

Hands on Introduction to Spark



Power of data. Simplicity of
design. Speed of innovation.

Bernie Beekman
Asad Mahmood
Michael Cronk
Brett Coffman

Agenda

9:00 - 10:00 - Kick off

Introduction to Spark
Introduction to the Data Science Experience
Register for Platform – datascience.ibm.com

10:00 - 11:30 - Lab 1 – Introduction to Spark

Overview of Lab
Hands on Exercises

11:30 – 1:45 Lab 2 – Spark SQL

Overview of Lab
Lunch/Hands on Exercises

1:45 - 2:00 – Break

2:00 - 3:45 - Lab 3 – Spark Machine Learning

Overview of Lab
Hands on Exercises

3:45 - 4:00 – Break

4:00 - 4:30 – Questions/Wrap Up

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”

Why Spark?

Performant



- In-memory architecture greatly reduces disk I/O
- Improved performance compared to MapReduce

Productive



- **Concise and expressive syntax**, especially compared to prior approaches
- **Single programming model** across a range of use cases and steps in data lifecycle
- **Integrated with common programming languages** – Java, Python, Scala, R
- **New tools** continually reduce skill barrier for access (e.g. SQL for analysts)

Leverages existing investments



- Works well within **existing Hadoop ecosystem**

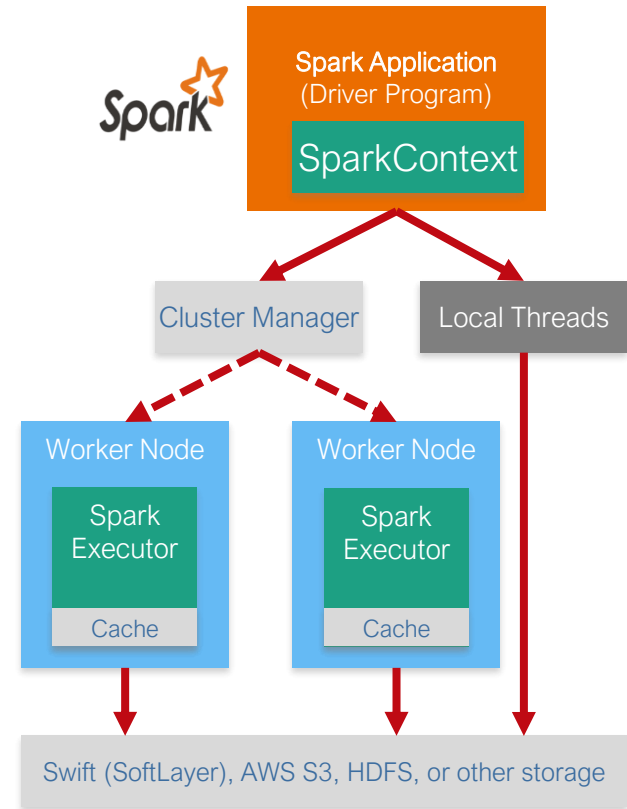
Improves with age



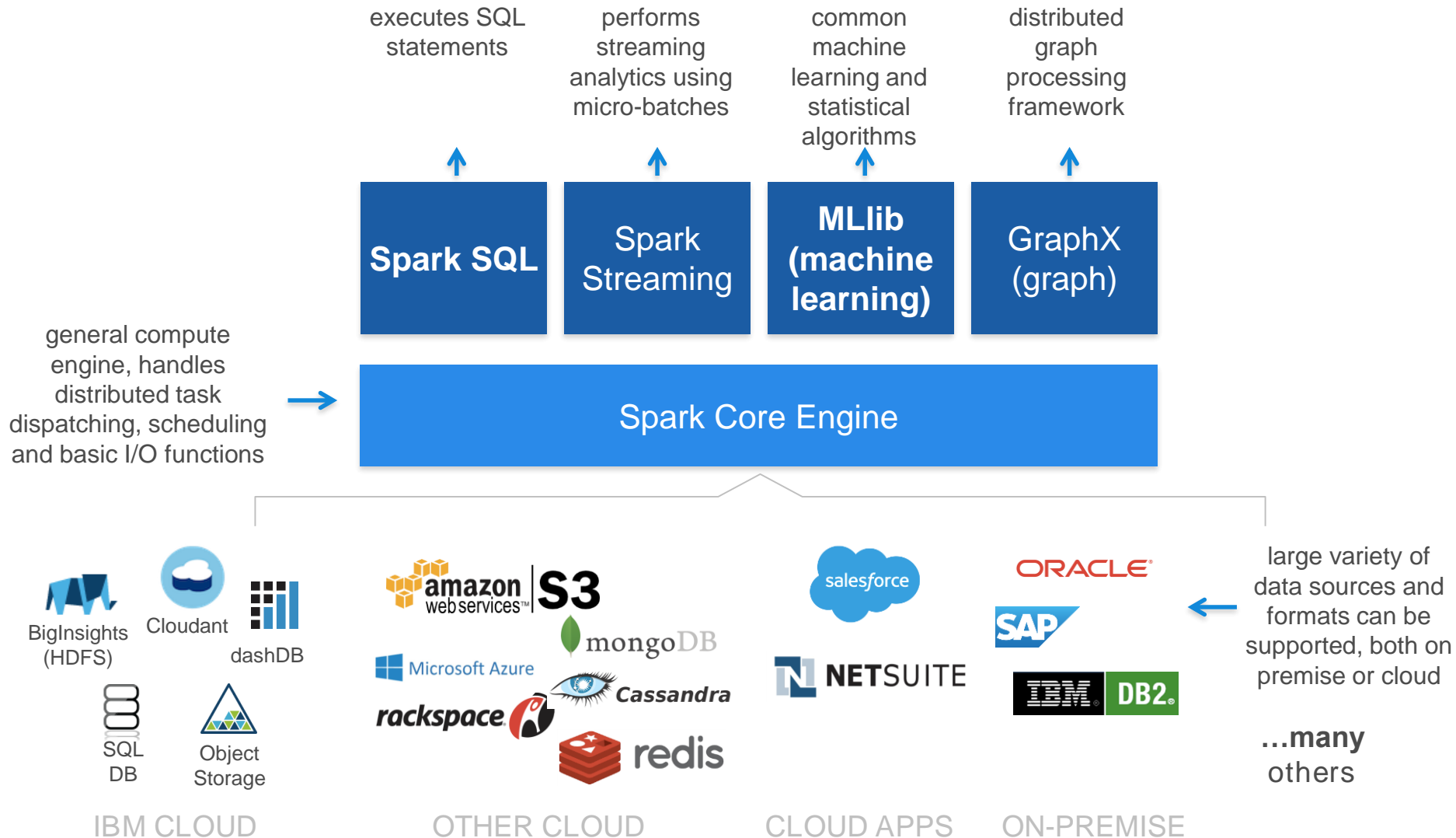
- **Large and growing community** of contributors continuously improve full analytics stack and extend capabilities

Apache Spark Scales

- Spark programs generally consist of two components: **driver program** and **worker program(s)**
 - **Driver Program** manages the division of computations (Task) that are sent to worker nodes
 - **Worker programs** run smaller portions of computations
- The **SparkContext** object instructs Spark on how & where to access a cluster
- **Cluster Manager** manages the physical resources needed to run driver and worker programs.



Spark blends multiple data types, sources, and workloads



Spark Programming Languages

▪ Scala

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

▪ Java

- New features in Java 8 makes for more compact coding (lambda expressions)

▪ Python

- Most widely used API with Spark today

▪ R

- Open source language and programming environment to support statistical analysis and visualization

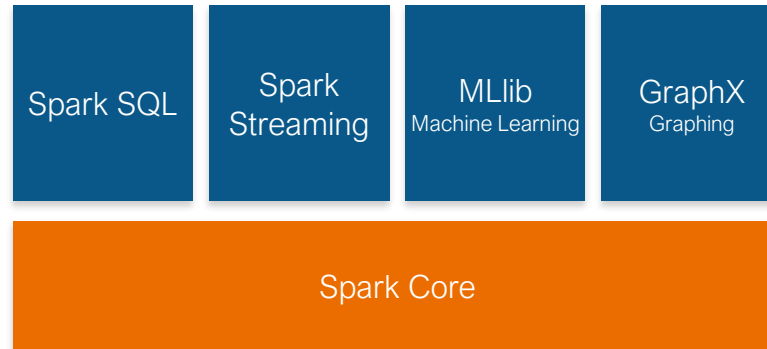
| Language | 2014 | 2015 | 2016 |
|----------|---------|------|------|
| Scala | 84% | 71% | 65% |
| Java | 38% | 31% | 29% |
| Python | 38% | 58% | 62% |
| R | unknown | 18% | 20% |

Survey done by Databricks,
Summer 2015, 2016

**This probably means that more “data scientists” are starting to use Spark
DataFrames make all languages equally performant**

Benefits of Spark for Data Science

- General compute engine
- Basic I/O functions
- Task dispatching
- Scheduling



Spark is Easy...

- Support multiple programming interfaces (Scala, Python, Java and R)
- Less lines of code to get answers.

Spark is Agile...

- Unified APIs (SQL, DataFrames, Streaming, Machine Learning, etc.)
- Supports Notebooks (Jupyter, Zeppelin, Etc.)

Spark is Fast...

- In-Memory processing that scales in a distributed architecture.
- Application follows lazy evaluations architecture w/ optimized execution.

IBM is all-in on Spark

Contribute to the Core

Launch Spark Technology Center (STC), 300 engineers

Open source SystemML

Partner with Databricks

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

Ben Horowitz,
Andreessen Horowitz

Foster Community

Educate 1M+ data scientists and engineers via online courses

Sponsor AMPLab, creators and evangelists of Spark

Infuse the Portfolio

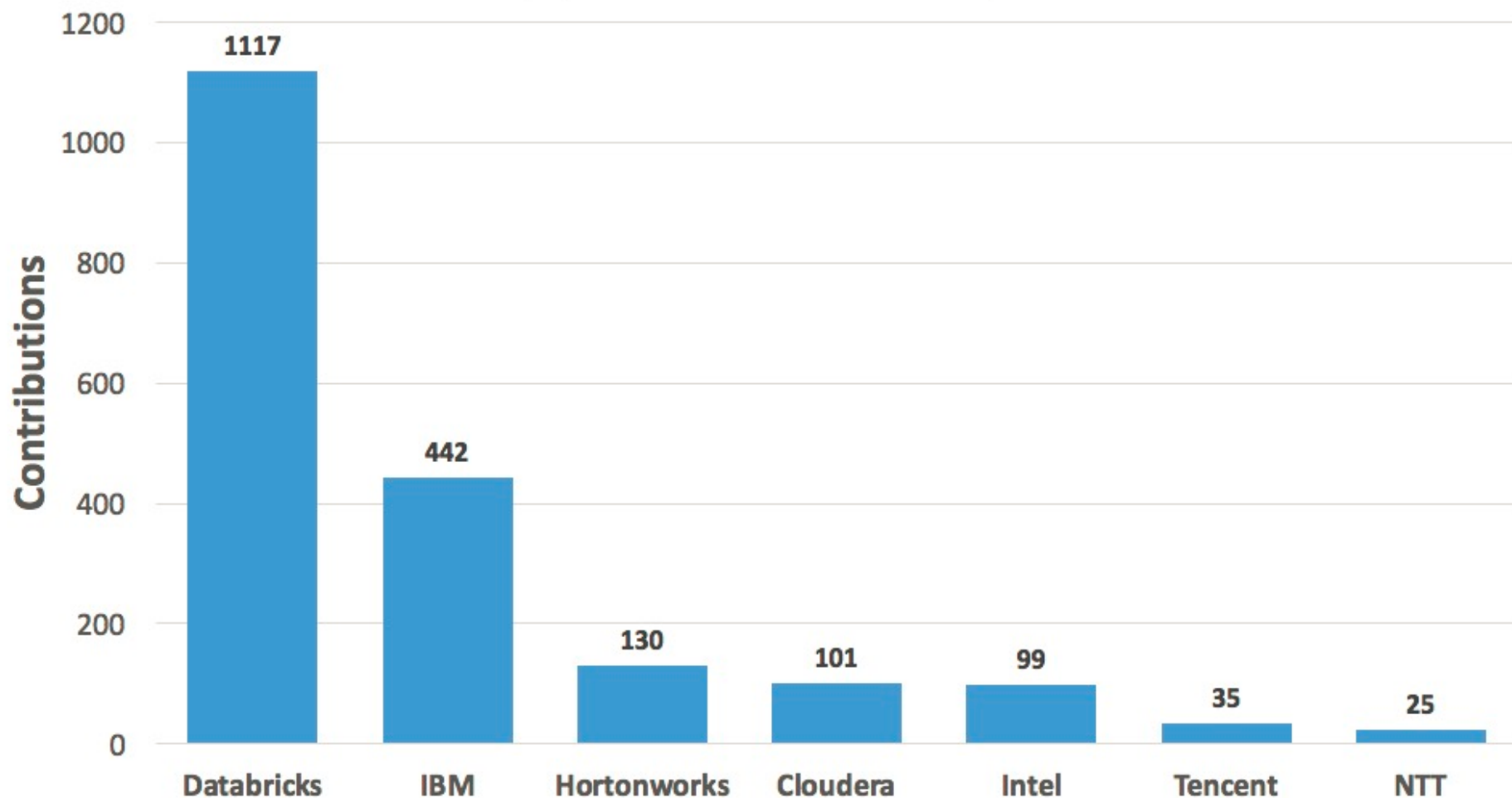
Integrate Spark throughout portfolio

3,500 employees working on Spark-related topics

Spark however customers want it – standalone, platform or products

Top 7 Contributing Companies to Spark 2.0.0

(Represents 70% of total contributions)



Apache Spark Under the hood!

Spark Terminology

- **Spark Context**

- Represents a connection to the Spark cluster. The Application which initiated the context can submit one or several jobs, sequentially or in parallel, batch or interactively.

- **Driver Process**

- The program or process running the Spark context. Responsible for running jobs over the cluster and converting the App into a set of tasks

- **Job**

- A piece of logic (code) which will take some input, perform some transformations and an action (computes a result, writes output)

- **Stage**

- Jobs are divided into stages

- **Tasks**

- Each stage is made up of tasks. One task per partition. One task is executed on one partition (of data) by one executor

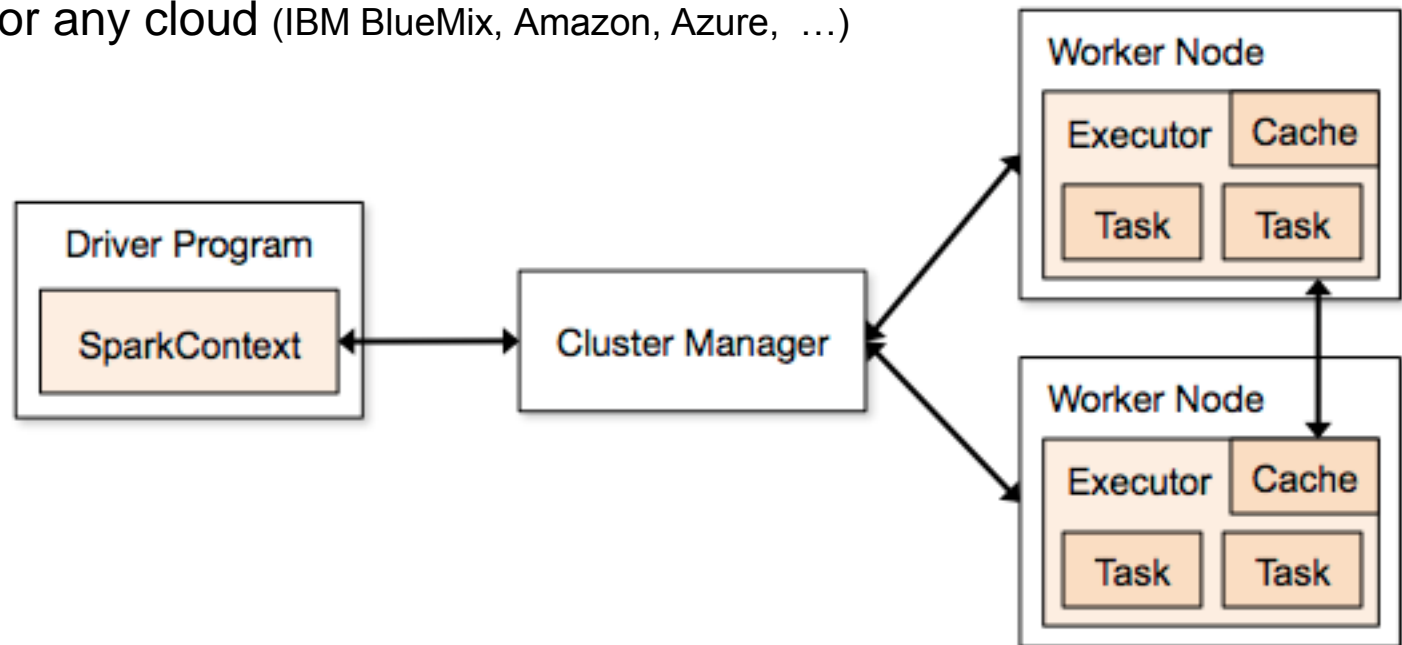
- **Executor**

- The process responsible for executing a task on a worker node

Spark Application Architecture

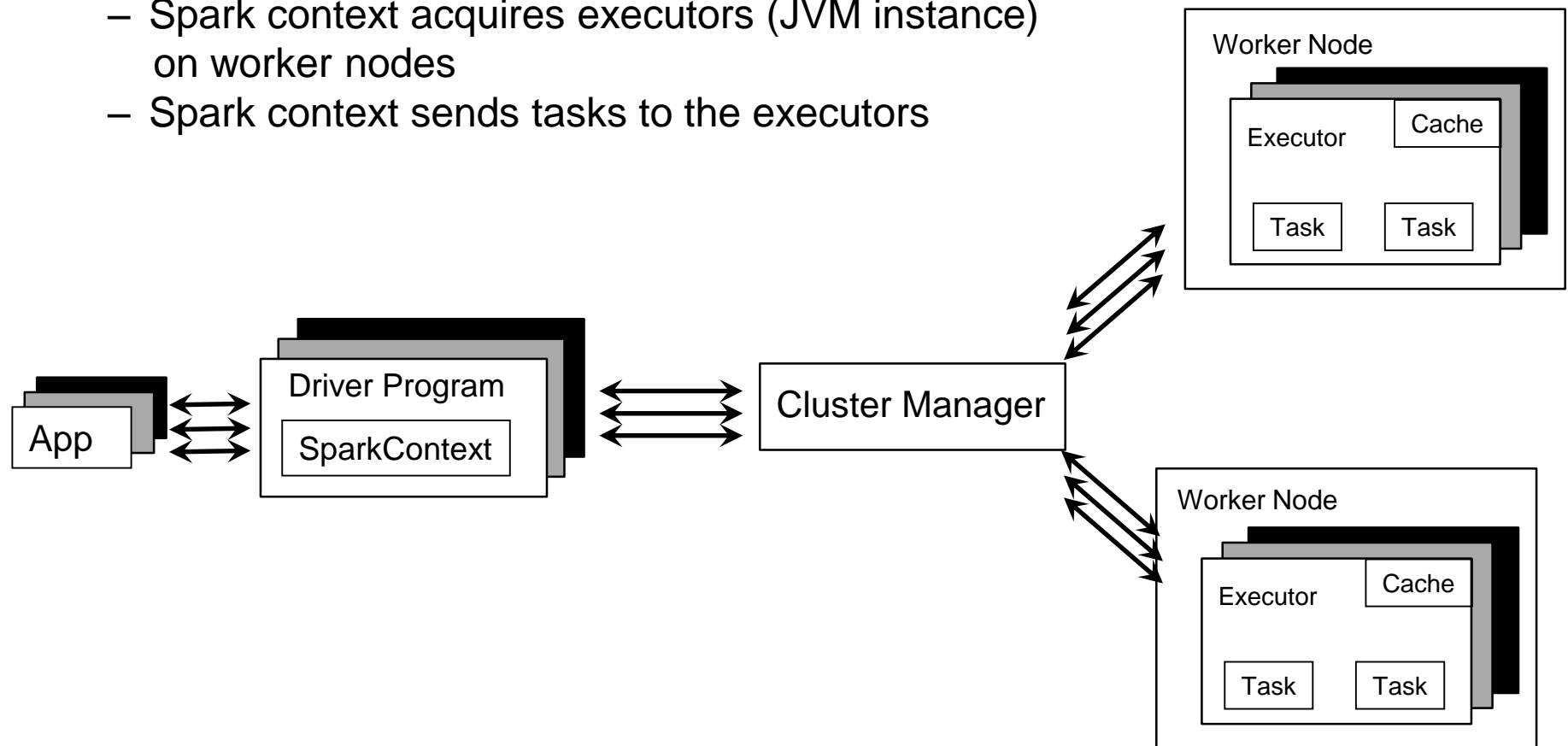


- A Spark application is initiated from a driver program
- **Spark execution modes:**
 - Standalone with the built-in cluster manager
 - Use Mesos as the cluster manager
 - Use YARN as the cluster manager
 - Kubernetes (experimental)
 - On-premise or any cloud (IBM BlueMix, Amazon, Azure, ...)



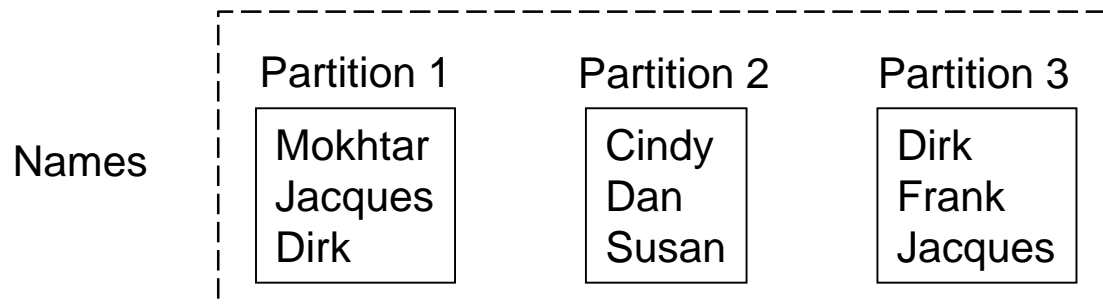
Showing multiple applications

- 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



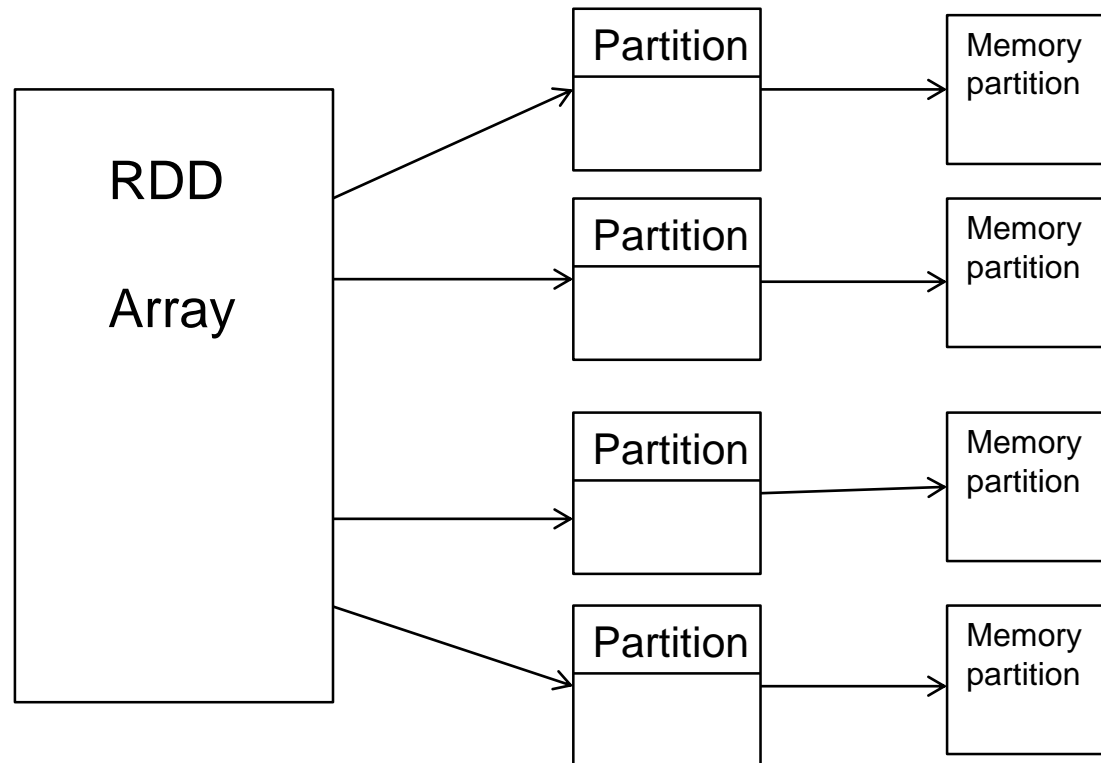
Resilient Distributed Datasets

- **An RDD is a distributed collection of Scala/Python/Java/R objects of the same type:**
 - RDD of strings
 - RDD of integers
 - RDD of (key, value) pairs
 - RDD of class Java/Python/Scala/R objects
- **An RDD is physically distributed across the cluster, but manipulated as one logical entity:**
 - Spark will “distribute” any required processing to all partitions where the RDD exists and perform necessary redistributions and aggregations as well.
 - Example: Consider a distributed RDD “Names” made of names



Resilient Distributed Dataset

- RDDs are **immutable**
 - Modifications create new RDDs
- 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
- RDDs keep information on their lineage
 - Fault tolerant
 - If data in memory is lost it will be recreated



Resilient Distributed Datasets

▪ Two types of operations

– Transformations

- `val rddNumbers = sc.parallelize(1 to 10): Numbers from 1 to 10`
- `val rddNumbers2 = rddNumbers.map (x => x+1): Numbers from 2 to 11`
- LINEAGE on how to obtain rddNumbers2 from rddNumber is recorded
- It's a Directed Acyclic Graph (DAG)
- No actual data processing does take place → **Lazy evaluations**

– Actions

- `rddNumbers2.collect(): Array [2, 3, 4, 5, 6, 7, 8, 9, 10, 11]`
- Performs list of transformations and THE action
- Returns a value (or write to a file)

▪ Fault tolerance

- If data in memory is lost it will be recreated from lineage

Code Execution (1)

- 'spark-shell' provides Spark context as 'sc'

```
// Create RDD
quotes = sc.textFile("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.filter(lambda x:x.startsWith("DAN"))
danSpark = danQuotes.map(lambda x:x.split(" ")).map(lambda x: x[1])
// Action
danSpark.filter(lambda x: "Spark" in x).count(1)
```

File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

Code Execution (2)

```
// Create RDD
quotes = sc.textFile("hdfs:/sparkdata/sparkQuotes.txt")

// Transformations
danQuotes = quotes.filter(lambda x:x.startsWith("DAN"))
danSpark = danQuotes.map(lambda x:x.split(" ")).map(lambda x: x[1])

// Action
danSpark.filter(lambda x: "Spark" in x).count(1)
```

File: sparkQuotes.txt

RDD: quotes

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```



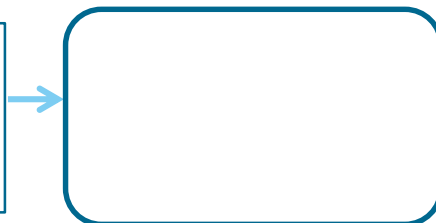
Code Execution (3)

```
// Create RDD
quotes = sc.textFile("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.filter(lambda x:x.startsWith("DAN"))
danSpark = danQuotes.map(lambda x:x.split(" ")).map(lambda x: x[1])
// Action
danSpark.filter(lambda x: "Spark" in x).count(1)
```

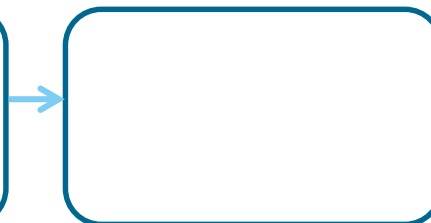
File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

RDD: quotes



RDD: danQuotes



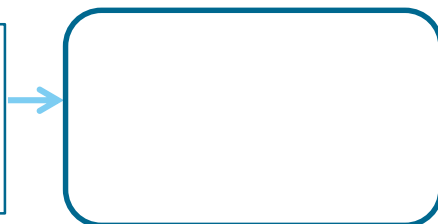
Code Execution (4)

```
// Create RDD
quotes = sc.textFile("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.filter(lambda x:x.startsWith("DAN"))
danSpark = danQuotes.map(lambda x:x.split(" ")).map(lambda x: x[1])
// Action
danSpark.filter(lambda x: "Spark" in x).count(1)
```

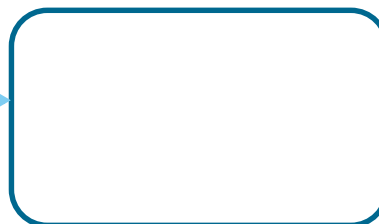
File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

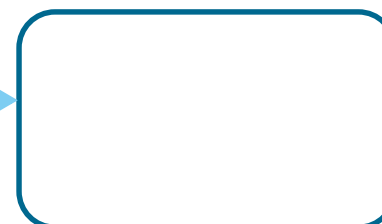
RDD: quotes



RDD: danQuotes

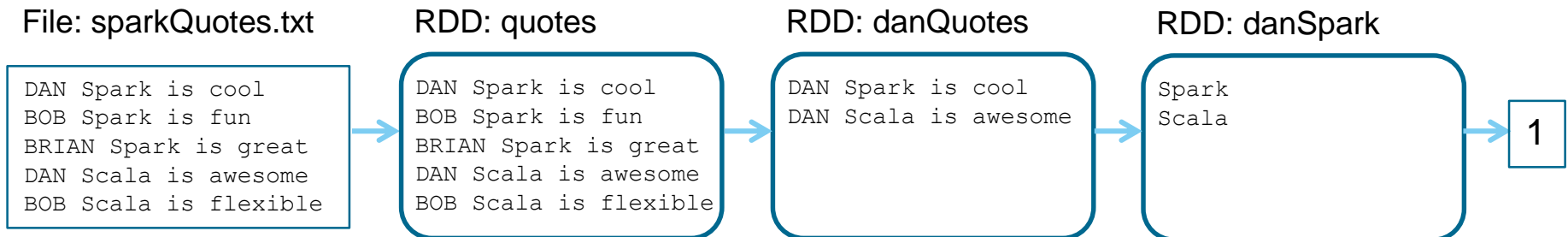


RDD: danSpark

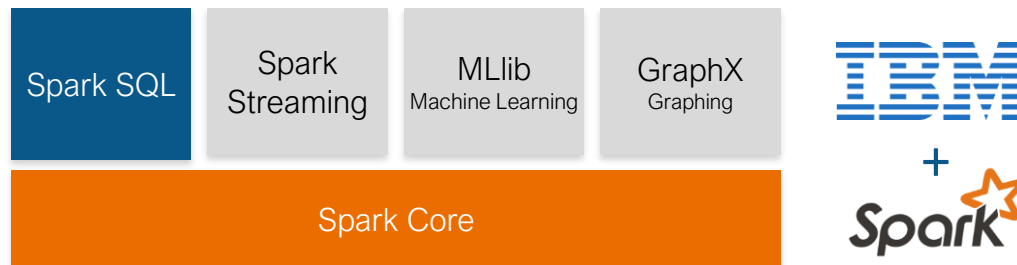


Code Execution (5)

```
// Create RDD
quotes = sc.textFile("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.filter(lambda x:x.startsWith("DAN"))
danSpark = danQuotes.map(lambda x:x.split(" ")).map(lambda x: x[1])
// Action
danSpark.filter(lambda x: "Spark" in x).count()
```



Closer Look at APIs – Spark SQL



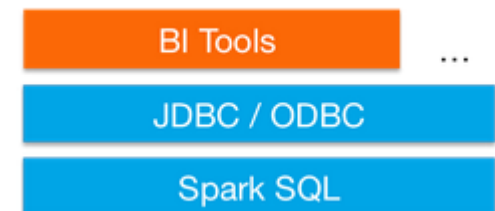
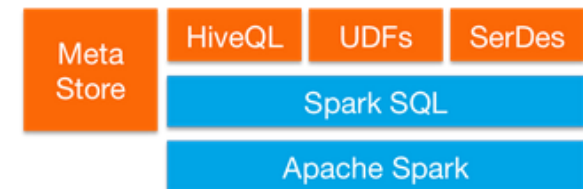
- Unified data access: Query structured data sets with SQL or DataFrame/Dataset APIs
- Fast, familiar query language across all of your enterprise data
- Use BI tools to connect and query via JDBC or ODBC drivers

Spark DataFrames & Datasets

- **DataFrame API announce in February 2015**
- **Dataset: Generalization of DataFrame**
 - Available in Scala and Java (DataFrame for Python and R)
- **Distributed collection of data organized in named columns**
 - Conceptually equivalent to a relational table, R/Python data frames
- **Supported format and sources**
 - From sources such as: JSON, Hive, JDBC, parquet, csv, etc.
- **Benefits:**
 - Easier manipulation interface (similar to SQL)
 - Higher abstraction for possible optimization
 - Unified interface for working with structured data

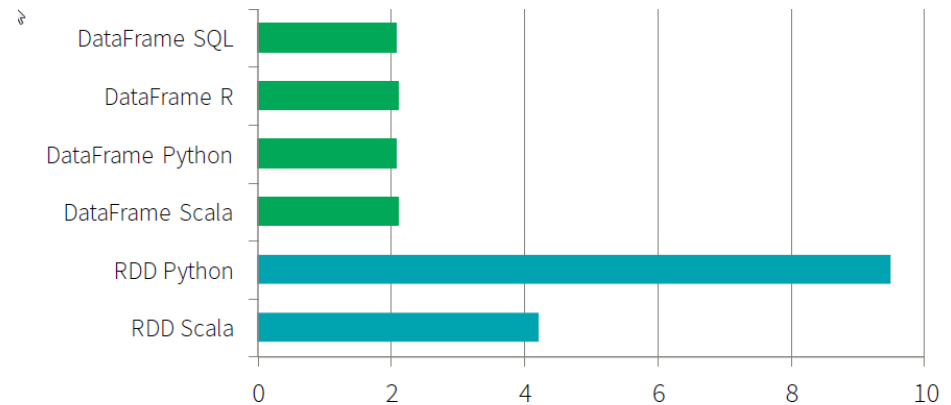
SparkSQL

- Provide for relational queries expressed in SQL, HiveQL using Scala, Python, and Java API's
- Seamlessly mix SQL queries with Spark programs
- Leverages Hive frontend and metastore
 - Compatibility with Hive data, queries, and UDFs
- Graduated from alpha status with Spark 1.3
 - DataFrames API marked as experimental in 2013
- Standard connectivity through JDBC/ODBC



SparkSQL, DataFrames and DataSets

- **A rich set of functionality that allows “Database-like” processing**
- **Share single optimizer, called “Catalyst”** (at the driver)
 - An open-source extensible query optimizer
- **Because it is the same engine, it has exactly the same performance for different APIs**
 - And performance is much better than for RDD
- **Much less code**
- **All SparkSQL, DF, and DataSets are essentially using the same engine**



Time to aggregate 10 million integer pairs (in seconds)

Picture credit: databricks.com

Code Execution (1)

- 'spark-shell' provides SparkSession as 'spark'

```
// Create DataFrame
quotes = spark.read.text("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.where(expr("value").startswith("DAN"))
danSpark = danQuotes.select(split(expr("value"), " ").alias("value")).select(expr("value")[1].alias("value"))
// Action
danSpark.where(expr("value").like("%Spark%")).count()
```

File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

Code Execution (2)

```
// Create DataFrame
quotes = spark.read.text("hdfs:/sparkdata/sparkQuotes.txt")

// Transformations
danQuotes = quotes.where(expr("value").startswith("DAN"))
danSpark = danQuotes.select(split(expr("value"), " ").alias("value")).select(expr("value")[1].alias("value"))

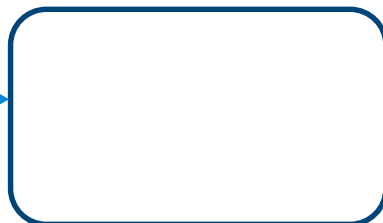
// Action
danSpark.where(expr("value").like("%Spark%")).count()
```

File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```



DataFrame: quotes



Code Execution (3)

```
// Create DataFrame
quotes = spark.read.text("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.where(expr("value").startswith("DAN"))
danSpark = danQuotes.select(split(expr("value"), " ").alias("value")).select(expr("value")[1].alias("value"))
// Action
danSpark.where(expr("value").like('%Spark%')).count()
```

File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

DataFrame: quotes



DataFrame: danQuotes



Code Execution (4)

```
// Create DataFrame
quotes = spark.read.text("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.where(expr("value").startswith("DAN"))
danSpark = danQuotes.select(split(expr("value"), " ")
).alias("value")).select(expr("value")[1].alias("value"))
// Action
danSpark.where(expr("value").like('%Spark%')).count()
```

File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

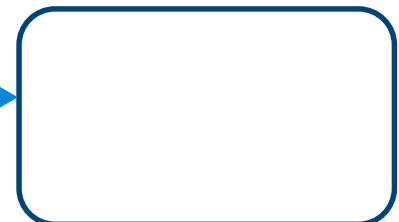
DataFrame: quotes



DataFrame: danQuotes



DataFrame: danSpark



Code Execution (5)

```
// Create DataFrame
quotes = spark.read.text("hdfs:/sparkdata/sparkQuotes.txt")
// Transformations
danQuotes = quotes.where(expr("value").startswith("DAN"))
danSpark = danQuotes.select(split(expr("value"), " ").alias("value")).select(expr("value")[1].alias("value"))
// Action
danSpark.where(expr("value").like('%Spark%')).count()
```

File: sparkQuotes.txt

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

DataFrame: quotes

```
DAN Spark is cool
BOB Spark is fun
BRIAN Spark is great
DAN Scala is awesome
BOB Scala is flexible
```

DataFrame: danQuotes

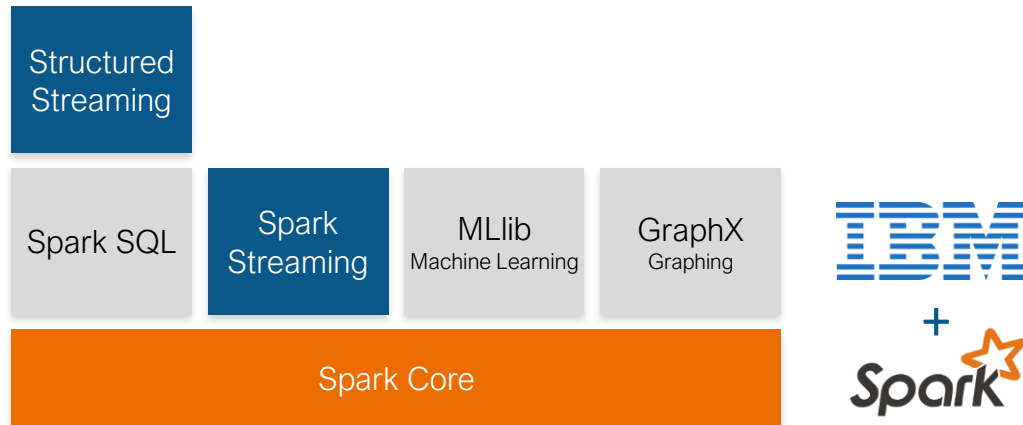
```
DAN Spark is cool
DAN Scala is awesome
```

DataFrame: danSpark

```
Spark
Scala
```

1

Closer Look at APIs - Streaming



- Micro-batch event processing for near-real time analytics
- e.g. Internet of Things (IoT) devices, Twitter feeds, Kafka (event hub), etc.
- Spark's engine drives some action or outputs data in batches to various data stores
- No multi-threading or parallel process programming required

Spark Streaming



- **Component of Spark**
 - Project started in 2012
 - First alpha release in Spring 2013
 - Out of alpha with Spark 0.9.0
 - More enhancements targeted for 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**



Structured Streaming



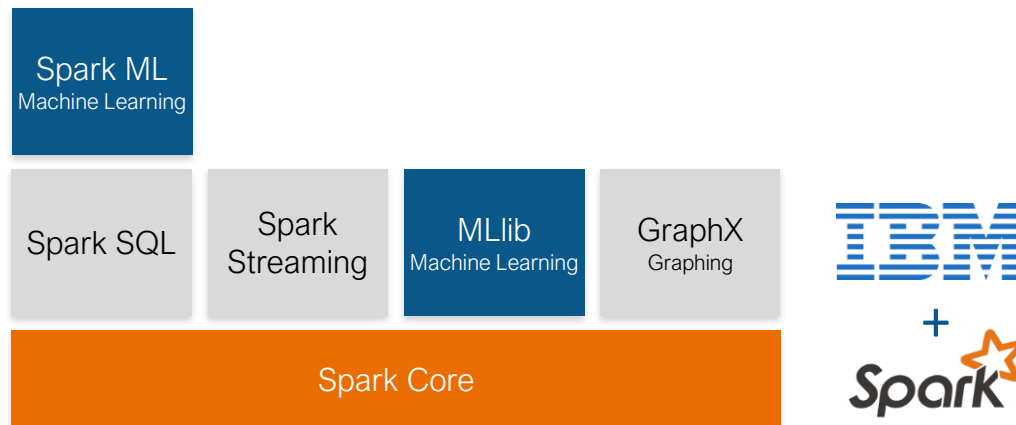
- **Component of Spark**

- Alpha - Spark 2.0
- Stabilized Spark 2.2

- **Advantages**

- Built on Spark SQL Engine
- Enables consistent API on streams as on batch processing
- Event Time Processing is handled
- Higher Level interface than with DStreams
- Easier to build end-to-end continuous applications

Closer Look at APIs – Machine Learning



- Feature Engineering
- Machine learning algorithms for:
 - Clustering
 - Classification
 - Regression
 - Recommendation
 - Etc.

Spark Machine Learning

- **Spark MLlib/ML for machine learning**

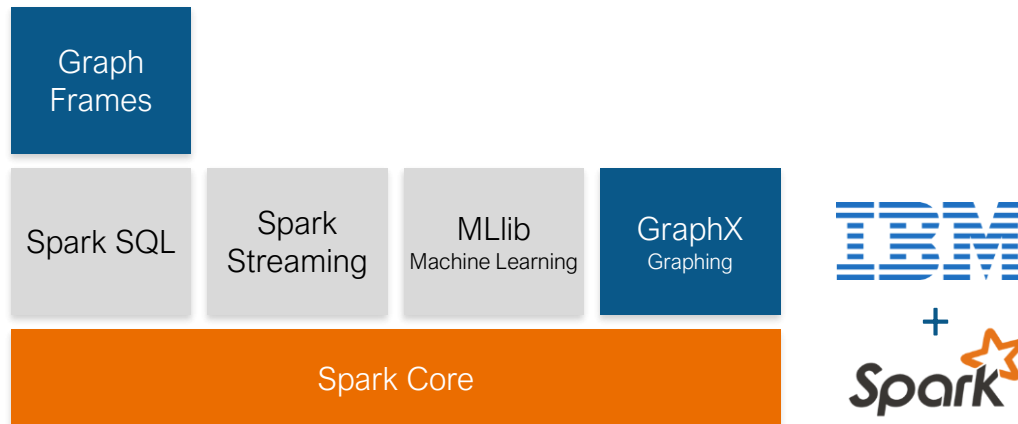
- RDD-based package `spark.mllib` now in maintenance mode
- The primary API is now the DataFrame-based package `spark.ml`

–

- **Provides common algorithm and utilities**

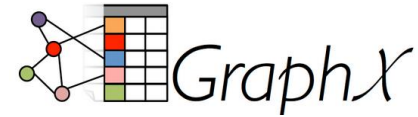
- Classification
- Regression
- Clustering
- Collaborative filtering
- Dimensionality reduction
- Etc.

Closer Look at APIs – Graph



- Represent and analyze systems represented by graph nodes and edges
- Trace interconnections between graph nodes
- Applicable to use cases in transportation, telecommunications, road networks, modeling personal relationships, social media, etc.

Spark GraphX, GraphFrames



▪ 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 – higher level interface built on top of GraphX to provide a DataFrame API

▪ 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 algorithms

