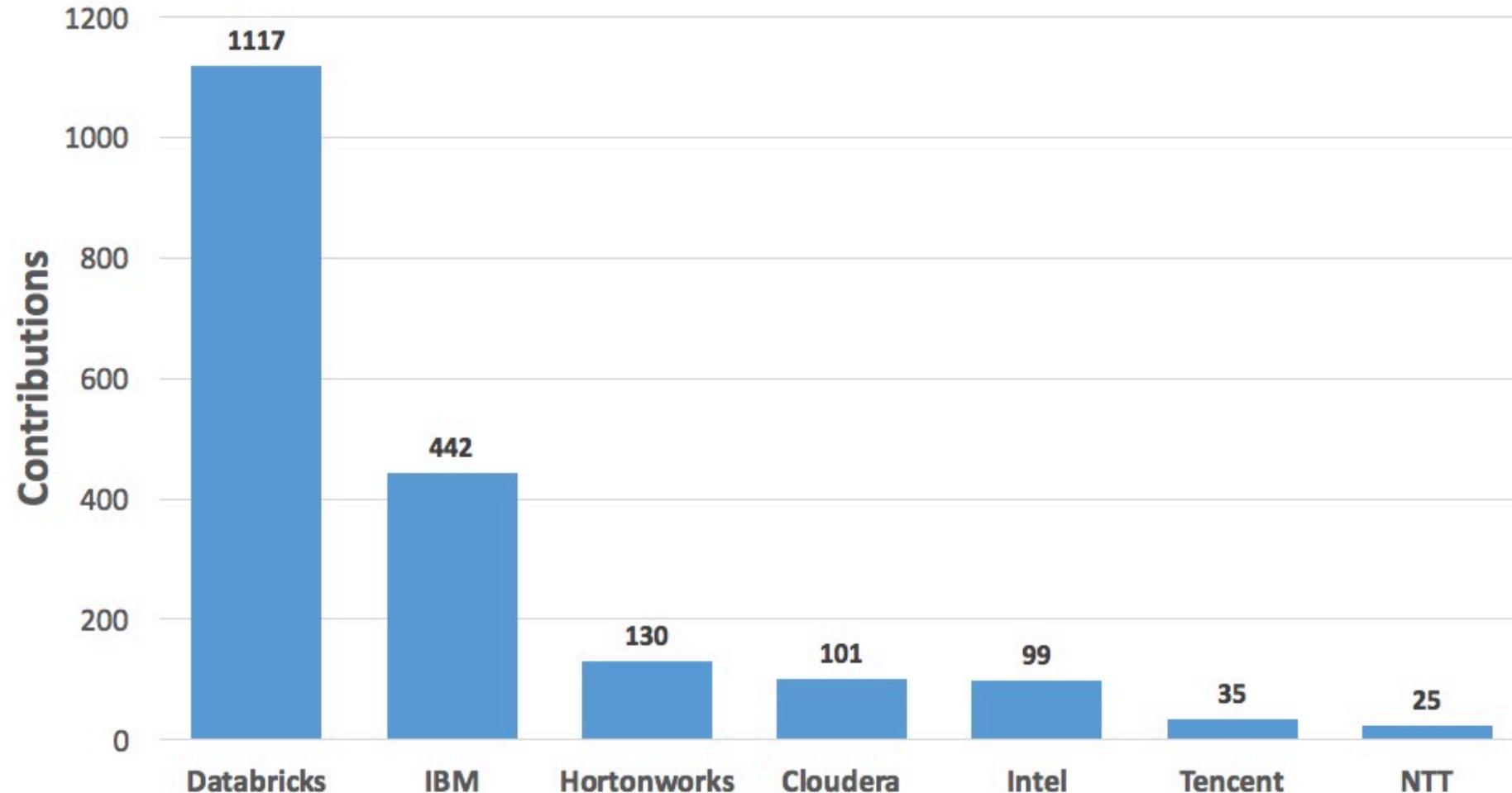
"It's like Spark just got blessed by the enterprise rabbi."

Ben Horowitz,  
Cofounder of Venture Capital  
firm Andreessen Horowitz

# Contribute to the Core: Actual Code!

## Top 7 Contributing Companies to Spark 2.0.0

(Represents 70% of total contributions)



# Apache Hadoop Ecosystem Committers

Committers	Projects
Juanjo Marron	Ambari
Eric Yang	Ambari, Hadoop
Di Li	Ambari
Tim Thorpe	Ambari
Tanping Wang	Hadoop
Songqing (Richardd) Ding	Pig
Shai Erera	Lucene
Tuong Truong	Sentry
Jing He	HBase
Doug Miel	HBase
Prashant Sharma	Spark
Trevor Grant	Mahout
Holden Karau	Spark

## New Hires



# Spark Technology Center



## Spark Summit: Platinum Sponsor



## Local Level: Meetups & Universities



# Data Science Tools: What do you use?



# Foster the Community: Data Science Experience

 IBM Data Science Experience

Log In

Sign Up

## Master the art of data science.

Enjoy the best of open source and collaborate in a social environment, built for data scientists by data scientists.

Get Started

Watch the video

This is a preview of Data Science Experience, stay tuned for additional features.

Use what you know, learn  
what you don't

Start a course, start from a sample, or

---

## Data Science Tools: We built it out of leaders



# Data Science Tools: We're actually a leader



# What is Apache Spark?



## What is Spark?

Spark is an **open** source  
**in-memory**  
**application** framework for  
**distributed** data processing and  
**iterative** analysis  
on **massive** data volumes

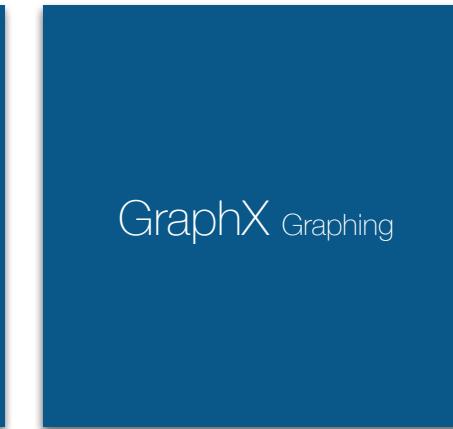
**“Analytic Operating System”**

## Introducing Apache Spark

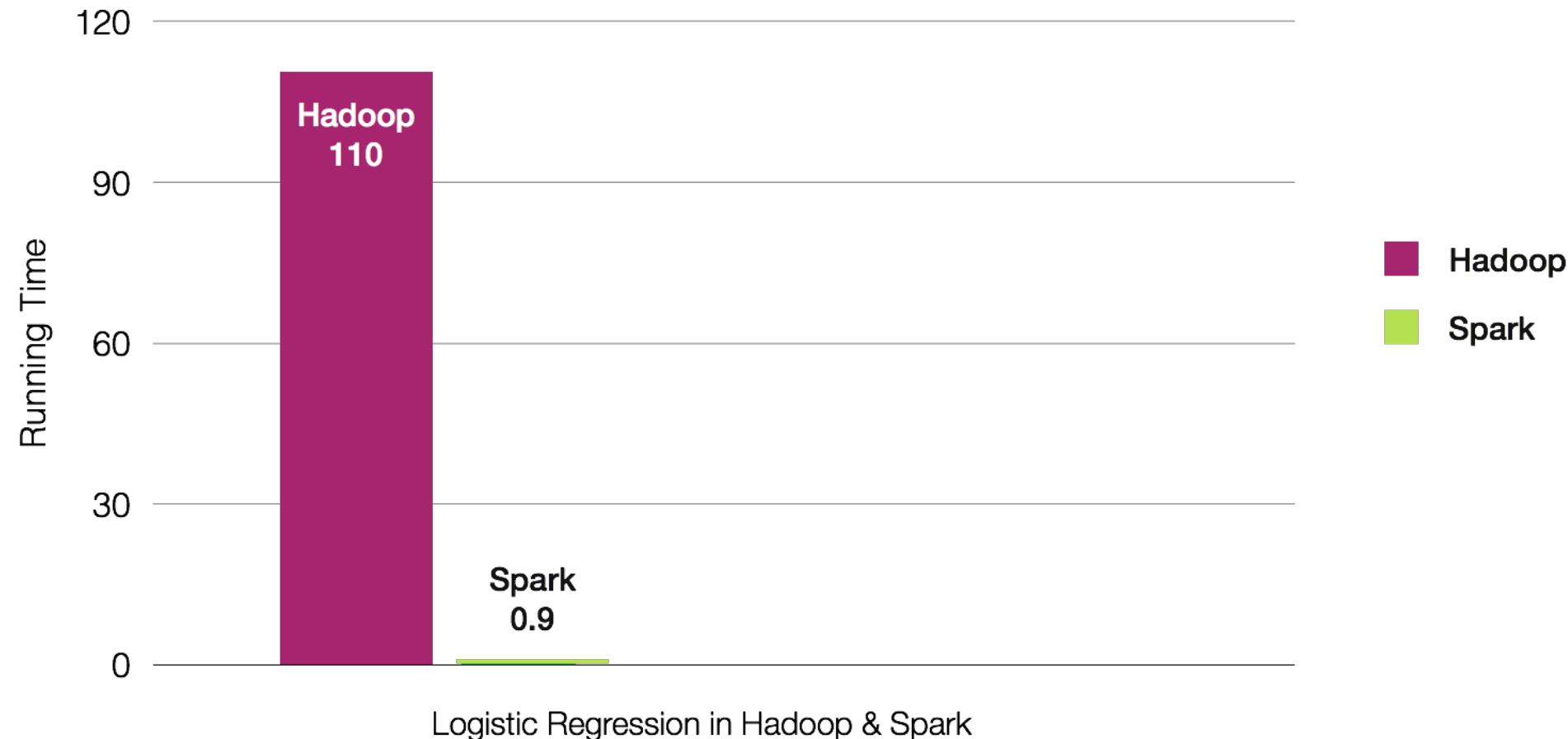
- **General purpose distributed data processing engine**
- **AMPLab at UC Berkeley**
  - Started by Matei Zaharia in 2009
  - Creators founded Databricks
    - IBM partners with DB
- **Apache open source project**
  - Initial release May 2014
  - Current version 2.1 released December 2016



# Spark and the Key Libraries



# Why Spark? Performance



# Why Spark? Productive

- Production ready code
- Concise
- Focus on the business problems
- and not coding

## Word Count

In this example, we use a few transformations to build a dataset of (String, Int) pairs called counts and then save it to a file.

Python    Scala    Java

```
text_file = sc.textFile("hdfs://...")
counts = text_file.flatMap(lambda line: line.split(" "))
    .map(lambda word: (word, 1))
    .reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("hdfs://...")
```

### Example: WordCount v1.0

Before we jump into the details, let's walk through an example MapReduce application to get a flavour for how they work.

WordCount is a simple application that counts the number of occurrences of each word in a given input set.

This works with a local-standalone, pseudo-distributed or fully-distributed Hadoop installation ([Single Node Setup](#)).

#### Source Code

```
1. package org.myorg;
2.
3. import java.io.IOException;
4. import java.util.*;
5.
6. import org.apache.hadoop.fs.Path;
7. import org.apache.hadoop.conf.*;
8. import org.apache.hadoop.io.*;
9. import org.apache.hadoop.mapred.*;
10. import org.apache.hadoop.util.*;
11.
12. public class WordCount {
13.
14.     public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
15.         private final static IntWritable one = new IntWritable(1);
16.         private Text word = new Text();
17.
18.         public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
19.             String line = value.toString();
20.             StringTokenizer tokenizer = new StringTokenizer(line);
21.             while (tokenizer.hasMoreTokens()) {
22.                 word.set(tokenizer.nextToken());
23.                 output.collect(word, one);
24.             }
25.         }
26.     }
27.
28.     public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
29.         public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
30.             int sum = 0;
31.             for (IntWritable val : values) {
32.                 sum += val.get();
33.             }
34.             output.collect(key, new IntWritable(sum));
35.         }
36.     }
37.
38.     public static void main(String[] args) throws Exception {
39.         JobConf conf = new JobConf(WordCount.class);
40.         conf.set("mapred.mapper.class", WordCount.Map.class);
41.         conf.set("mapred.reducer.class", WordCount.Reduce.class);
42.         conf.set("mapred.inputformat", TextInputFormat.class.getName());
43.         conf.set("mapred.outputformat", TextOutputFormat.class.getName());
44.         conf.setMapperClass(Map.class);
45.         conf.setCombinerClass(Reduce.class);
46.         conf.setReducerClass(Reduce.class);
47.         conf.setInputFormat(TextInputFormat.class);
48.         conf.setOutputFormat(TextOutputFormat.class);
49.         FileInputFormat.setInputPaths(conf, new Path(args[0]));
50.         FileOutputFormat.setOutputPath(conf, new Path(args[1]));
51.         JobClient.runJob(conf);
52.     }
53. }
```

[WordCount.java](#)

# Why Spark? Community

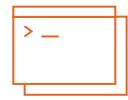
- Open source community
- Fastest growing
- Wide adoption



SPARK MEETUP  
MEMBERS

+240%

2015      2016  
66,000      225,000



CODE  
CONTRIBUTORS

+67%

2015      2016  
600      1000



SPARK SUMMIT  
ATTENDEES

+30%

2015      2016  
3912      5100



NUMBER OF COMPANIES  
AT SUMMITS

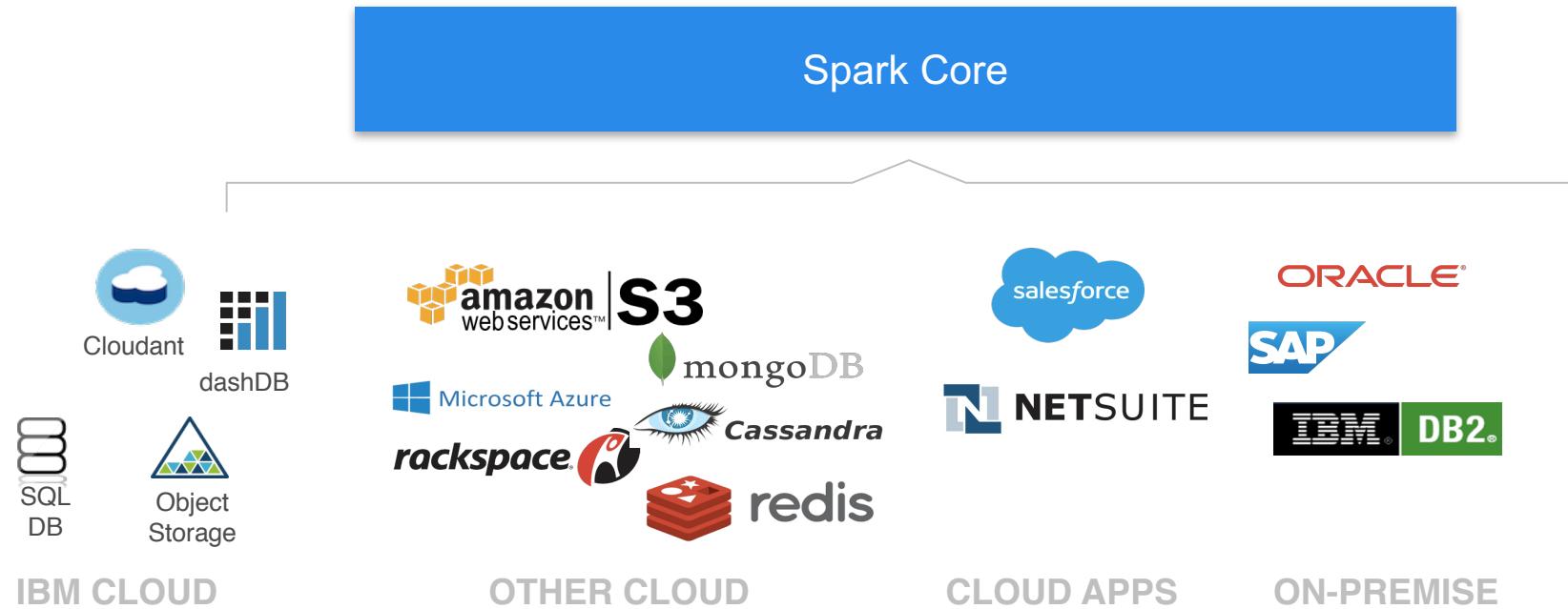
+57%

2015      2016  
1144      1800

source: DataBricks 2016 Survey

## Why Spark? Ease of Use

- Four APIs
- Documentation
- Data Platform Flexibility



## What Spark Is Not

- **Not only for Hadoop** – Spark can work with Hadoop (especially HDFS), but Spark is a standalone system
- **Not a data store** – Spark attaches to other data stores but does not provide its own
- **Not only for machine learning** – Spark includes machine learning and does it very well, but it can handle much broader tasks equally well
- **Not real time** – Spark Streaming is micro-batching, not true streaming, and cannot handle the real-time complex event processing



# Apache Spark Programming

# Spark Programming Languages

- **Scala**

- Functional programming
- Spark written in Scala
- Scala compiles into Java byte code

- **Python**

- Easily ports to Spark

- **Java**

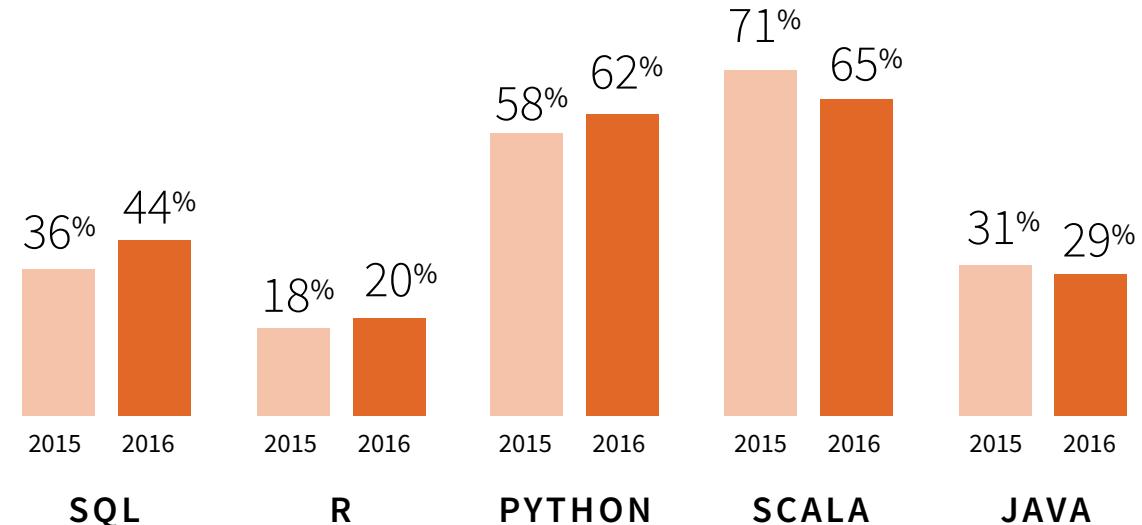
- New features in Java 8 makes for more compact coding

- **R**

- Statistical language used to create and manipulate data

LANGUAGES USED IN SPARK YEAR-OVER-YEAR

% of respondents who use each language (more than one language could be selected)



source: DataBricks 2016 Survey

# Functional Programming

- **Spark depends on functional programming**
- **Closures or in line anonymous functions**
  - Functions have inputs and outputs only
  - One can pass functions as inputs to other functions as parameters

## ▪ Python

–`function(lambda x: x)`

## ▪ Scala

–`function(x => x)`

# Functions!

- **Two types of operations**

- **Transformations**

- `map()`
    - `filter()`
    - `flatMap()`

- **Actions**

- `first()`
    - `collect()`
    - `take()`

- **Lazy evaluation**

## A Few Good Functions

**map** - Return a new RDD by applying a function to each element of this RDD.

**filter** - Return a new RDD containing only the elements that satisfy a predicate.

**reduce** - Reduces the elements of this RDD using the specified commutative and associative binary operator. Currently reduces partitions locally.

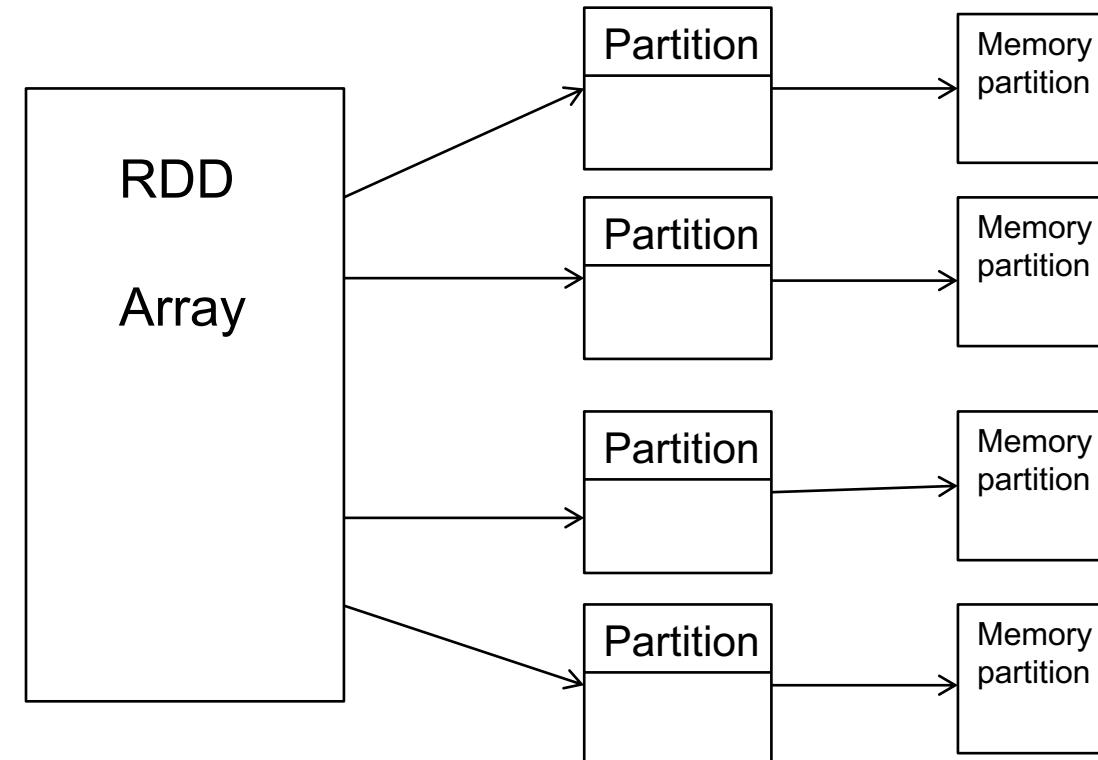
**flatMap** - Return a new RDD by first applying a function to all elements of this RDD, and then flattening the results.

**distinct** - Return a new dataset that contains the distinct elements of the source dataset.

**sortByKey** - When called on a dataset of (K, V) pairs where K implements Ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean **ascending** argument

# Resilient Distributed Dataset (RDD)

- **Immutable**
- **Holds references to partition objects**
- **Each partition is a subset of the overall data**
- **Partitions are assigned to nodes on the cluster**
- **Partitions are in memory by default**
- **Fault tolerance**
  - If data in memory is lost it will be recreated from lineage

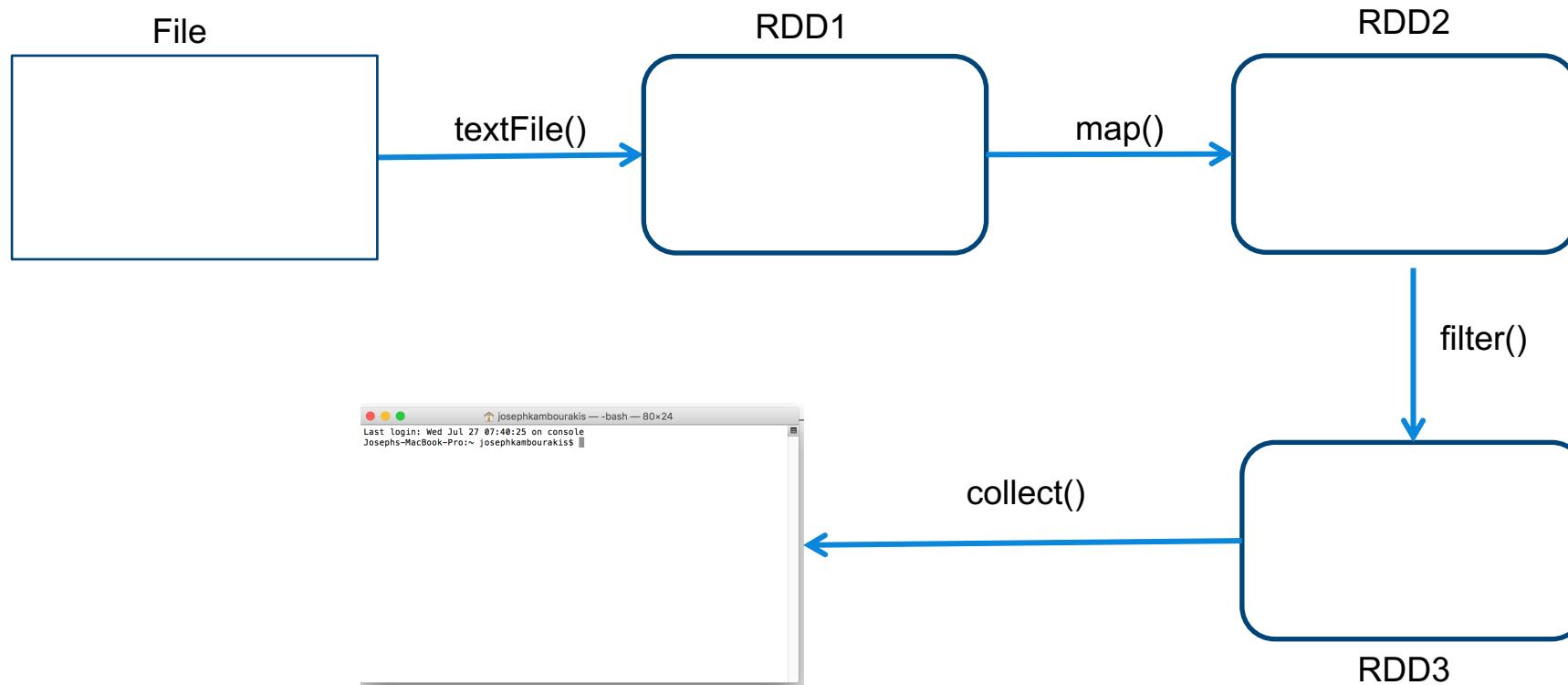


# Dataframes

- **Subset of RDDs**
- **Main data object abstraction in SparkSQL**
  - Organized into NAMED columns
- **Can be created from:**
  - Existed structured data source
    - json, database table, Apache Parquet
  - Another RDD
  - Transformation of another dataframe
  - Programmatically

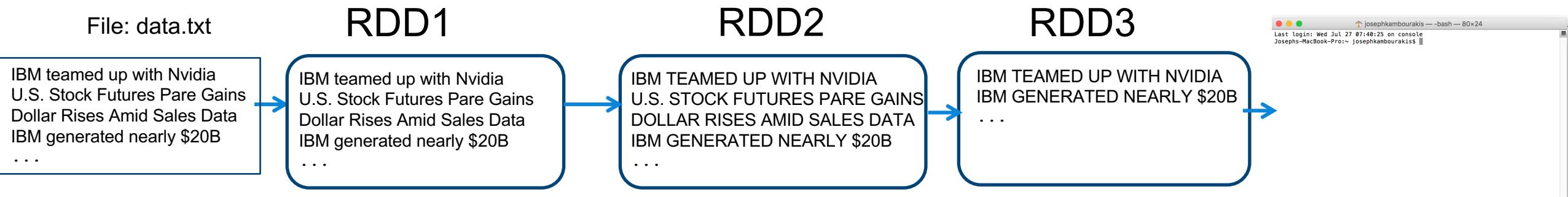
# Apache Spark Code Execution

# Directed Acyclic Graph - DAG



## Sample Code

```
▪ RDD1 = sc.textFile(data.txt)  
RDD2 = RDD1.map(lambda line: line.upper())  
RDD3 = RDD2.filter(lambda line: 'IBM' in line)  
RDD3.collect()
```

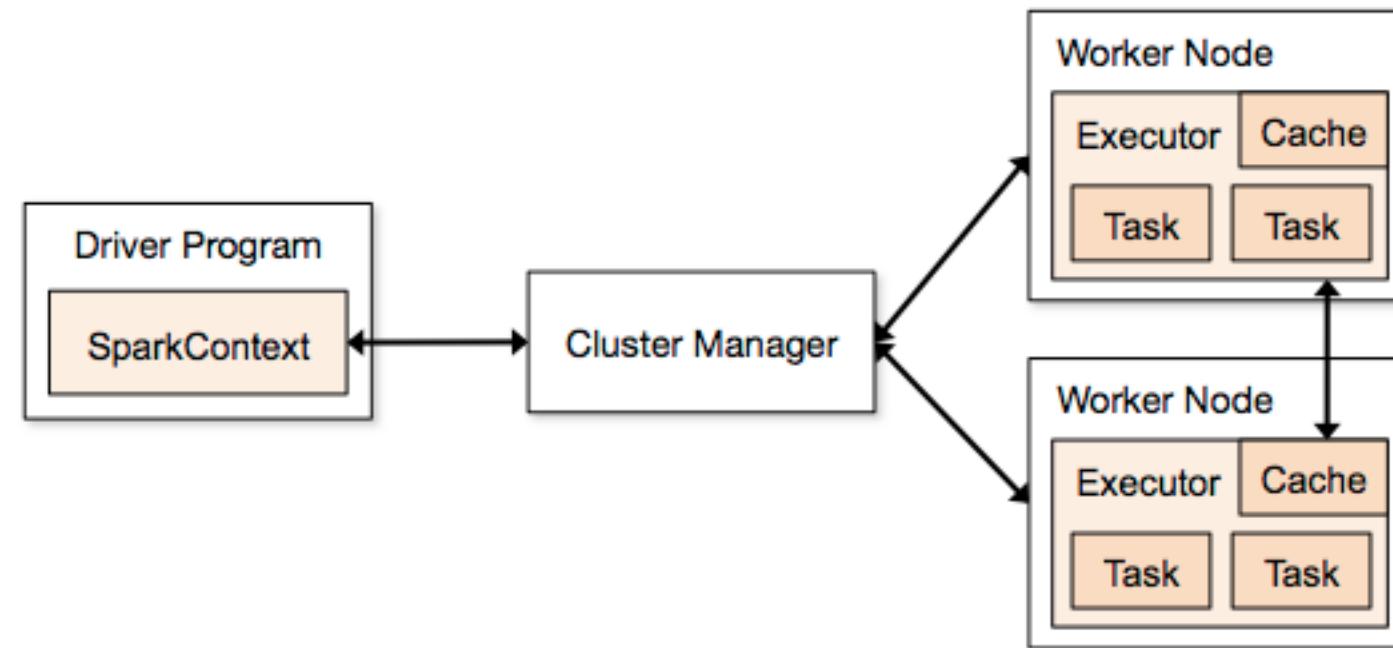


# Apache Spark Under the hood!

# Spark Application Architecture

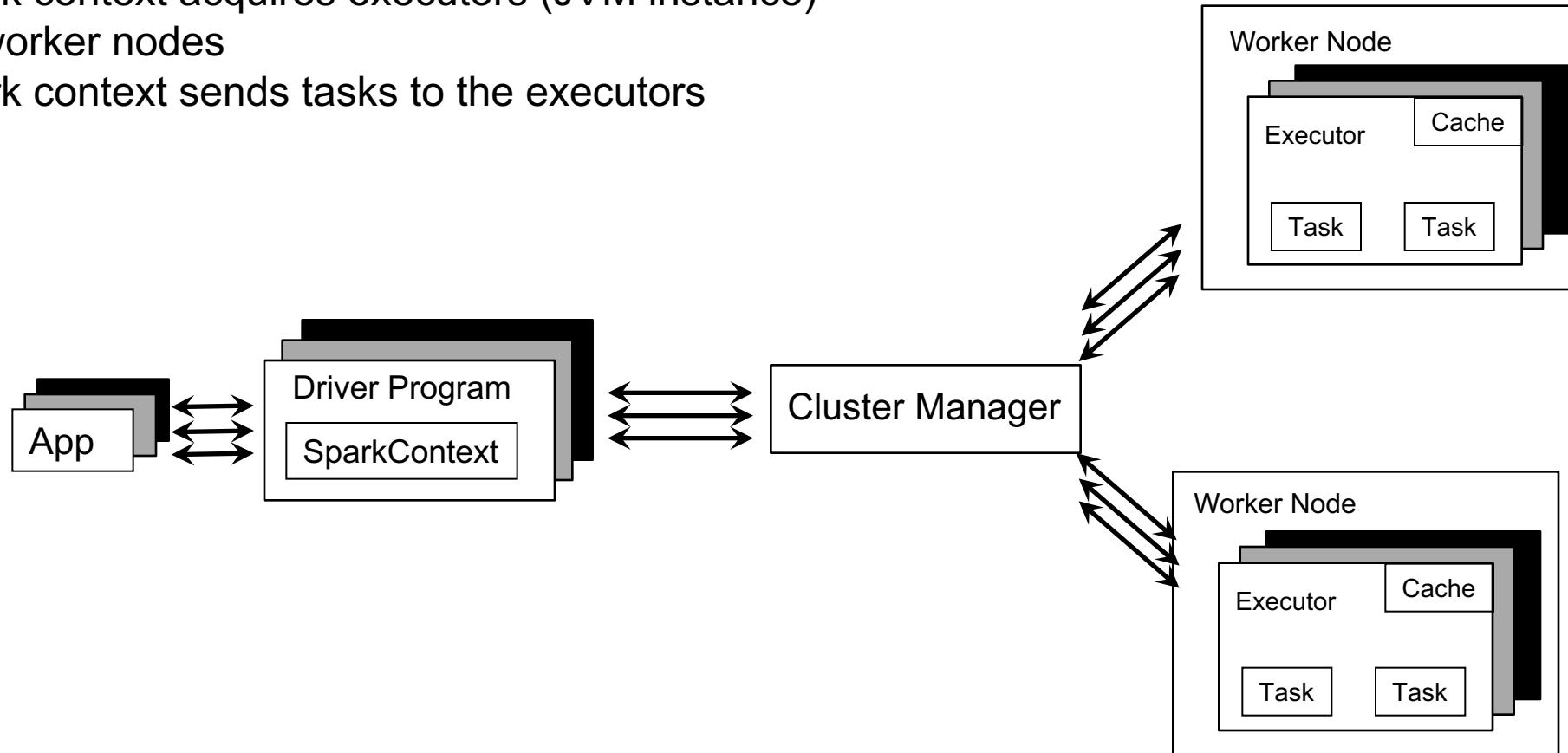


- A Spark application is initiated from a driver program
- Spark execution modes:
  - Standalone with the built-in cluster manager
  - Use Mesos as the cluster manager
  - Use YARN as the cluster manager
  - Standalone cluster on any cloud (BlueMix, IBM Softlayer, Amazon, Azure, ...)



# Showing multiple applications

- Each Spark application runs as a set of processes coordinated by the Spark context object (driver program)
  - Spark context connects to Cluster Manager (standalone, Mesos/Yarn)
  - Spark context acquires executors (JVM instance) on worker nodes
  - Spark context sends tasks to the executors



# Apache Spark Libraries

# SparkSQL

- **Component of Spark**
  - Project started in 2012
  - First alpha release in Spring 2013
  - Out of alpha with Spark 0.9.0
  - New features in Spark 2.0
- **Allows you to interact with Apache Hive**
- **More in next section!**

# Spark Streaming

- **Component of Spark**

- Project started in 2012
  - First alpha release in Spring 2013
  - Out of alpha with Spark 0.9.0
  - New features in Spark 2.0



- **Discretized Stream (DStream) programming abstraction**

- Represented as a sequence of RDDs (micro-batches)
  - RDD: set of records for a specific time interval
  - Supports Scala, Java, and Python (with limitations)

- **Fundamental architecture: batch processing of datasets**

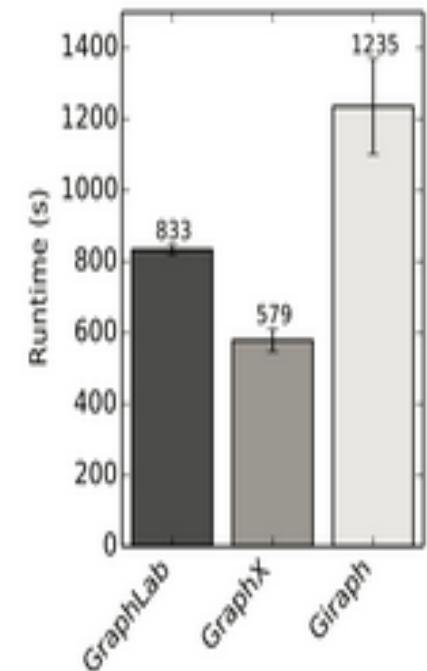


## Spark MLlib

- **Spark MLlib for machine learning library**
- **Provides common algorithm and utilities**
  - Classification
  - Regression
  - Clustering
  - Collaborative filtering
  - Dimensionality reduction
- **Leverages iteration and yields better results than one-pass approximations sometimes used with MapReduce**
- **More on this later!**

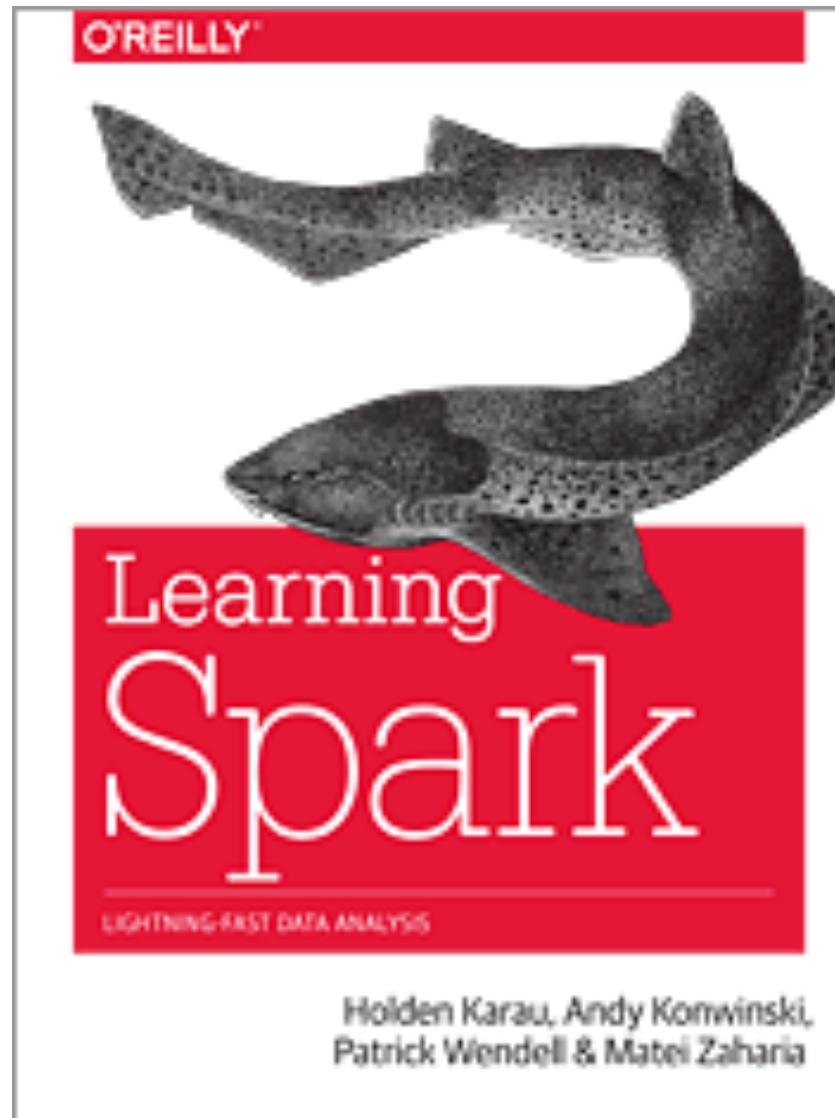
# Spark GraphX

- **Scala only!**
- **Flexible Graphing**
  - GraphX unifies ETL, exploratory analysis, and iterative graph computation
  - You can view the same data as both graphs and collections, transform and join graphs with RDDs efficiently, and write custom iterative graph algorithms with the API
  - GraphFrames - graph library based on Data Frames
- **Speed**
  - Comparable performance to the fastest specialized graph processing systems.
- **Algorithms**
  - Choose from a growing library of graph algorithms
  - In addition to a highly flexible API, GraphX comes with a variety of graph algorithms



# Apache Further Resources

## Additional Resources



# Additional Resources

The screenshot shows the Apache Spark website on a web browser. The header includes the Apache logo and the text "Lightning-fast cluster computing". The main navigation bar has links for Download, Libraries, Documentation, Examples, Community, and FAQ, along with a link to the Apache Software Foundation. A sub-navigation bar on the left lists Speed, Ease of Use, and Built-in Libraries. The main content area features a section about Speed with a bar chart comparing Hadoop and Spark's running times for logistic regression, and another section about Ease of Use with sample code for reading a file and performing operations like flatMap, map, and reduceByKey.

Apache Spark™ - Lightning-fast cluster computing

APACHE Spark™ Lightning-fast cluster computing

Download Libraries Documentation Examples Community FAQ Apache Software Foundation

Apache Spark™ is a fast and general engine for large-scale data processing.

## Speed

Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

Apache Spark has an advanced DAG execution engine that supports cyclic data flow and in-memory computing.

## Ease of Use

Write applications quickly in Java, Scala, Python, R.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it interactively from the Scala, Python and R shells.

Running time (s)

System	Running time (s)
Hadoop	110
Spark	0.9

Logistic regression in Hadoop and Spark

```
text_file = spark.textFile("hdfs://...")  
text_file.flatMap(lambda line: line.split())  
    .map(lambda word: (word, 1))  
    .reduceByKey(lambda a, b: a+b)
```

## Latest News

Spark 2.0.0 released (Jul 26, 2016)  
Spark 1.6.2 released (Jun 25, 2016)  
Call for Presentations for Spark Summit EU is Open (Jun 16, 2016)  
Preview release of Spark 2.0 (May 26, 2016)

Archive

## Download Spark

## Built-in Libraries:

SQL and DataFrames  
Spark Streaming  
MLlib (machine learning)  
GraphX (graph)  
Third-Party Packages

# Notebook Intro

- **Data Science Experience**
  - Jupyter Notebook
- **Student Code**
  - Notebook and solutions will be provided after the lab and we'll walk through them.

---

## Small Demo

- **datascience.ibm.com**