

# Apache Spark: Paralelismo y Big Data

Antonio J. Nebro

- Dr. Ingeniero en Informática
- Grupo de Investigación Khaos
- Depto. Lenguajes y Ciencias de la Computación
- Universidad de Málaga

# Contenidos del taller

- Introducción
- ¿Qué es Apache Spark?
- Spark y paralelismo
- Programación en Spark
- Conclusiones del taller

# Introducción



- **Ponente**
  - Antonio Jesús Nebra Urbaneja
    - Doctor Ingeniero en Informática
    - Inicio de estudios en Informática: 1984
  - Afiliación
    - Departamento de Lenguajes y Ciencias de la Computación
      - E.T.S. De Ingeniería Informática
    - Grupo de investigación Khaos
      - Edificio de investigación Ada Byron
  - Contacto
    - [antonio@lcc.uma.es](mailto:antonio@lcc.uma.es)
    - @AJNebro



# Introducción



## Paralelismo y Big Data

- Investigación (desde 1990)
  - Paralelismo
  - Optimización
- Software libre
  - Framework de optimización jMetal
  - <http://jmetal.sourceforge.net>
  - <https://github.com/jMetal/jMetal>

- Máster en Ingeniería Web (RIAtec)
  - <http://riatec.lcc.uma.es>



- Máster en Big Data Analytics (BDA)
  - <http://bigdata.lcc.uma.es>
- Curso de extensión universitaria “Introducción a Big Data y Hadoop”
  - <http://bigdata.lcc.uma.es/curso/curso.html>

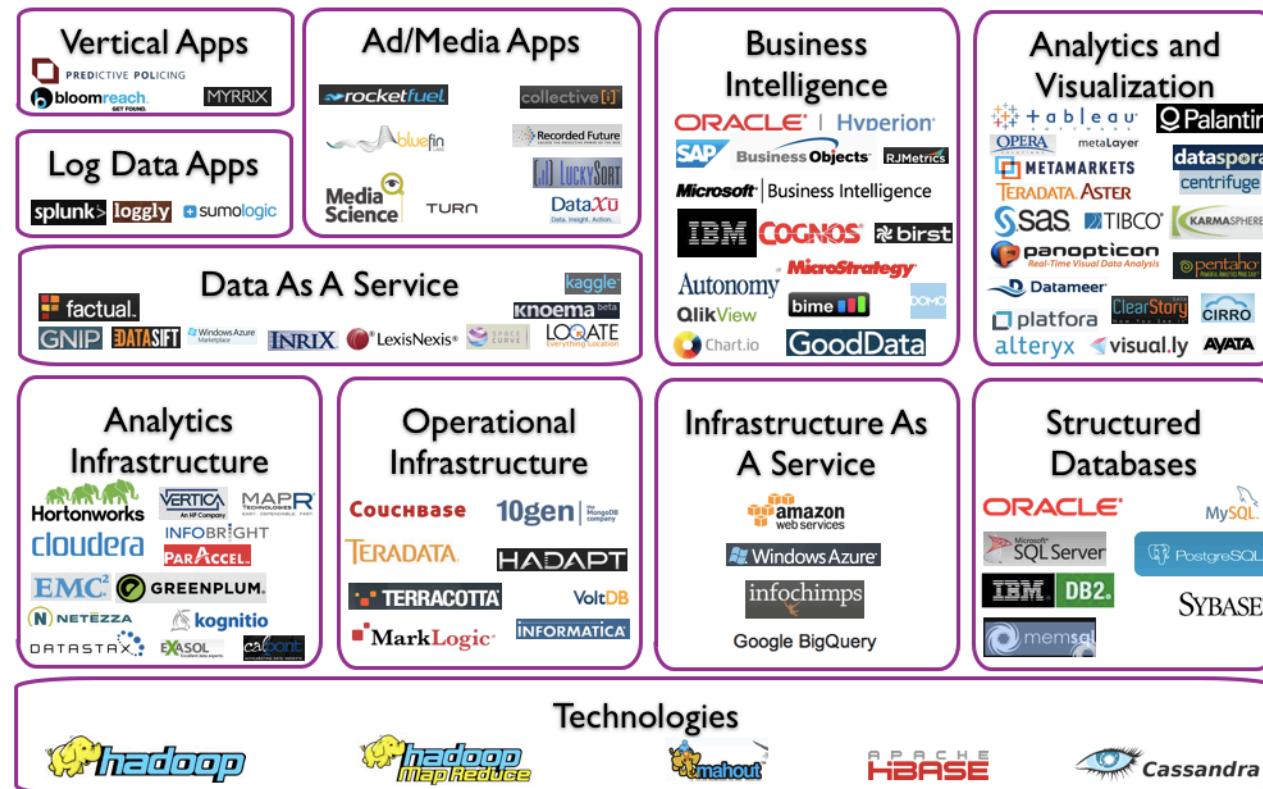


- El paralelismo intenta sacar partido a computadores paralelos para
  - Ejecutar programas más rápidamente
  - Ejecutar programas con altos requisitos de CPU, memoria, almacenamiento
- ¿Qué se estudia en las titulaciones de Ingeniería Informática?
  - Fundamentos: secciones críticas, interbloqueos, ...
  - Soluciones de bajo nivel
    - Memoria compartida: semáforos, monitores
    - Memoria distribuida: paso de mensajes
  - ¿Qué tienen en común?
    - Complejidad, dificultad de reutilizar código

- Big Data
  - Procesar grandes cantidades de información cuando las tecnologías tradicionales de BBDD no son suficientes



# Big Data Landscape



Copyright © 2012 Dave Feinleib

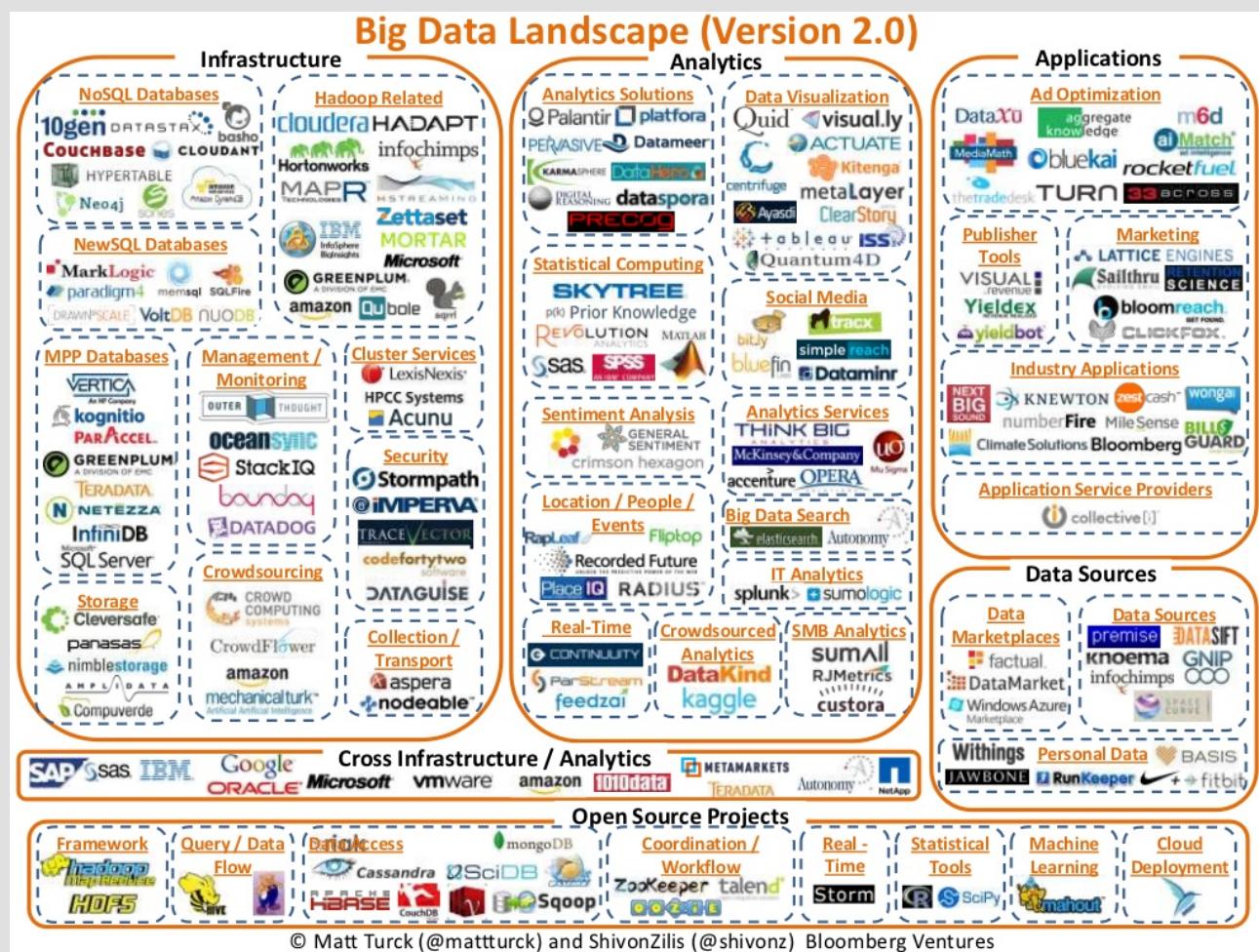
[dave@vcdave.com](mailto:dave@vcdave.com)

[blogs.forbes.com/davefeinleib](http://blogs.forbes.com/davefeinleib)

# Big data Landscape 2.0

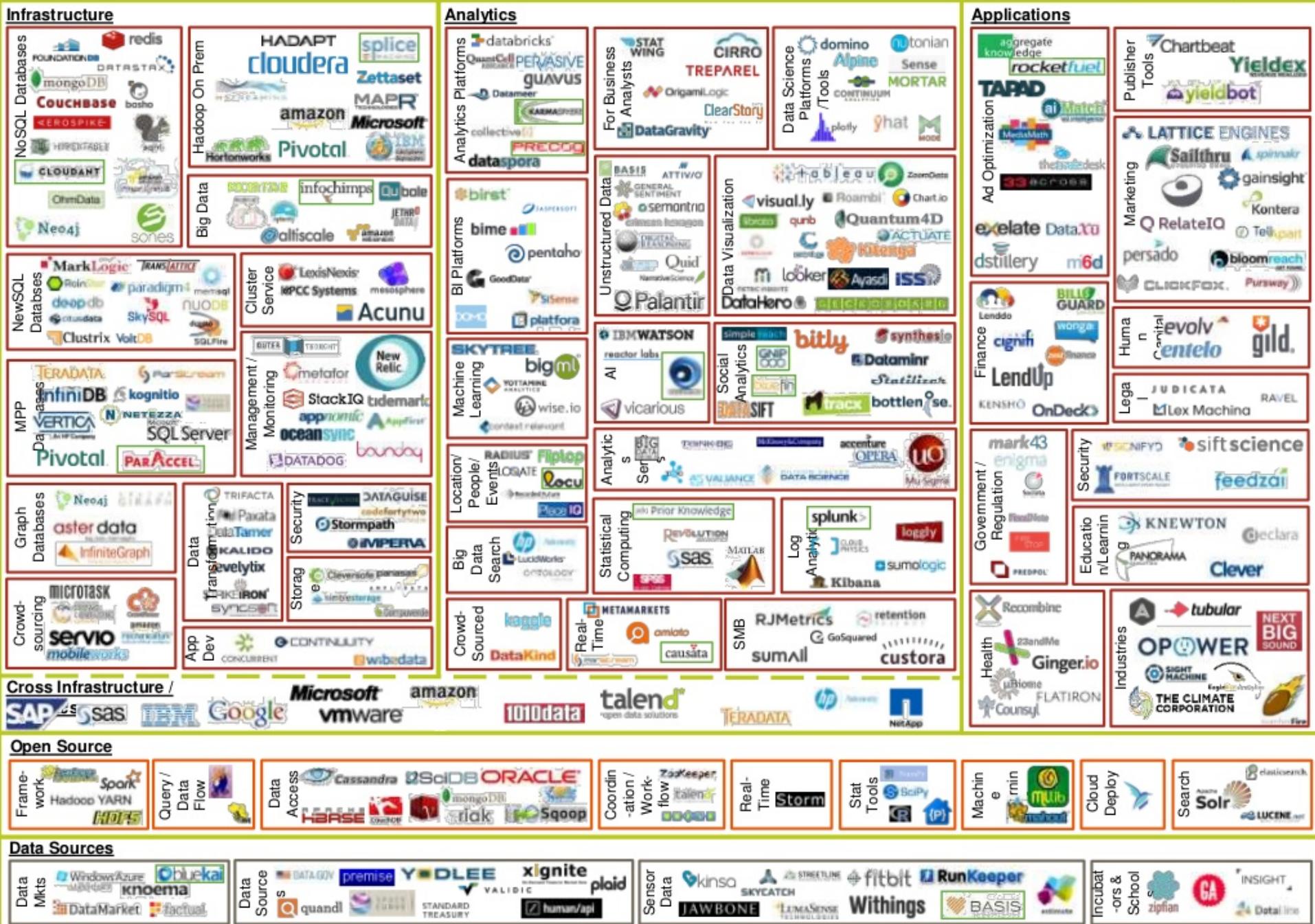


## Paralelismo y Big Data



# BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO





# ¿Qué es Apache Spark?

- Apache Spark
  - Is a fast and general-purpose cluster computing system
    - Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.
  - It provides high-level APIs in Java, Scala and Python
  - It also supports a rich set of higher-level tools including
    - Spark SQL for SQL and structured data processing
    - MLlib for machine learning, GraphX for graph processing
    - Spark Streaming
  - It can access diverse data sources
    - HDFS, Cassandra, Hbase, S3

<https://spark.apache.org/docs/latest/index.html>



Paralelismo y Big Data

## Latest News

Spark 1.3.0 released (Mar 13, 2015)

Spark 1.2.1 released (Feb 09, 2015)

Spark Summit East agenda posted, CFP open for West (Jan 21, 2015)

Spark 1.2.0 released (Dec 18, 2014)

[Archive](#)

# Spark y paralelismo

- El paralelismo en Spark es a muy alto nivel
  - No hay que manejar semáforos, secciones críticas, etc.
- Un mismo programa Spark se puede ejecutar
  - En un PC con procesador multi-núcleo
  - En un cluster propio (standalone deploy mode)
  - Apache Mesos (<https://mesos.apache.org> )
  - Hadoop YARN

# Programación en Spark

- Objetivo del taller
  - Aplicación (simplificada) para obtener *trending topics* de Twitter
  - Fuente datos
    - Fichero de tweets de un partido de baloncesto de la liga universitaria de los EE.UU

443955894288535552	@SebastianH007	RT BaylorMBB: That's a wrap from Sprint Center. Baylor moves on to the #Big12MBB Quarterfinals to face Oklahoma. #SicOU http://t.co/uhjn3C?	Thu Mar 13 04:45:08 CET 2014
443955835442429953	@BaylorMBB	Cory J presenting coach Drew with game ball from his 200th win at Baylor. #SicEm #Big12MBB http://t.co/JnbwHZIxZc Wac	Thu Mar 13 04:44:54 CET 2014
443955799417552896	@DerBaylorBaer	RT HeathNielsen: It's unbelievable how far this man has taken Baylor basketball since we hit rock bottom in 2003. Congrats BUDREW 200 ht?	Thu Mar 13 04:44:45 CET 2014
443955790194282496	@wingoz	RT VoiceofBears: RT BaylorAthletics: RT BaylorMBB: Congratulations to Coach BUDREW on career win No. 200. #SicEm http://t.co/Z74PQydzU5	Thu Mar 13 04:44:43 CET 2014
443955733457944576	@HeathNielsen	It's unbelievable how far this man has taken Baylor basketball since we hit rock bottom in 2003. Congrats BUDREW 200 http://t.co/l2dsT5Psxg Heart of Texas	Thu Mar 13 04:44:29 CET 2014
443955697273671681	@austin_anderson	RT BaylorMBB: Congratulations to Coach BUDREW on career win No. 200. #SicEm http://t.co/l7IBT0IlyR Dallas	Thu Mar 13 04:44:21 CET 2014
443955693758857217	@EmbracinChaos	RT BaylorMBB: That's a wrap from Sprint Center. Baylor moves on to the #Big12MBB Quarterfinals to face Oklahoma. #SicOU http://t.co/uhjn3C?	Thu Mar 13 04:44:20 CET 2014

# Programación en Spark

- Herramientas de desarrollo:

- Spark: <https://spark.apache.org/downloads.html>
- IDE de desarrollo: Eclipse Luna, IntelliJ Idea 14
- Maven (gestión de dependencias):  
<http://maven.apache.org/download.cgi?Preferred=ftp://mirror.reverse.net/pub/apache/>



# Ejemplo 1: sumar números



```
1 package es.uma.hackersweek.sparkdemo;
2
3+import org.apache.spark.SparkConf;□
10
12+ * Example to show how to start a Spark application and create and work with a JavaRDD (Res
19
20 public class SparkDemoSumAListOfNumbers {
21+ public static void main(String[] args) {
22     // STEP 1: create a SparkConf object
23     SparkConf conf = new SparkConf().setAppName("SparkDemo") ;
24
25     // STEP 2: create a Java Spark context
26     JavaSparkContext sparkContext = new JavaSparkContext(conf) ;
27
28     // An array of Integer objects
29     Integer[] values = new Integer[]{1, 2, 3, 4, 5, 6, 8, 9} ;
30
31     // A list of Integers
32     List<Integer> data = Arrays.asList(values);
33
34     // STEP 3: create a JavaRDD
35     JavaRDD<Integer> distributedData = sparkContext.parallelize(data);
36
37     // STEP 4: compute the sum
38     int sum = distributedData.reduce(new Function2<Integer, Integer, Integer>() {
39         @Override public Integer call(Integer integer, Integer integer2) throws Exception {
40             return integer+integer2;
41         }
42     }) ;
43
44     // STEP 5: print the result
45     System.out.println("The sum is: " + sum) ;
46
47     // STEP 6: stop de spark context
48     sparkContext.stop();
49 }
50 }
```

# Fundamentos de Spark: RDD

- JavaRDD
  - *Resilient distributed dataset*
  - Colección de elementos que pueden ser procesados en paralelo
  - Tolerantes a fallos
- Creación de RDDs
  - Paralelizando una colección de datos que ya existe
  - A partir de datos almacenados en un sistema de ficheros o base de datos externa

# Fundamentos de Spark: tipos de operaciones

- **Transformaciones**
  - Crean datos a partir de otros
- Acciones
  - Devuelven un valor tras procesar un conjunto de datos
- Ejemplos

Transformación	Significado
<b>map(func)</b>	Return a new distributed dataset formed by passing each element of the source through a function func.
<b>filter(func)</b>	Return a new dataset formed by selecting those elements of the source on which func returns true.
<b>flatMap(func)</b>	Similar to map, but each input item can be mapped to 0 or more output items (so func should return a Seq rather than a single item).
<b>union(otherDataset)</b>	Return a new dataset that contains the union of the elements in the source dataset and the argument.

# Fundamentos de Spark: tipos de operaciones

- Transformaciones
  - Crean datos a partir de otros
- **Acciones**
  - **Devuelven un valor tras procesar un conjunto de datos**
- Ejemplos

Acción	Significado
<code>reduce(func)</code>	Aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
<code>collect()</code>	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
<code>count()</code>	Return the number of elements in the dataset.
<code>first()</code>	Return the first element of the dataset.
<code>saveAsTextFile(path)</code>	Write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system.

# Ejemplo 1: sumar números. Compilación



- Maven POM.XML

<http://search.maven.org/>

The Central Repository

SEARCH

spark-core

New: About Central Advanced Search | API Guide | Help

## Artifact Details For [org.apache.spark:spark-core\\_2.10:1.3.0](#)

Click on a link above to browse the repository.

**Project Information**

GroupId:	org.apache.spark
ArtifactId:	spark-core_2.10
Version:	1.3.0

**Dependency Information**

**Apache Maven**

```
<dependency>
  <groupId>org.apache.spark</groupId>
  <artifactId>spark-core_2.10</artifactId>
  <version>1.3.0</version>
</dependency>
```

**Apache Buildr**

**Apache Ivy**

**Groovy Grape**

**Gradle/Grails**

**Scala SBT**

**Leiningen**

**Project Object Model (POM)**

```
<?xml version="1.0" encoding="UTF-8"?>
<project xmlns="http://maven.apache.org/POM/
<parent>
  <artifactId>spark-parent_2.10</artifactId>
  <groupId>org.apache.spark</groupId>
  <version>1.3.0</version>
</parent>
<modelVersion>4.0.0</modelVersion>
<groupId>org.apache.spark</groupId>
<artifactId>spark-core_2.10</artifactId>
<name>Spark Project Core</name>
<url>http://spark.apache.org/</url>
<build>
  <outputDirectory>target/scala-${scala.bi
  <testOutputDirectory>target/scala-${scal
  <resources>
    <resource>
      <directory>src/main/resources</direc
    </resource>
    <resource>
      <directory>..../python</directory>
      <includes>
        <include>pyspark/*.py</include>
      </includes>
    </resource>
  <resources>
```

# Ejemplo 1: sumar números. Compilación



Parallelismo y Big Data

- Compilación

```
MacBook-Pro-de-Antonio:hackersWeek ajnebro$ mvn package
[INFO] Scanning for projects...
[INFO]
[INFO] -----
[INFO] Building Spark 1.0-SNAPSHOT
[INFO] -----
[INFO]
[INFO] --- maven-resources-plugin:2.6:resources (default-resources) @ Spark ---
[WARNING] Using platform encoding (UTF-8 actually) to copy filtered resources, i.e. build is platform dependent!
[INFO] Copying 0 resource
[INFO]
[INFO] --- maven-compiler-plugin:3.1:compile (default-compile) @ Spark ---
[INFO] Nothing to compile - all classes are up to date
[INFO]
[INFO] --- maven-resources-plugin:2.6:testResources (default-testResources) @ Spark ---
[WARNING] Using platform encoding (UTF-8 actually) to copy filtered resources, i.e. build is platform dependent!
[INFO] skip non existing resourceDirectory /Users/ajnebro/Softw/hackersWeek/src/test/resources
[INFO]
[INFO] --- maven-compiler-plugin:3.1:testCompile (default-testCompile) @ Spark ---
[INFO] Nothing to compile - all classes are up to date
[INFO]
[INFO] --- maven-surefire-plugin:2.12.4:test (default-test) @ Spark ---
[INFO]
[INFO] --- maven-jar-plugin:2.4:jar (default-jar) @ Spark ---
[INFO] Building jar: /Users/ajnebro/Softw/hackersWeek/target/Spark-1.0-SNAPSHOT.jar
[INFO]
[INFO] --- maven-assembly-plugin:2.4:single (make-assembly) @ Spark ---
[INFO] Building jar: /Users/ajnebro/Softw/hackersWeek/target/Spark-1.0-SNAPSHOT-jar-with-dependencies.jar
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 24.151 s
[INFO] Finished at: 2015-03-19T18:46:57+01:00
[INFO] Final Memory: 30M/756M
[INFO] -----
```

MacBook-Pro-de-Antonio:hackersWeek ajnebro\$

# Ejemplo 1: sumar números. Ejecución

```
MacBook-Pro-de-Antonio:hackersWeek ajnebro$ cat runSparkDemoSumAListOfNumbers.sh  
spark-submit --class "es.uma.hackersweek.sparkdemo.SparkDemoSumAListOfNumbers" --master local[8]  
target/Spark-1.0-SNAPSHOT-jar-with-dependencies.jar  
MacBook-Pro-de-Antonio:hackersWeek ajnebro$
```

```
MacBook-Pro-de-Antonio:hackersWeek ajnebro$ ./runSparkDemoSumAListOfNumbers.sh  
The sum is: 38  
MacBook-Pro-de-Antonio:hackersWeek ajnebro$
```



## Paralelismo y Big Data

```
1 package es.uma.hackersweek.sparkdemo;
2
3 import org.apache.spark.SparkConf;
4
5 * Example showing how to read data from files. The applications sum a list of integer values
6
7 public class SparkDemoSumAListOfNumbersStoredInAFile {
8     public static void main(String[] args) {
9
10        // STEP 1: argument checking
11        if (args.length == 0) {
12            throw new RuntimeException("The number of args is 0. Usage: "
13                + "SparkDemoSumAListOfNumberStoredInAFile fileOrDirectoryName") ;
14        }
15
16        // STEP 2: create a SparkConf object
17        SparkConf conf = new SparkConf().setAppName("SparkDemo") ;
18
19        // STEP 3: create a Java Spark context
20        JavaSparkContext sparkContext = new JavaSparkContext(conf) ;
21
22        // STEP 4: read the lines from the file(s)
23        JavaRDD<String> lines = sparkContext.textFile(args[0]) ;
24
25        // STEP 5: map the string coded numbers into integers
26        JavaRDD<Integer> listOfNumbers = lines.map(new Function<String, Integer>() {
27            @Override public Integer call(String s) throws Exception {
28                return Integer.valueOf(s);
29            }
30        });
31
32        // STEP 6: compute the sum
33        int sum = listOfNumbers.reduce(new Function2<Integer, Integer, Integer>() {
34            @Override public Integer call(Integer integer, Integer integer2) throws Exception {
35                return integer+integer2;
36            }
37        });
38
39        // STEP 7: print the result
40        System.out.println("The sum is: " + sum) ;
41
42        // STEP 8: stop de spark context
43        sparkContext.stop();
44    }
45
46
47
48
49
50
51
52
53
54 }
```

# Ejemplo 2.1: sumar números de ficheros

```

1 package es.uma.hackersweek.sparkdemo;
2
3 import org.apache.spark.SparkConf;
4
5 /**
6  * Example showing how to read data from files. The applications sum a list of integer values.
7  * Each data file must contain an integer per line. Compact version using lambda expression
8  * and fluent code.
9  *
10 * @author Antonio J. Nebro <antonio@lcc.uma.es>
11 */
12
13
14 public class SparkDemoSumAListOfNumbersStoredInAFileCompactVersion {
15     public static void main(String[] args) {
16
17         // STEP 1: argument checking
18         if (args.length == 0) {
19             throw new RuntimeException("The number of args is 0. Usage: "
20                     + "SparkDemoSumAListOfNumbersStoredInAFileCompactVersion fileOrDirectoryName") ;
21         }
22
23
24         // STEP 2: create a SparkConf object
25         SparkConf conf = new SparkConf().setAppName("SparkDemo") ;
26
27         // STEP 3: create a Java Spark context
28         JavaSparkContext sparkContext = new JavaSparkContext(conf) ;
29
30         // STEP 4: read the lines from the file(s)
31         JavaRDD<String> lines = sparkContext.textFile(args[0]) ;
32
33         // STEP 5: compute the sum
34         int sum = lines.map(s -> Integer.valueOf(s)).reduce((integer, integer2) -> integer+integer2) ;
35
36         // STEP 6: print the result
37         System.out.println("The sum is: " + sum) ;
38
39         // STEP 8: stop de spark context
40         sparkContext.stop();
41     }
42 }
```

# Fundamentos de Spark: tipos de operaciones

- Transformaciones
  - Crean datos a partir de otros
- Acciones
  - Devuelven un valor tras procesar un conjunto de datos
- **Carácterística CLAVE:**
  - Las transformaciones son LAZY (se ejecutan de forma perezosa)
  - Únicamente se computan cuando una acción requiere devolver un valor



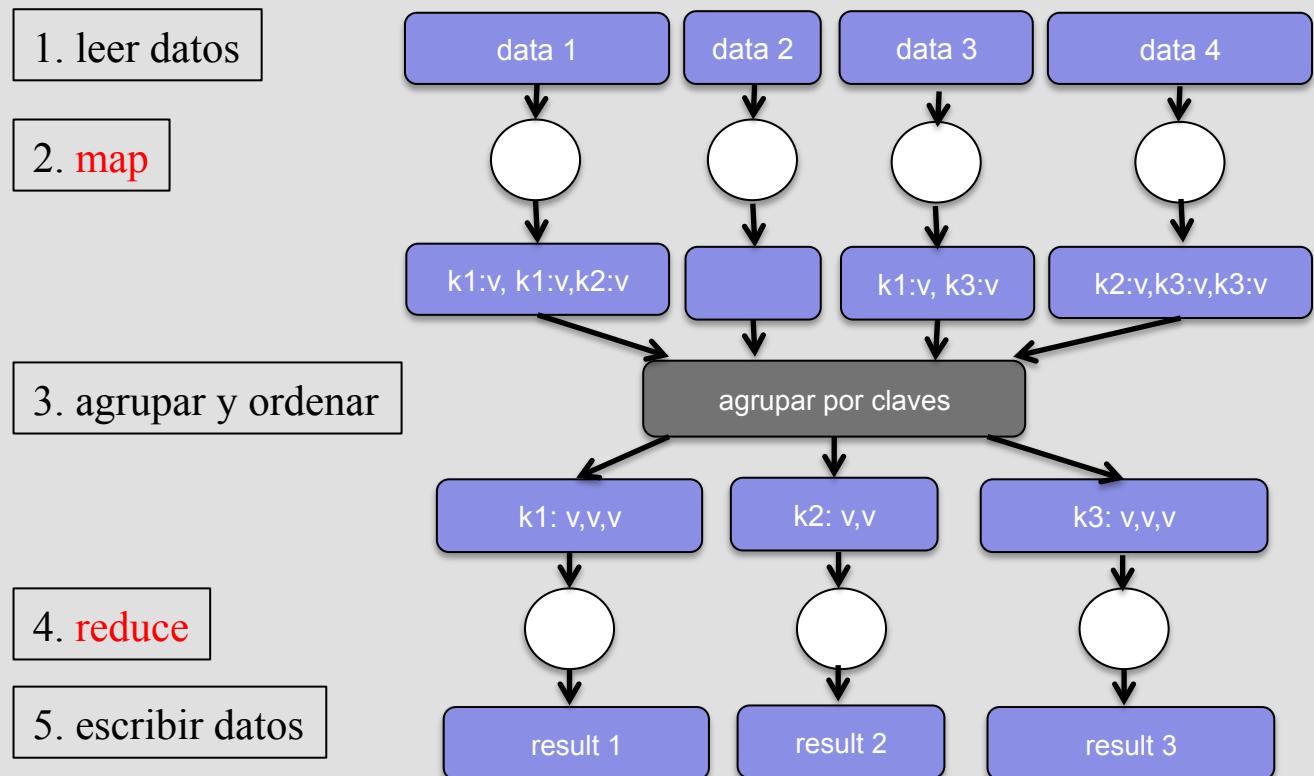
# Ejemplo 3: contar caracteres de ficheros

## Parallelismo y Big Data

```
1 package es.uma.hackersweek.sparkdemo;
2
3 import org.apache.log4j.Logger;
4
5 * Spark applications which counts the number of lines of a file (or the files in a di
6 public class SparkDemoFindCharactersInFiles {
7     static Logger log = Logger.getLogger(SparkDemoFindCharactersInFiles.class.getName());
8
9     public static void main(String[] args) throws ClassNotFoundException {
10
11         // STEP 1: argument checking
12         if (args.length < 1) {
13             log.fatal("Syntax Error: there must be one argument (a file or directory name)");
14             throw new RuntimeException();
15         }
16
17         String fileName = args[0];
18
19         // STEP 2: create a SparkConf object
20         SparkConf sparkConf = new SparkConf().setAppName("FindCharacters");
21
22         // STEP 3: create a Java Spark context
23         JavaSparkContext sparkContext = new JavaSparkContext(sparkConf);
24
25         // STEP 4: read lines of files
26         JavaRDD<String> lines = sparkContext.textFile(fileName).cache();
27
28         // STEP 5: filter and count the number of lines containing the character 'a'
29         long numAs = lines.filter(new Function<String, Boolean>() {
30             public Boolean call(String s) {
31                 return s.contains("a");
32             }
33         }).count();
34
35         // STEP 6: filter and count the number of lines containing the character 'b'
36         long numBs = lines.filter(new Function<String, Boolean>() {
37             public Boolean call(String s) {
38                 return s.contains("b");
39             }
40         }).count();
41
42         // STEP 8: print the results
43         System.out.println("Lines with a: " + numAs + ", lines with b: " + numBs);
44
45         // STEP 9: stop the spark context
46     }
47 }
```

# Ejemplo 4: contar palabras de ficheros

- Programa “Hello Word” en Big Data
  - Contar todas las ocurrencias de todas las palabras de un conjunto de ficheros
  - Paradigma Map/Reduce



# Ejemplo 4: contar palabras de ficheros

```
package es.uma.hackersweek.sparkdemo;

import org.apache.log4j.Logger;

 * Example implementing the Spark version of the "Hello World" Big Data program: counting the words in a file.

public class SparkDemoCountWordsFromFiles {
    static Logger log = Logger.getLogger(SparkDemoCountWordsFromFiles.class.getName());

    public static void main(String[] args) {

        // STEP 1: create a SparkConf object
        if (args.length < 1) {
            log.fatal("Syntax Error: there must be one argument (a file name or a directory)");
            throw new RuntimeException();
        }

        // STEP 2: create a SparkConf object
        SparkConf sparkConf = new SparkConf().setAppName("Spark Word count");

        // STEP 3: create a Java Spark context
        JavaSparkContext sparkContext = new JavaSparkContext(sparkConf);

        // STEP 4: read lines of files
        JavaRDD<String> lines = sparkContext.textFile(args[0]);
    }
}
```

# Ejemplo 4: contar palabras de ficheros

```

43     // STEP 5: split the lines into words
44     JavaRDD<String> words = lines.flatMap(
45         new FlatMapFunction<String, String>() {
46             public Iterable call(String s) throws Exception {
47                 return Arrays.asList(s.split(" "));
48             }
49         } );
50
51     // STEP 6: map operation to create pairs <word, 1> por each word
52     JavaPairRDD<String, Integer> ones = words.mapToPair(
53         new PairFunction<String, String, Integer>() {
54             public Tuple2<String, Integer> call(String string) {
55                 return new Tuple2<String, Integer>(string, 1);
56             }
57         }
58     );
59
60     // STEP 6: reduce operation that sum the values of all the pairs having the same key (word),
61     // generating a pair <key, sum>
62     JavaPairRDD<String, Integer> counts = ones.reduceByKey(
63         new Function2<Integer, Integer, Integer>() {
64             public Integer call(Integer integer, Integer integer2) throws Exception {
65                 return integer + integer2 ;
66             }
67         }
68     );
69
70     // STEP 7: sort the results by key
71     List<Tuple2<String, Integer>> output = counts.sortByKey().collect() ;
72
73     // STEP 8: print the results
74     for (Tuple2<?, ?> tuple : output) {
75         System.out.println(tuple._1() + ":" + tuple._2());
76     }
77
78     // STEP 9: stop the spark context
79     sparkContext.stop();
80 }
81 }
```

# Ejemplo 5: trending topics in tweets

- Aplicación (muy simplificada) para obtener *trending topics* de Twitter
- Fuente datos
  - Fichero de tweets de un partido de baloncesto de la liga universitaria de los EE.UU

443955894288535552	@SebastianH007	RT BaylorMBB: That's a wrap from Sprint Center. Baylor moves on to the #Big12MBB Quarterfinals to face Oklahoma. #SicOU http://t.co/uhjn3C?	Thu Mar 13 04:45:08 CET 2014
443955835442429953	@BaylorMBB	Cory J presenting coach Drew with game ball from his 200th win at Baylor. #SicEm #Big12MBB http://t.co/JnbwHZIxZc Wac	Thu Mar 13 04:44:54 CET 2014
443955799417552896	@DerBaylorBaer	RT HeathNielsen: It's unbelievable how far this man has taken Baylor basketball since we hit rock bottom in 2003. Congrats BUDREW 200 ht?	Thu Mar 13 04:44:45 CET 2014
443955790194282496	@wingoz	RT VoiceofBears: RT BaylorAthletics: RT BaylorMBB: Congratulations to Coach BUDREW on career win No. 200. #SicEm http://t.co/Z74PQydzU5	Thu Mar 13 04:44:43 CET 2014
443955733457944576	@HeathNielsen	It's unbelievable how far this man has taken Baylor basketball since we hit rock bottom in 2003. Congrats BUDREW 200 http://t.co/l2dsT5Psxg Heart of Texas	Thu Mar 13 04:44:29 CET 2014
443955697273671681	@austin_anderson	RT BaylorMBB: Congratulations to Coach BUDREW on career win No. 200. #SicEm http://t.co/l7IBT0IlyR Dallas	Thu Mar 13 04:44:21 CET 2014
443955693758857217	@EmbracinChaos	RT BaylorMBB: That's a wrap from Sprint Center. Baylor moves on to the #Big12MBB Quarterfinals to face Oklahoma. #SicOU http://t.co/uhjn3C?	Thu Mar 13 04:44:20 CET 2014

# Ejemplo 5: trending topics in tweets



```
1 package es.uma.hackersweek.twitter;
2
3 import org.apache.spark.SparkConf;
4
5 /**
6  * Example showing how to count the most used words in files containing tweets
7  *
8  * @author Antonio J. Nebro <antonio@lcc.uma.es>
9  */
10
11 public class TwitterTrendingTopics {
12     public static void main(String[] args) {
13
14         // STEP 1: argument checking
15         if (args.length == 0) {
16             throw new RuntimeException("The number of args is 0. Usage: "
17                     + "TwitterTrendingTopics fileOrDirectoryName");
18         }
19
20         // STEP 2: create a SparkConf object
21         SparkConf conf = new SparkConf().setAppName("TwitterTrendingTopics");
22
23         // STEP 3: create a Java Spark context
24         JavaSparkContext sparkContext = new JavaSparkContext(conf);
25
26         // STEP 4: read the lines from the file(s)
27         JavaRDD<String> lines = sparkContext.textFile(args[0]);
28
29         // STEP 5: map to select the strings in the column 2, which contains the tweet message
30         JavaRDD<String> messages = lines.map(new Function<String, String>() {
31             public String call(String s) {
32                 String [] split = s.toString().split("\t+");
33
34                 return split[2];
35             }
36         });
37
38         // STEP 5: split the lines into words
39         JavaRDD<String> words = messages.flatMap(
40             new FlatMapFunction<String, String>() {
41                 public Iterable call(String s) throws Exception {
42                     return Arrays.asList(s.split(" "));
43                 }
44             });
45
46         // STEP 6: count the words
47         JavaPairRDD<String, Integer> wordCount = words.countByValue();
48
49         // STEP 7: print the results
50         wordCount.saveAsTextFile("output");
51
52         System.out.println("Number of words: " + wordCount.count());
53
54     }
55 }
```

# Ejemplo 5: trending topics in tweets



```
56     // STEP 6: map operation to create pairs <word, 1> por each word
57     JavaPairRDD<String, Integer> ones = words.mapToPair(
58         new PairFunction<String, String, Integer>() {
59             public Tuple2<String, Integer> call(String string) {
60                 return new Tuple2<String, Integer>(string, 1);
61             }
62         }
63     );
64
65     // STEP 6: reduce operation that sum the values of all the pairs having the same key (word),
66     // generating a pair <key, sum>
67     JavaPairRDD<String, Integer> counts = ones.reduceByKey(
68         new Function2<Integer, Integer, Integer>() {
69             public Integer call(Integer integer, Integer integer2) throws Exception {
70                 return integer + integer2 ;
71             }
72         }
73     );
74
75     // STEP 7: invert the keys and values to use sortByKey() in step 8 on the values instead of the
76     // keys
77     JavaPairRDD<Integer, String> reverse = counts.mapToPair(
78         new PairFunction<Tuple2<String, Integer>, Integer, String>() {
79             @Override public Tuple2<Integer, String> call(Tuple2<String, Integer> stringIntegerTuple2)
80                 throws Exception {
81                 return new Tuple2<Integer, String>(stringIntegerTuple2._2, stringIntegerTuple2._1);
82             }
83         }
84     );
85
86     /*// STEP 8: sort the results by key
87     List<Tuple2<Integer, String>> output = reverse.sortByKey().collect() ;
88
89     // STEP 9.1: print the results to screen
90     for (Tuple2<?, ?> tuple : output) {
91         System.out.println(tuple._1() + ":" + tuple._2());
92     }*/
93
94     // STEP 9.2 print the results to file
95     reverse.sortByKey().saveAsTextFile("data/tweets/output");
96
97     // STEP 10: stop de spark context
98     sparkContext.stop(); }
```

# Conclusiones del taller

- En los últimos 30 años las tecnologías informáticas han estado en continua evolución
  - Pero en los últimos 5 años los avances se suceden a ritmo vertiginoso
  - Ámbitos:
    - Web: HTML 5, CSS3, Javascript, AngularJS, NodeJS, Bootstrap, etc.
    - Dispositivos móviles: Android, iOS, Windows 10?
    - Big Data: Hadoop, MapReduce, Hive, Hbase, Spark, etc.
- Spark es un ejemplo de tecnología que ha aparecido recientemente y que tiene una gran proyección
  - Nuevas versiones casi cada mes
  - Adopción por muchas empresas
  - Herramienta de alto nivel y elevado rendimiento
  - Puede ser la nueva tecnología principal sobre la que se sustenten las aplicaciones de Big Data
- A saber qué ocurre el año que viene por estas fechas ...

# Gracias por asistir



# Paralelismo y Big Data

**GRACIAS**  
ARIGATO  
SHUKURIA  
JUSPAXAR  
GOZAIMASHITA  
EFCHARISTO  
KOMAPSUMMIDA  
MAAKE  
LAH  
GRAZIE  
MEHRBANI  
PALDIES  
TASHAKKUR ATU  
SUKSAMA  
EKHMET  
HATUR GUI  
TINGKI  
BİYAN  
SHUKRIA

SPASSIBO  
NUHUN  
SNACHALHUYA  
CHALTU  
WABEEJA  
MAITEKA  
YUSPIGARATAM  
DHANYABRAD  
AHNA  
HUI  
ATTO  
NERSI  
DENKAUJA  
NEHACHALHYA  
UNALCHEEESH  
HATUR GUI  
EKOU  
SIKOMO  
MAKETAI  
MINMONCHAR