

Arquitecturas Big Data

Mario Pérez Esteso



@_Mario_Perez

mario@geekytheory.com

Índice

1. Introducción
2. Tecnologías de procesamiento
3. Arquitectura
4. Casos de estudio
5. Conclusiones

A vibrant blue background featuring the text "BIG DATA" in large, white, outlined letters. Surrounding the text are various geometric shapes and icons representing data analysis, including circles, squares, triangles, lines, and charts. The design is modern and abstract, with long shadows cast by the shapes.



A DAY IN THE INTERNET

In one day, enough information is
consumed by internet traffic to fill

168 MILLION DVDS.



= 1 MILLION



294 BILLION
emails are sent.



It would take
2 years to process
that many pieces
of mail in the US.



2 MILLION BLOG POSTS are written.

Enough posts to fill
Time Magazine for 770 years.



172 MILLION
different people
visit Facebook.

Twitter: **40 MILLION**
LinkedIn: **22 MILLION**
Google+: **20 MILLION**
Pinterest: **17 MILLION**

864,000 HOURS OF VIDEO

are uploaded to YouTube.



That's 98 years of
non-stop cat videos.



Internet users spend
14.6 MINUTES
viewing porn online.

The average fap session
is 12 minutes.



378,000

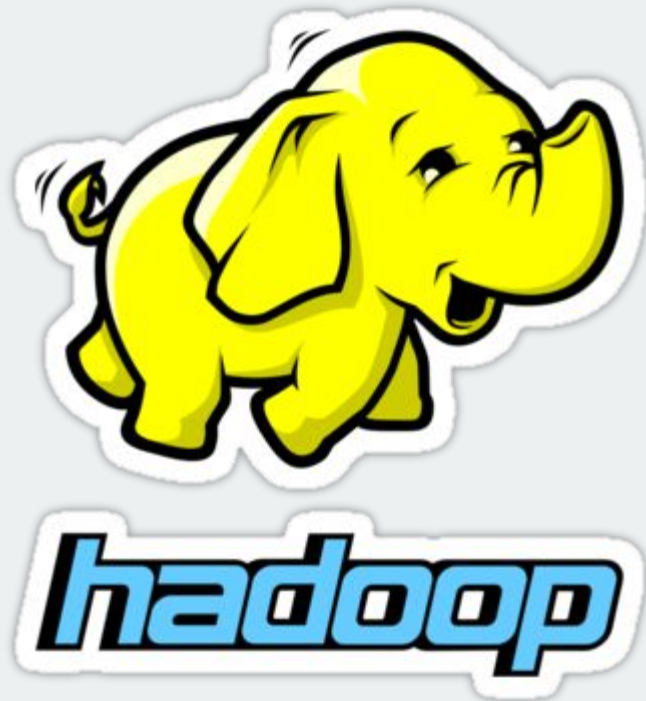
Number of iPhones Sold



371,000

Number of babies born





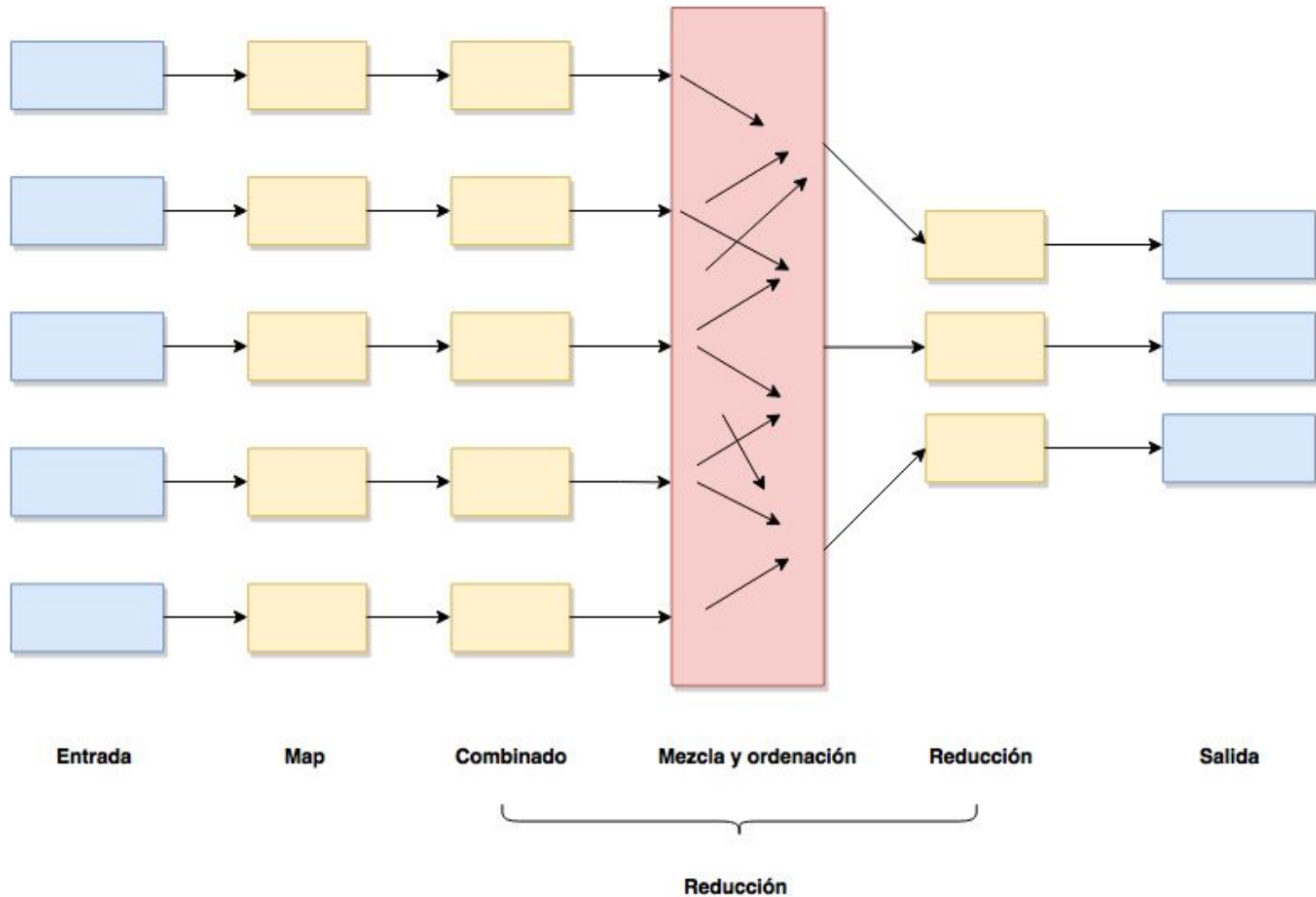
Spark[★]

Apache Hadoop

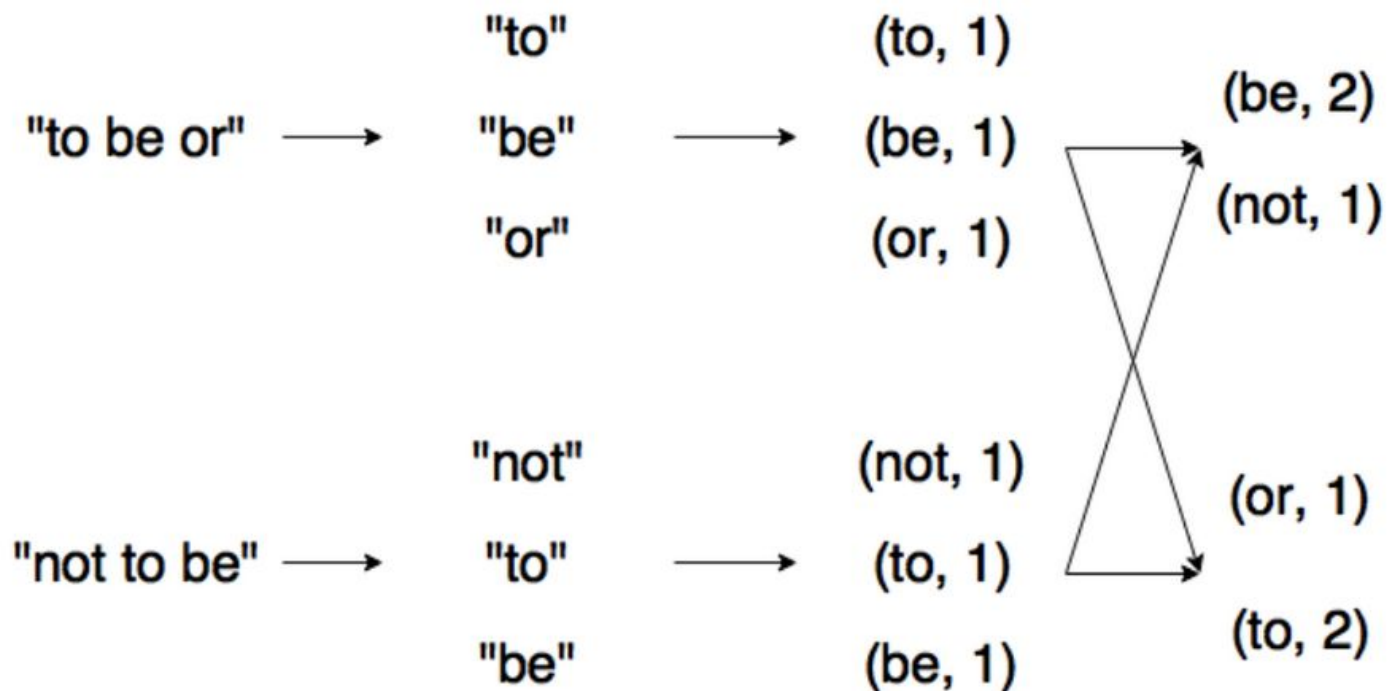
Hadoop está formado por los siguientes módulos:

- Hadoop Common.
- Hadoop Distributed Filesystem (HDFS).
- Hadoop YARN (Yet Another Resource Negotiator).
- MapReduce.

Apache Hadoop MapReduce



Apache Hadoop MapReduce



Apache Spark

Es considerado el primer software de código abierto que hace la programación distribuida realmente accesible a los científicos de datos.

Spark mantiene la escalabilidad lineal y la tolerancia a fallos de MapReduce, pero amplía sus bondades gracias a varias funcionalidades: **DAG** y **RDD**.

Spark
SQL

Spark
Streaming

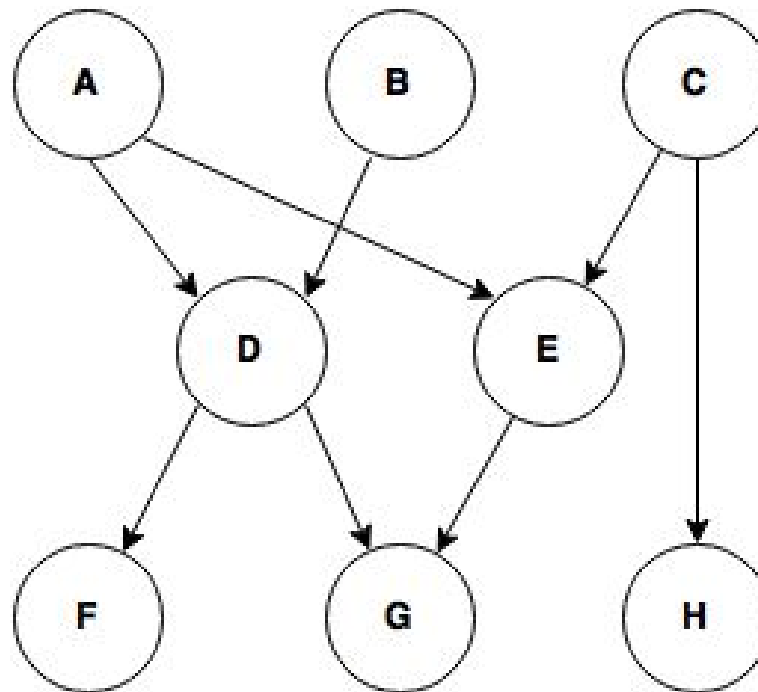
MLlib
(machine
learning)

GraphX
(graph)

Apache Spark

A diagram showing the components of the Apache Spark ecosystem. At the top, four dark blue rectangular boxes are arranged horizontally, each containing white text. From left to right, the boxes are labeled 'Spark SQL', 'Spark Streaming', 'MLlib (machine learning)', and 'GraphX (graph)'. Below these four boxes is a single, wider light blue rectangular box that spans the width of the four boxes above it, containing the text 'Apache Spark' in white. This layout visually represents the four main interfaces or libraries of Spark sitting on top of the core Spark engine.

Apache Spark - DAG (Directed Acyclic Graph)



Apache Spark - RDD (Resilient Distributed Dataset)

Un objeto RDD permite a los programadores realizar operaciones sobre grandes cantidades de datos en clusters de una manera rápida y tolerante a fallos.

Mantener los datos en memoria puede mejorar el rendimiento de una aplicación considerablemente.



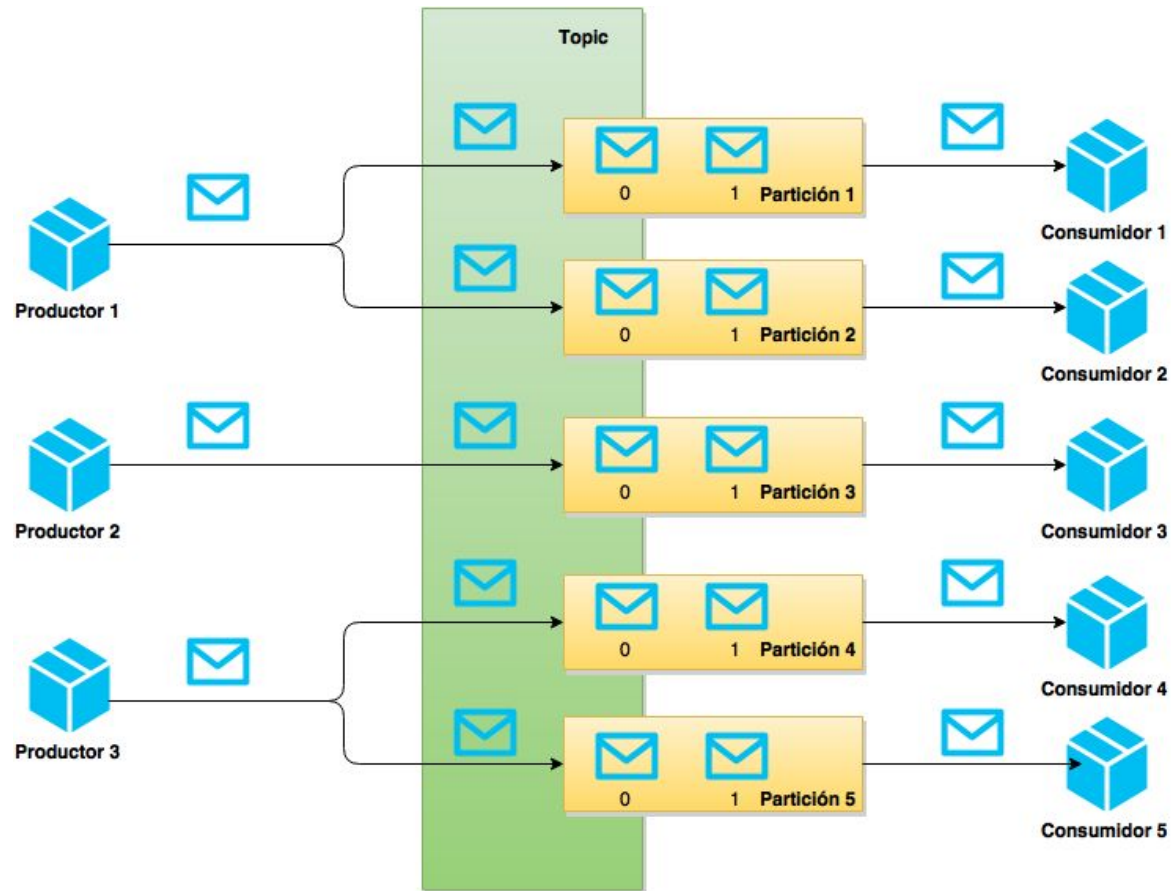
Sistemas de streaming



Apache Kafka

Sistema de mensajería distribuido basado en publicación-suscripción.

Apache Kafka



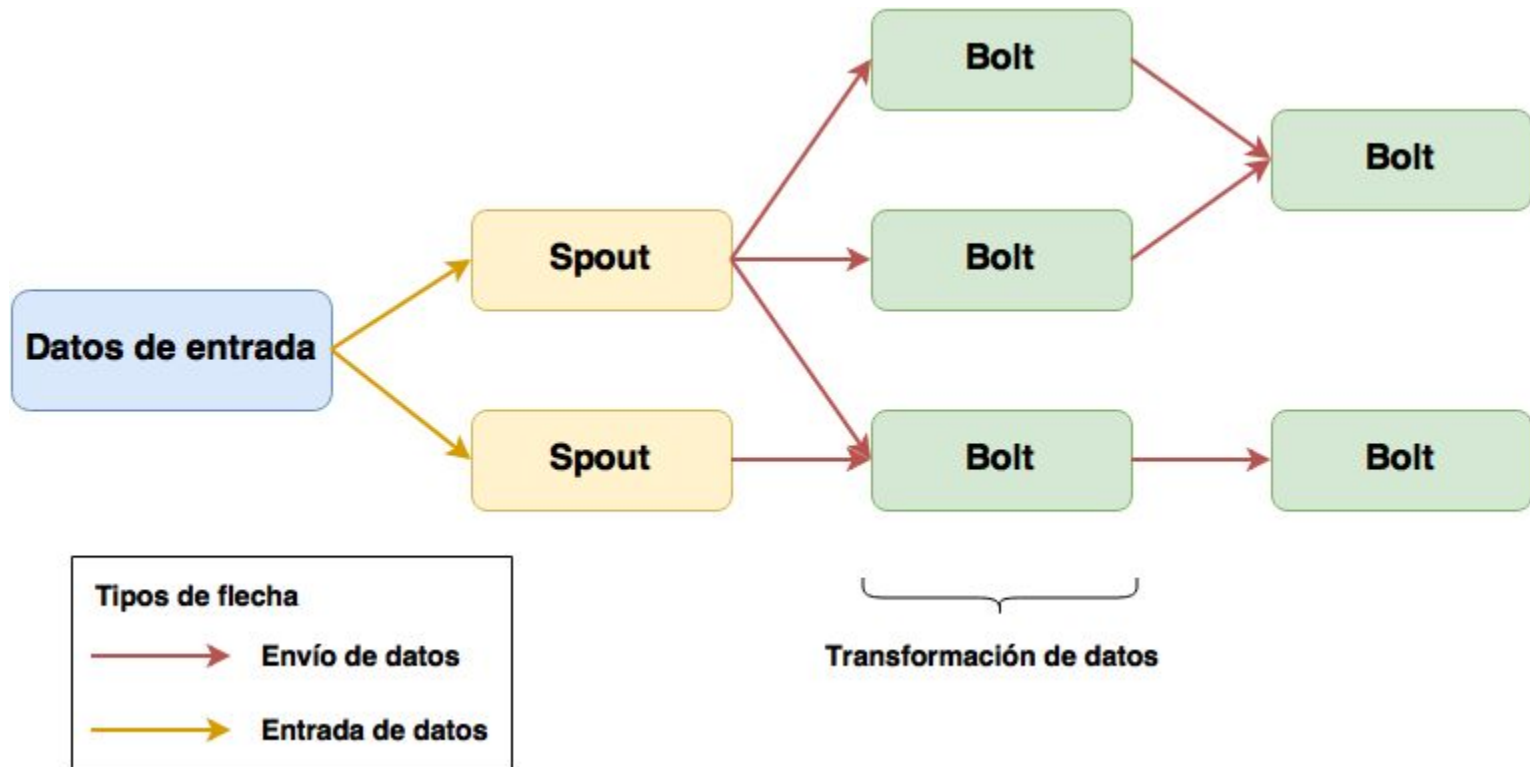
Apache Storm

Se basa en una arquitectura maestro-esclavo y su objetivo es procesar datos en tiempo real.

Se compone de dos partes principales:

- Spout
- Bolt

Apache Storm

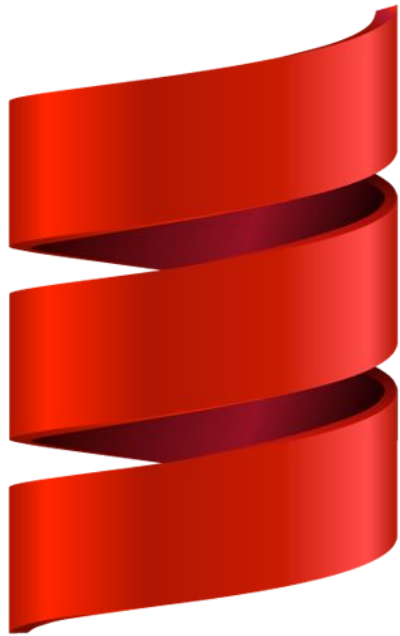


Spark Streaming

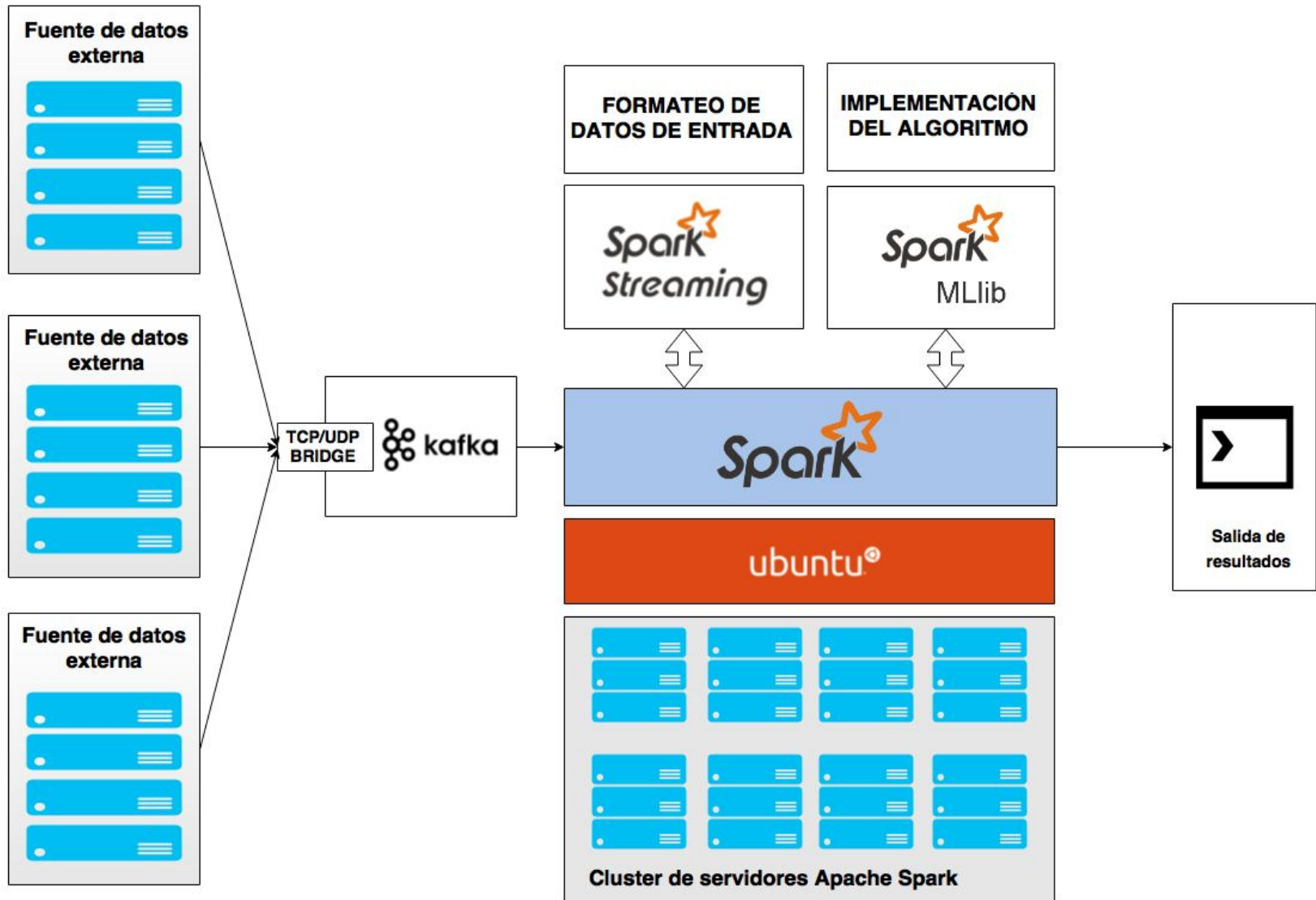


Spark Streaming





Arquitectura



Detección de anomalías en redes

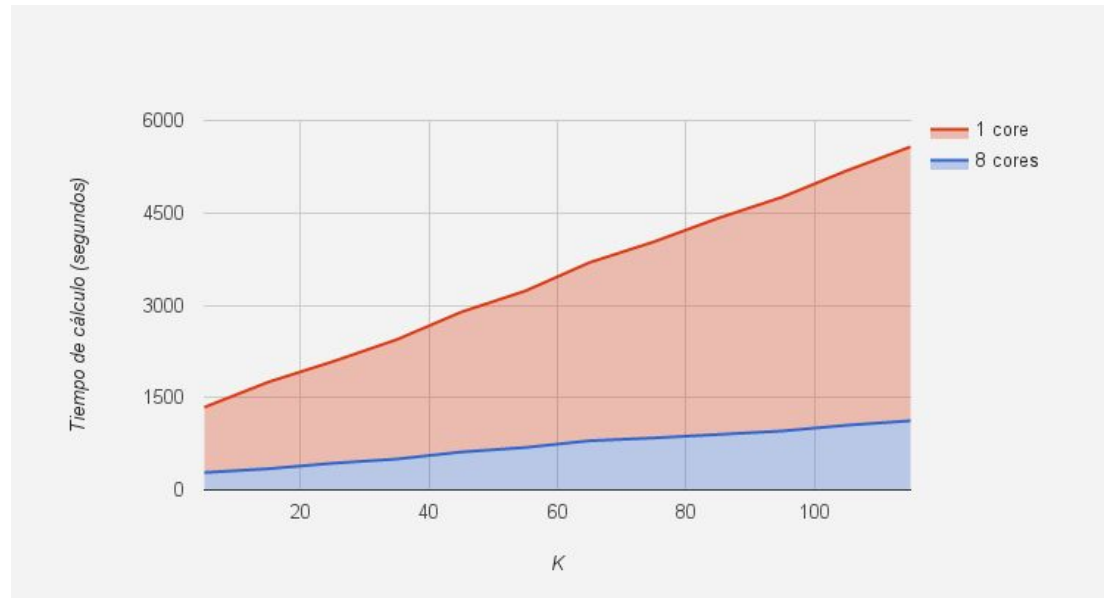
Caso de estudio - Detección de anomalías

Conjunto de datos:

4.9 millones de registros.

Algoritmo:

KMeans Clustering.



```
1 from pyspark.mllib.clustering import KMeans
2 model = KMeans.train(data,
3                       k,
4                       maxIterations=10,
5                       runs=5,
6                       initializationMode="random")
```


Predicción de fallos online

Caso de estudio - Predicción de fallos online

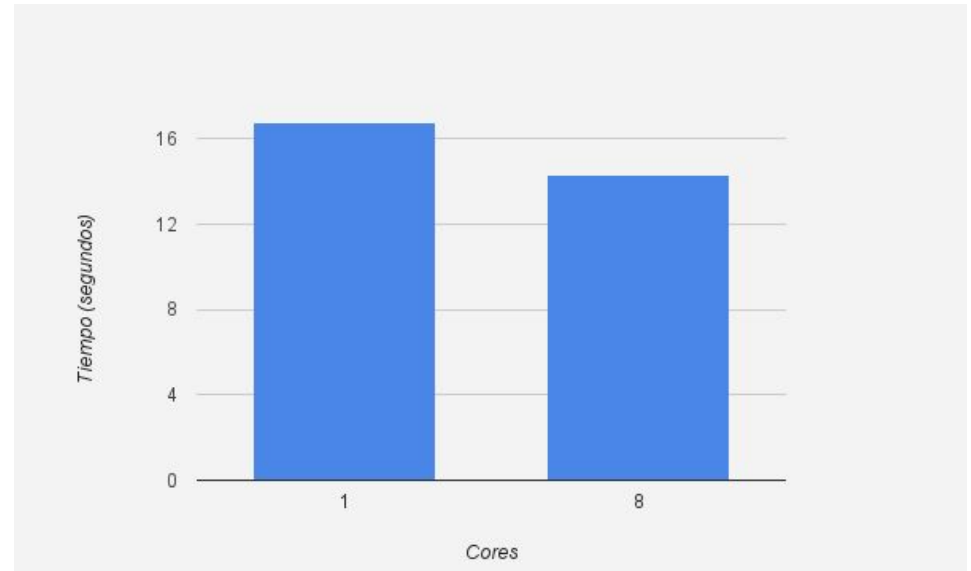
Conjunto de datos:

10661 líneas

88 columnas

Algoritmo:

Random Forests.



```
1 model = RandomForest.trainClassifier(  
2     trainingData,  
3     numClasses=10,  
4     categoricalFeaturesInfo={},  
5     numTrees=5,  
6     featureSubsetStrategy="auto",  
7     impurity='gini',  
8     maxDepth=5,  
9     maxBins=32)
```

Análisis de sentimientos en Twitter

Caso de estudio - Análisis de sentimientos

Conjunto de datos:

Mensajes en Twitter.

Algoritmo:

Recuento de palabras positivas y negativas.

```
1 | oauth.consumerKey = XXXXXXXXXXXXXXXX
2 | oauth.consumerSecret = XXXXXXXXXXXXXXXX
3 | oauth.accessToken = XXXXXXXXXXXXXXXX
4 | oauth.accessTokenSecret = XXXXXXXXXXXXXXXX
```

```
1 | val sc = new SparkContext(conf)
2 | val ssc = new StreamingContext(sc, Seconds(5))
3 | val filters = new Array[String](1)
4 | filters(0) = "greece"
5 | val twitterStream = TwitterUtils.createStream(ssc,
6 |                                             None,
7 |                                             filters)
```

Análisis de sentimientos en Twitter

Puntuación de mensajes

Time: 1436353630000 ms

(6187380632,RT @IvanaKottasova: Tried to buy a train ticket in Athens. It's free", I am told. "The banks are closed."#Greece)
(6187380678,Greece has the opportunity to repudiate a lot of inconvenient foreign debt.)

```
1 | twitterStream.  
2 |   map{tweet => (tweet.getId(), tweet.getText())}.  
3 |   map{case (id, text) => (id, clean(text))}.  
4 |   map{case (id, words) => (id, rateWordList(words))}.  
5 |   print()
```

Time: 1436361275000 ms

(618770130890133505,1)
(618770131334918144,0)
(618770132014362624,0)
(618770132144377856,0)
(618770133599825920,1)
(618770135097192448,0)
(618770136187727873,-1)
(618770136149790720,0)
(618770137676517376,1)
(618770138267914240,-2)

Conclusiones

Apache Spark es entre 10 y 100 veces más rápido que Hadoop MapReduce.

La programación es muy similar en Scala, Java y Python.

La arquitectura propuesta es válida para cualquier caso de estudio.



@_Mario_Perez

mario@geekytheory.com