# Diseño de Sistemas Distribuidos

## Máster en Ciencia y Tecnología Informática
## Curso 2018-2019

# Sistemas escalables en entornos distribuidos.
# Introducción a Spark

Alejandro Calderón Mateos  &  Jaime Pons Bailly-Bailliere

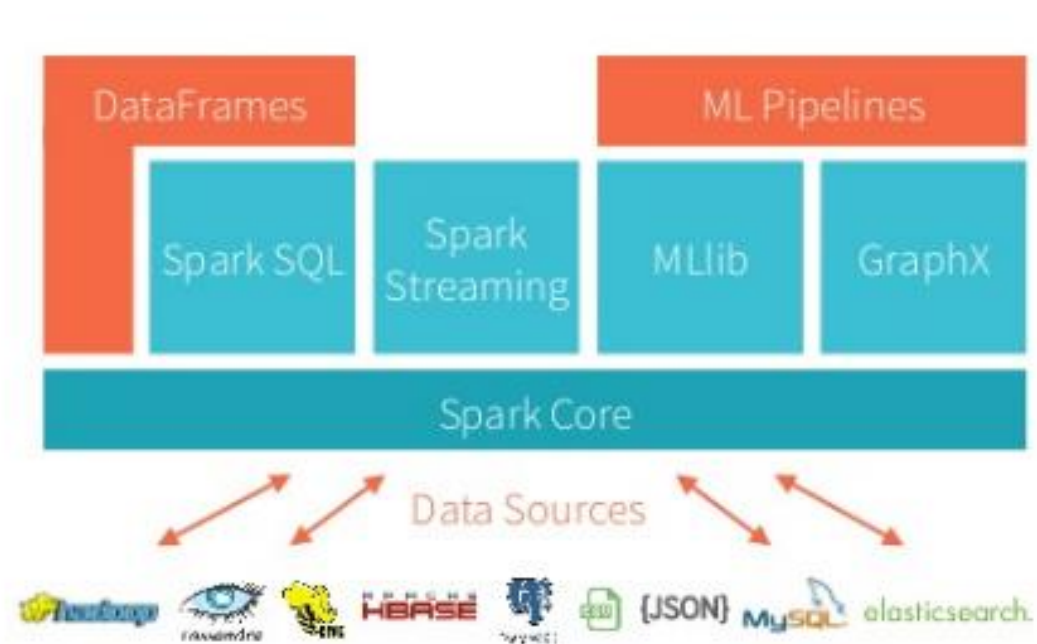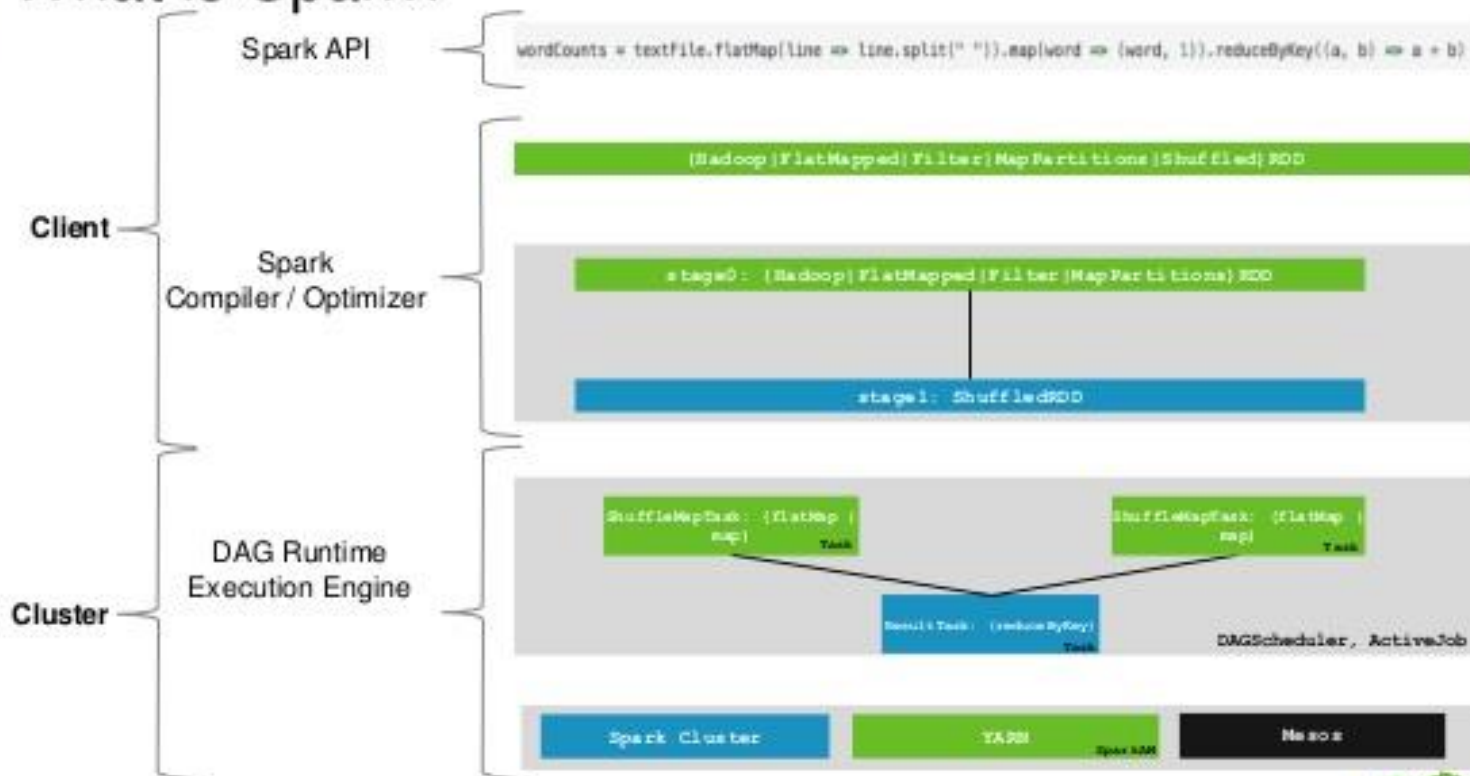acaldero@inf.uc3m.es          jaime@lab.inf.uc3m.es

# Contenidos



- **Introducción**
- *Hand-on*
  - Pre-requisitos e instalación
  - Nodo autónomo
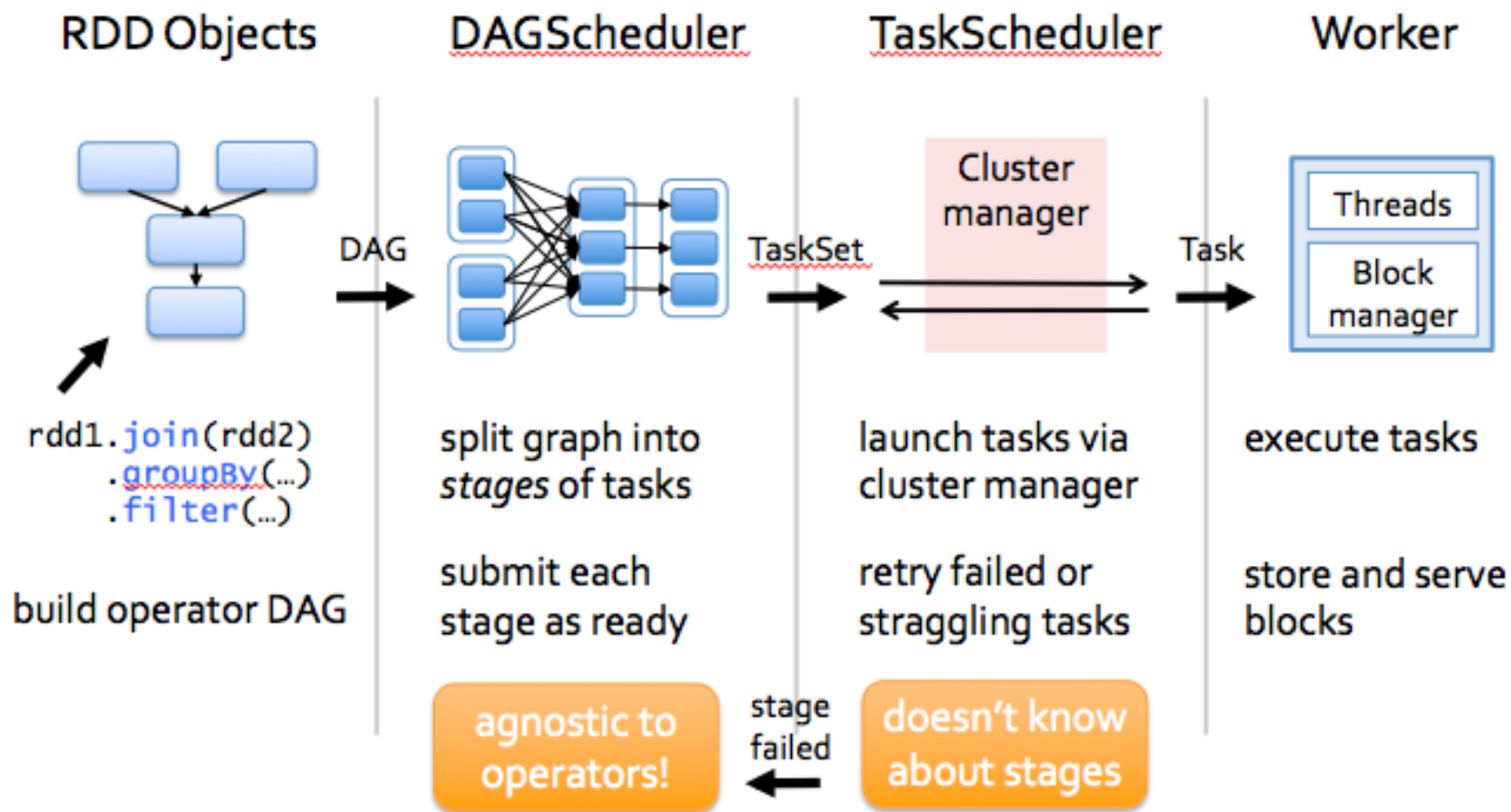  - Cluster
- *Benchmarking*

# Arquitectura

# Arquitectura

# Arquitectura



## RDD Objects

```
rdd1.join(rdd2)
     .groupBy(…)
     .filter(…)
```

build operator DAG

DAG →

## DAGScheduler

split graph into *stages* of tasks

submit each stage as ready

agnostic to operators!

TaskSet →

## TaskScheduler

Cluster manager

launch tasks via cluster manager

retry failed or straggling tasks

stage failed

doesn't know about stages

Task →

## Worker

Threads

Block manager

execute tasks

store and serve blocks

https://sigmoid.com/wp-content/uploads/2015/03/Apache_Spark1.png

# Contenidos



- Introducción
- *Hand-on*
  - **Pre-requisitos e instalación**
  - Nodo autónomo
  - Cluster
- *Benchmarking*

# Spark, Anaconda y Jupyter

**Prerequisitos**      Instalación      Prueba básica

```
acaldero@h1:~$ du –mh –s .
2,8G .
```

# Spark

```
acaldero@h1:~$ sudo apt-get install ssh rsync
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following NEW packages will be installed:
  rsync ssh
…
```

```
acaldero@h1:~$ sudo apt-get install default-jdk
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following extra packages will be installed:
  libice-dev libpthread-stubs0-dev libsm-dev libx11-dev libx11-doc
  libxau-dev libxcb1-dev libxdmcp-dev libxt-dev openjdk-7-jdk
…
```

# Spark

Prerequisitos    **Instalación**    Prueba básica

APACHE
# Spark™

*Lightning-fast unified analytics engine*

**Download**    **Libraries ▾**    **Documentation ▾**    **Examples**    **Community ▾**    **Developers ▾**

## Download Apache Spark™

1. Choose a Spark release: [ 2.4.0 (Nov 02 2018) ⌄ ]

2. Choose a package type: [ Pre-built for Apache Hadoop 2.7 and later                    ⌄ ]

3. Download Spark: spark-2.4.0-bin-hadoop2.7.tgz

4. Verify this release using the 2.4.0 signatures and checksums and project release KEYS.

*Note: Starting version 2.0, Spark is built with Scala 2.11 by default. Scala 2.10 users should download the Spark source package and build with Scala 2.10 support.*

http://spark.apache.org/downloads.html

# Spark

Prerequisitos   **Instalación**   Prueba básica

```
acaldero@h1:~$ wget https://www.apache.org/dyn/closer.lua/spark/spark-2.4.0/spark-2.4.0-bin-hadoop2.7.tgz
…
2018-11-18 12:40:44 (6,02 MB/s) - "spark-2.2.0-bin-hadoop2.7.tgz" guardado [...]


acaldero@h1:~$ tar zxf spark-2.4.0-bin-hadoop2.7.tgz
acaldero@h1:~$ ls -las spark-2.4.0-bin-hadoop2.7
total 96
 4 drwxr-xr-x 12 acaldero acaldero  4096 jul  1 01:09 .
 4 drwx------ 39 acaldero acaldero  4096 oct 17 00:50 ..
 4 drwxr-xr-x  2 acaldero acaldero  4096 jul  1 01:09 bin
 4 drwxr-xr-x  2 acaldero acaldero  4096 jul  1 01:09 conf
 0 drwxr-xr-x  5 acaldero acaldero    47 jul  1 01:09 data
 0 drwxr-xr-x  4 acaldero acaldero    27 jul  1 01:09 examples
12 drwxr-xr-x  2 acaldero acaldero  8192 jul  1 01:09 jars
20 -rw-r--r--  1 acaldero acaldero 17881 jul  1 01:09 LICENSE
 4 drwxr-xr-x  2 acaldero acaldero  4096 jul  1 01:09 licenses
28 -rw-r--r--  1 acaldero acaldero 24645 jul  1 01:09 NOTICE
 4 drwxr-xr-x  8 acaldero acaldero  4096 jul  1 01:09 python
 0 drwxr-xr-x  3 acaldero acaldero    16 jul  1 01:09 R
 4 -rw-r--r--  1 acaldero acaldero  3809 jul  1 01:09 README.md
 4 -rw-r--r--  1 acaldero acaldero   128 jul  1 01:09 RELEASE
 4 drwxr-xr-x  2 acaldero acaldero  4096 jul  1 01:09 sbin
 0 drwxr-xr-x  2 acaldero acaldero    41 jul  1 01:09 yarn
```

# Spark

Prerequisitos      Instalación      **Prueba básica**

```
acaldero@h1:~$ ./bin/run-example SparkPi 5
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
17/10/17 01:02:41 INFO SparkContext: Running Spark version 2.2.0
17/10/17 01:02:42 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using
     builtin-java classes where applicable
17/10/17 01:02:42 INFO SparkContext: Submitted application: Spark Pi
17/10/17 01:02:42 INFO SecurityManager: Changing view acls to: acaldero
17/10/17 01:02:42 INFO SecurityManager: Changing modify acls to: acaldero
17/10/17 01:02:42 INFO SecurityManager: Changing view acls groups to:
17/10/17 01:02:42 INFO SecurityManager: Changing modify acls groups to:
17/10/17 01:02:42 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users
     with view permissions: Set(acaldero); groups with view permissions: Set(); users  with modify
     permissions: Set(acaldero); groups with modify permissions: Set()
17/10/17 01:02:42 INFO Utils: Successfully started service 'sparkDriver' on port 39281.
17/10/17 01:02:42 INFO SparkEnv: Registering MapOutputTracker
17/10/17 01:02:42 INFO SparkEnv: Registering BlockManagerMaster
…
```
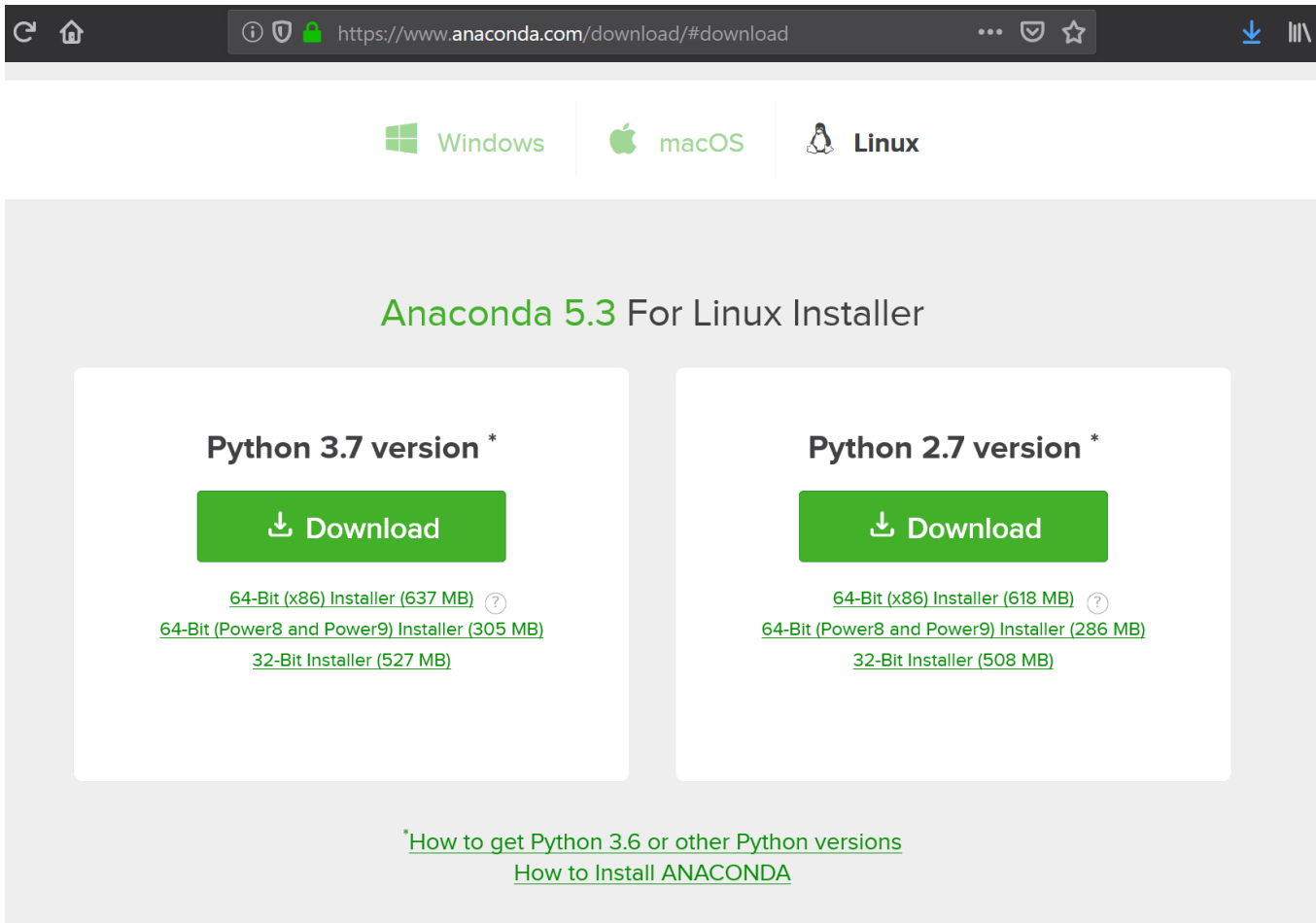
# Anaconda

## Instalación    Prueba básica

https://www.anaconda.com/download/

# Anaconda

## Instalación        Prueba básica

```
acaldero@h1:~$ wget https://repo.continuum.io/archive/Anaconda3-5.3.0-Linux-x86_64.sh

…
2018-11-18 15:12:23 (5,57 MB/s) - "Anaconda3-5.3.0-Linux-x86_64.sh" guardado [...]
```

```
acaldero@h1:~$ chmod a+x Anaconda3-5.3.0-Linux-x86_64.sh
acaldero@h1:~$ ./ Anaconda3-5.3.0-Linux-x86_64.sh
Welcome to Anaconda3 5.3.0 (by Continuum Analytics, Inc.)

In order to continue the installation process, please review the license
agreement.
Please, press ENTER to continue
>>>
…
```

```
acaldero@h1:~$ bash
acaldero@h1:~$ conda update --all
Fetching package metadata .......
Solving package specifications: ..........
…
```

http://jupyter.readthedocs.io/en/latest/install.html#existing-python-new-jupyter

# Spark, Anaconda y Jupyter

## Configuración

```
acaldero@h1:~$ ln -s spark-2.2.0-bin-hadoop2.7 spark
acaldero@h1:~$ echo "export PATH=$PATH:/home/acaldero/spark/bin"              >> .profile
acaldero@h1:~$ echo "export PYSPARK_DRIVER_PYTHON=ipython"                    >> .profile
acaldero@h1:~$ echo "export PYSPARK_DRIVER_PYTHON_OPTS='notebook' pyspark">> .profile
acaldero@h1:~$ source .profile

…
```
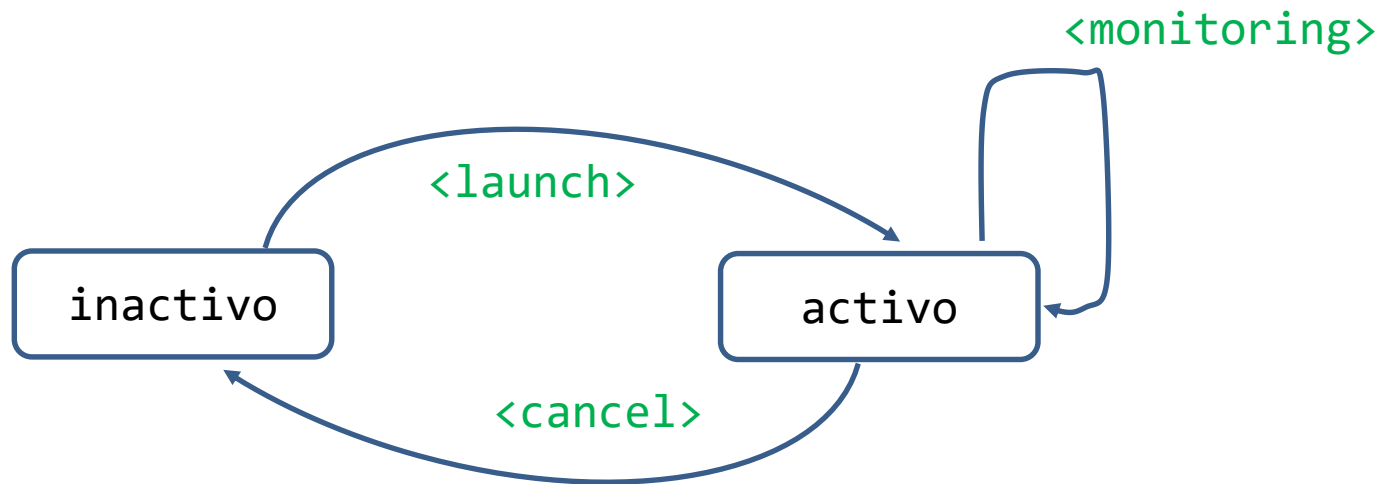
https://gist.github.com/tommycarpi/f5a67c66a8f2170e263c

# Contenidos

&mdash; Introducción

&mdash; ***Hand-on***

- Pre-requisitos e instalación

- **Nodo autónomo**

- Cluster

&mdash; *Benchmarking*

# Spark

## Funcionamiento General
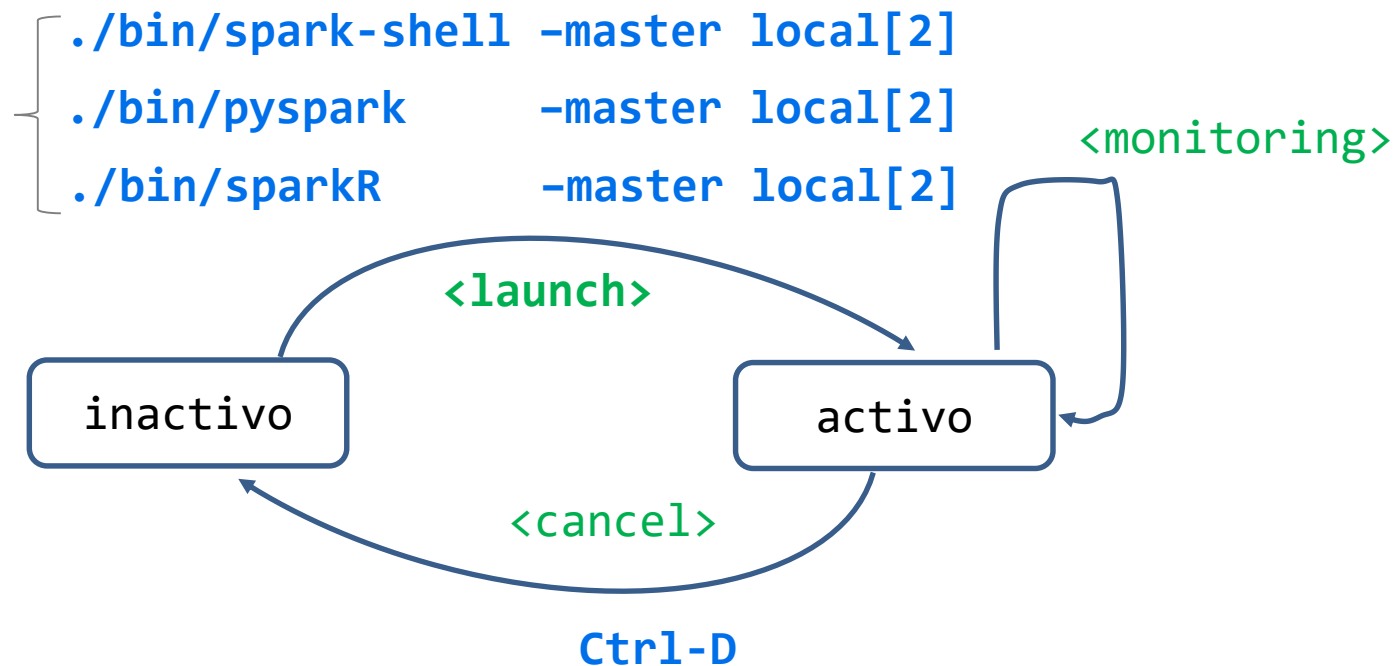
<monitoring>

<launch>

inactivo

activo

# Spark: nodo autónomo

**shell-interactivo**          submit          libro-interactivo

```
./bin/spark-shell –master local[2]
./bin/pyspark     –master local[2]
./bin/sparkR      –master local[2]
```
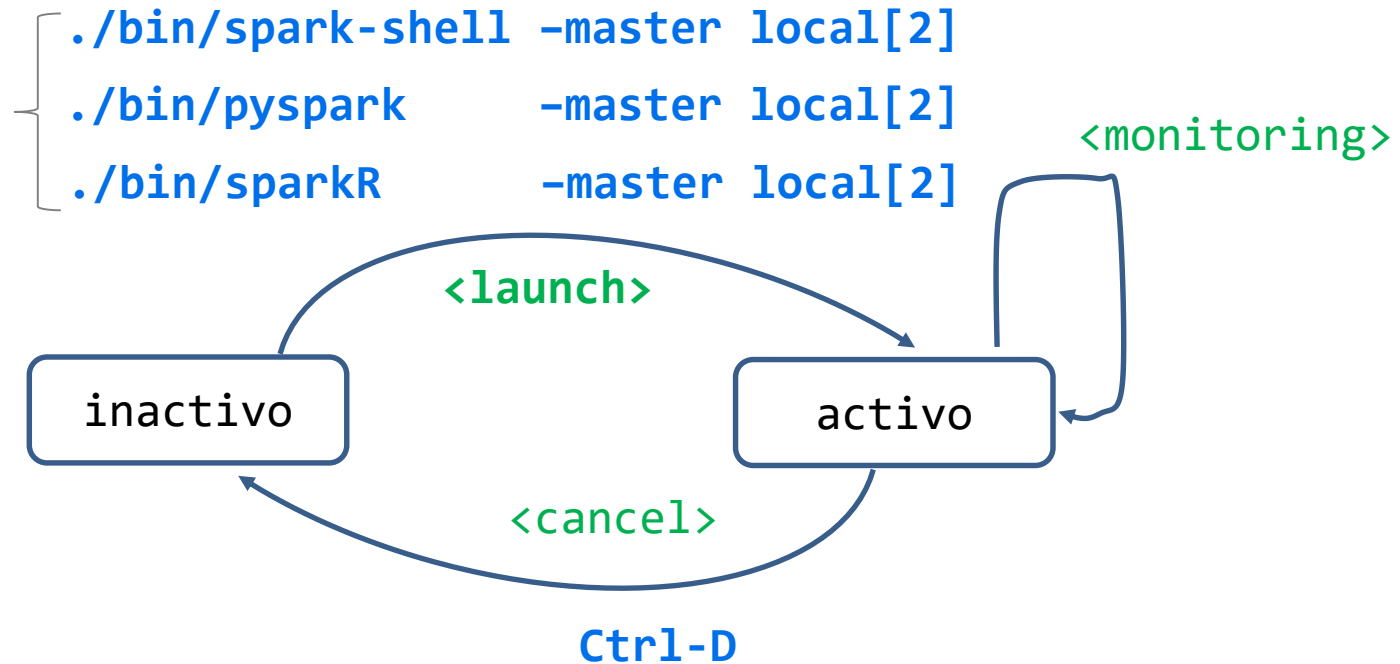
# Spark: nodo autónomo

**shell-interactivo**          submit          libro-interactivo

```
local        -> 1 thread
local[N]     -> N threads
local[*]     -> as many threads as cores are
```

```
./bin/spark-shell  –master local[2]

./bin/pyspark      –master local[2]

./bin/sparkR       –master local[2]
```

**<monitoring>**

**<launch>**

inactivo          activo

**<cancel>**

**Ctrl-D**

# Spark: nodo autónomo

## shell-interactivo          submit          libro-interactivo

```
acaldero@h1:~$ ./bin/pyspark
Python 2.7.13 (default, Jan 19 2017, 14:48:08)
[GCC 6.3.0 20170118] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
17/10/17 01:08:04 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using
    builtin-java classes where applicable
17/10/17 01:08:12 WARN ObjectStore: Version information not found in metastore.
    hive.metastore.schema.verification is not enabled so recording the schema version 1.2.0
17/10/17 01:08:12 WARN ObjectStore: Failed to get database default, returning NoSuchObjectException
17/10/17 01:08:13 WARN ObjectStore: Failed to get database global_temp, returning NoSuchObjectException
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 2.2.0
      /_/

Using Python version 2.7.13 (default, Jan 19 2017 14:48:08)
SparkSession available as 'spark'.
>>>
```

# Spark: nodo autónomo

**shell-interactivo**                submit                libro-interactivo

```
Using Python version 2.7.13 (default, Jan 19 2017 14:48:08)
SparkSession available as 'spark'.
>>> import sys
>>> from random import random
>>> from operator import add
>>> from pyspark.sql import SparkSession
>>>
>>> partitions = 2
>>> n = 100000 * partitions
>>> def f(_):
...     x = random() * 2 - 1
...     y = random() * 2 - 1
...     return 1 if x ** 2 + y ** 2 < 1 else 0
...
>>> spark = SparkSession.builder.appName("PythonPi").getOrCreate()
>>> count = spark.sparkContext.parallelize(range(1, n + 1), partitions).map(f).reduce(add)
16/11/27 14:08:13 WARN TaskSetManager: Stage 0 contains a task of very large size (368 KB). The maximum
    recommended task size is 100 KB.
>>> print("Pi is roughly %f" % (4.0 * count / n))
Pi is roughly 3.139500
>>> spark.stop()
>>>
```
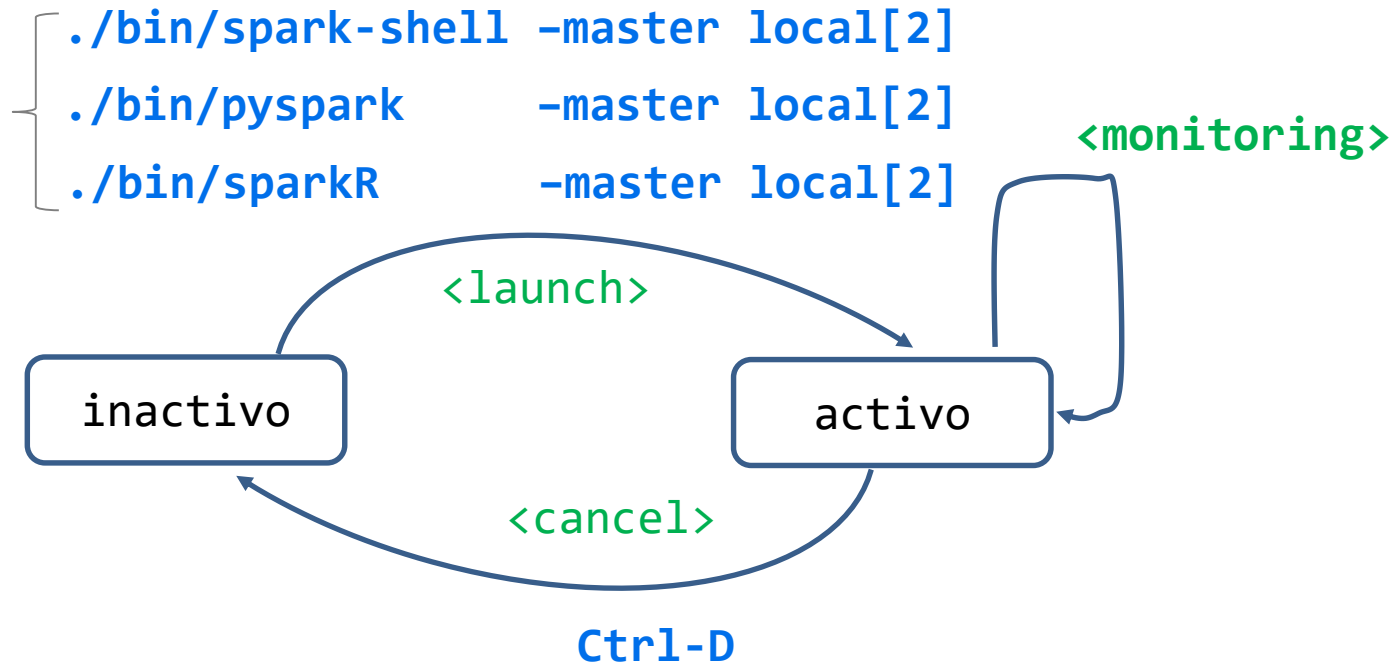
# Spark: nodo autónomo

**shell-interactivo**          submit          libro-interactivo

```
./bin/spark-shell –master local[2]
./bin/pyspark     –master local[2]
./bin/sparkR      –master local[2]
```

<monitoring>

<launch>

inactivo

activo

<cancel>

Ctrl-D

# Spark: nodo autónomo

**shell-interactivo**          submit          libro-interactivo

```
Using Python version 2.7.9 (default, Jun 29 2016 13:08:31)
SparkSession available as 'spark'.
>>> import sys
>>> from random import random
>>> from operator import add
>>> from pyspark.sql import SparkSession
>>>
>>> partitions = 2
>>> n = 100000 * partitions
>>> def f(_):
...     x = random() * 2 - 1
...     y = random() * 2 - 1
...     return 1 if x ** 2 + y ** 2 < 1 else 0
...
>>> spark = SparkSession.builder.appName("PythonPi").getOrCreate()
>>> count = spark.sparkContext.parallelize(range(1, n + 1), partitions).map(f).reduce(add)
16/11/27 14:08:13 WARN TaskSetManager: Stage 0 contains a task of very large size (368 KB).
    The maximum recommended task size is 100 KB.
>>> print("Pi is roughly %f" % (4.0 * count / n))
Pi is roughly 3.139500
>>> spark.stop()
>>>
```

http://<ip>:4040
http://<ip>:4041
...

http://spark.apache.org/docs/latest/monitoring.html

# Spark: nodo autónomo

**shell-interactivo**          submit          libro-interactivo

# Spark: nodo autónomo

**shell-interactivo**          submit          libro-interactivo

# Spark: nodo autónomo

**shell-interactivo**      submit      libro-interactivo

# Spark: nodo autónomo

shell-interactivo **submit** libro-interactivo

```
./bin/spark-submit --master local[8] \
  --class org.apache.spark.examples.SparkPi \
  ./examples/jars/spark-examples_2.11-2.2.0.jar \
5
```

<monitoring>

<launch>

inactivo

activo

<cancel>

Ctrl-D

# Spark: nodo autónomo

shell-interactivo          submit          **libro-interactivo**

```
acaldero@h1:~$ mkdir work
acaldero@h1:~$ cd work
acaldero@h1:~$ wget http://www.gutenberg.org/cache/epub/2000/pg2000.txt


acaldero@h1:~$ pyspark
[TerminalIPythonApp] WARNING | Subcommand `ipython notebook` is deprecated and will be removed in future versions.
[TerminalIPythonApp] WARNING | You likely want to use `jupyter notebook` in the future
[I 18:48:14.980 NotebookApp] [nb_conda_kernels] enabled, 2 kernels found
[I 18:48:15.016 NotebookApp] ✓ nbpresent HTML export ENABLED
[W 18:48:15.016 NotebookApp] ✗ nbpresent PDF export DISABLED: No module named nbbrowserpdf.exporters.pdf
[I 18:48:15.018 NotebookApp] [nb_conda] enabled
…
```

https://gist.github.com/tommycarpi/f5a67c66a8f2170e263c

# Spark: nodo autónomo

shell-interactivo          submit          **libro-interactivo**

```
acaldero@h1:~$ firefox http://localhost:8888/
ps# sc + <shift + enter>
```

# Spark: nodo autónomo

shell-interactivo          submit          **libro-interactivo**

```
myRDD = sc.textFile("file:///home/acaldero/work/pg2000.txt")
words = myRDD.flatMap(lambda line : line.split(" ")).map(lambda word : (word,
    1)).reduceByKey(lambda a, b : a + b)
words.saveAsTextFile("file:///home/acaldero/work/pg2000-wc")
```

# Spark: nodo autónomo

shell-interactivo          submit          **libro-interactivo**

```
myRDD = sc.textFile("file:///home/acaldero/work/pg2000.txt")
words = myRDD.flatMap(lambda line : line.split(" ")).map(lambda word : (word,
    1)).reduceByKey(lambda a, b : a + b)
words.takeOrdered(10, key=lambda x: -x[1])
```



http://stackoverflow.com/questions/24656696/spark-get-collection-sorted-by-value

# Contenidos

- Introducción
- ***Hand-on***
  - Pre-requisitos e instalación
  - Nodo autónomo
  - **Cluster**
- *Benchmarking*

# Spark: cluster privado

**Prerequisitos**    Instalación    Uso básico

Allocates resources
(cores + memory)

Client

Submit App
(mode=cluster)

Driver    Executors    Executors    Executors    Application

Spark
Master    Spark
Worker    Spark
Worker    Spark
Worker    Spark
Worker    Spark

# Spark: cluster privado

Prerequisitos   **Instalación**   Uso básico

Allocates resources
(cores + memory)

Client

Submit App
(mode=cluster)

| Driver | Executors | Executors | Executors | Application |

| Spark Master | Spark Worker | Spark Worker | Spark Worker | Spark Worker | Spark |

```
acaldero@h1:~$ echo "127.0.0.1  master1" >> /etc/hosts
acaldero@h1:~$ echo "127.0.0.1  slave1"  >> /etc/hosts
acaldero@h1:~$ echo "127.0.0.1  slave2"  >> /etc/hosts
```

http://spark.apache.org/docs/latest/spark-standalone.html

# Spark: cluster privado

Prerequisitos   **Instalación**   Uso básico

Allocates resources
(cores + memory)

Client

Submit App
(mode=cluster)

| Driver | Executors | Executors | Executors | Application |

| Spark Master | Spark Worker | Spark Worker | Spark Worker | Spark Worker | Spark |

```
acaldero@h1:~$ echo "node1" >> spark/conf/slaves
acaldero@h1:~$ echo "node2" >> spark/conf/slaves

acaldero@h1:~$ : Spark en todos los nodos (si fuera necesario)
acaldero@h1:~$ scp –r spark acaldero@node1:~/
…
```

http://spark.apache.org/docs/latest/spark-standalone.html

# Spark: cluster privado

Prerequisitos     **Instalación**     Uso básico

```
acaldero@h1:/home/acaldero$ ssh-keygen -t rsa -P ""
Generating public/private rsa key pair.
Enter file in which to save the key (/home/acaldero/.ssh/id_rsa):
Created directory '/home/acaldero/.ssh'.
Your identification has been saved in /home/acaldero/.ssh/id_rsa.
Your public key has been saved in /home/acaldero/.ssh/id_rsa.pub.
The key fingerprint is:
f0:14:95:a1:0b:78:57:0b:c7:65:47:43:39:b2:2f:8a acaldero@ws1
The key's randomart image is:
+---[RSA 2048]----+
|        oo=+oo=. |
|     .   *oo..o. |
…
```

# Spark: cluster privado

Prerequisitos    **Instalación**    Uso básico

```
acaldero@h1:/home/acaldero$ scp .ssh/id_rsa.pub acaldero@node1:~/.ssh/authorized_keys
Password:

…


acaldero@h1:/home/acaldero$ ssh node1
The authenticity of host 'localhost (::1)' can't be established.
ECDSA key fingerprint is bb:85:4c:6a:ff:e4:34:f8:ac:82:bf:56:a6:79:d8:80.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'localhost' (ECDSA) to the list of known hosts.

…


acaldero@node1:~$ exit
logout
```

# Spark: cluster privado

Prerequisitos    Instalación    **Uso básico**

Allocates resources
(cores + memory)

Client

Submit App
(mode=cluster)

**Driver**  **Executors**  **Executors**  **Executors**    Application

Spark
Master  Spark
Worker  Spark
Worker  Spark
Worker  Spark
Worker    Spark

```
acaldero@h1:~$ : Ir al nodo master
acaldero@h1:~$ ssh acaldero@master1
acaldero@master1:~$ ./spark/sbin/start-all.sh
localhost: starting org.apache.spark.deploy.worker.Worker, logging to /home/acaldero/spark/logs/spark-acaldero-org.apache.spark.deploy.worker.Worker-1-ws1.out
localhost: starting org.apache.spark.deploy.worker.Worker, logging to /home/acaldero/spark/logs/spark-acaldero-org.apache.spark.deploy.worker.Worker-1-ws1.out
localhost: starting org.apache.spark.deploy.worker.Worker, logging to /home/acaldero/spark/logs/spark-acaldero-org.apache.spark.deploy.worker.Worker-1-ws1.out
…
```

http://spark.apache.org/docs/latest/spark-standalone.html

# Spark: cluster privado

Prerequisitos      Instalación      **Uso básico**

Allocates resources
(cores + memory)

Client

Submit App
(mode=cluster)

| Driver | Executors | Executors | Executors | Application |

| Spark Master | Spark Worker | Spark Worker | Spark Worker | Spark Worker | Spark |

```
acaldero@master1:~$ ./spark/sbin/stop-all.sh
acaldero@master1:~$ exit
acaldero@h1:~$ : Regresar al cliente
localhost: stopping org.apache.spark.deploy.worker.Worker
localhost: stopping org.apache.spark.deploy.worker.Worker
localhost: stopping org.apache.spark.deploy.worker.Worker
stopping org.apache.spark.deploy.master.Master
```

http://spark.apache.org/docs/latest/spark-standalone.html

# Spark: cluster privado

Prerequisitos     Instalación     **Uso básico**



**\<submit\>**

**\<monitoring\>**

**$SPARK_HOME/sbin/start-all.sh**

inactivo     activo

**$SPARK_HOME/sbin/stop-all.sh**

# Spark: cluster privado

Prerequisitos        Instalación        **Uso básico**

**Allocates resources**

(cores + memory)

Client

**Submit App**

(mode=cluster)

| Driver | Executors | Executors | Executors | Application |

| Spark Master | Spark Worker | Spark Worker | Spark Worker | Spark Worker | Spark |

```
acaldero@h1:~$ ./spark/bin/spark-shell --master spark://master1:7077
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel).
16/11/27 23:13:55 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform...

…
scala> exit
```
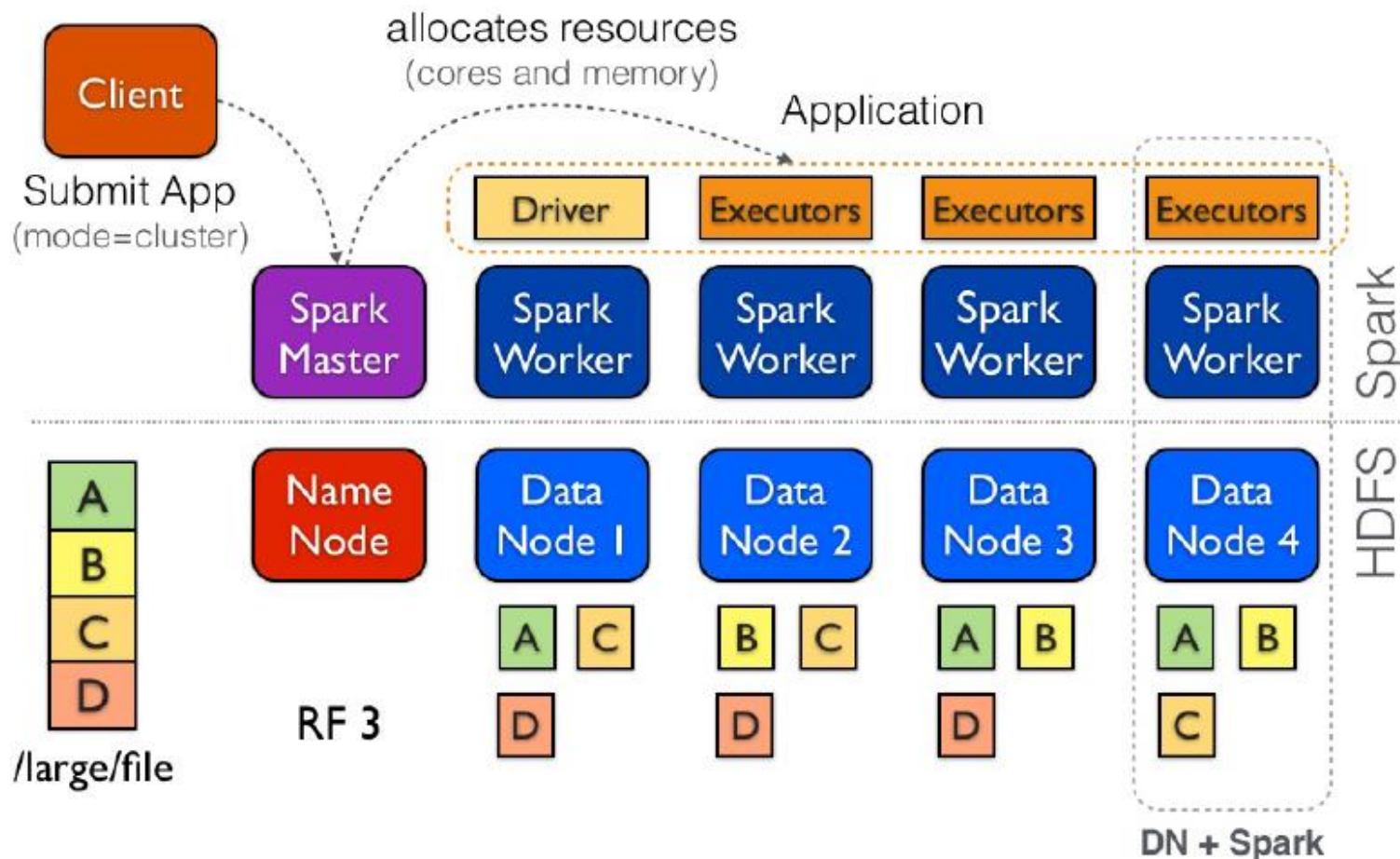
http://spark.apache.org/docs/latest/spark-standalone.html

# Spark: cluster privado

Prerequisitos    **Instalación**    Uso básico



https://trongkhoanguyenblog.wordpress.com/2014/11/23/a-gentle-introduction-to-apache-spark/

# Contenidos

– Introducción

– *Hand-on*

- Pre-requisitos e instalación

- Nodo autónomo

- Cluster

– ***Benchmarking***

http://Spark.apache.org

# Benchmarking

- HiBench
  - https://github.com/intel-hadoop/HiBench

- Spark-perf
  - https://github.com/databricks/spark-perf

http://01org.github.io/sparkscore/about.html

# Benchmarking

- TeraSort
  - Elevada entrada y salida, y comunicación intermedia

- WordCount, PageRank
  - Contar referencias de palabras, enlaces, etc.

- SQL
  - Scan, Join, Aggregate
  - …

- Machine Learning
  - Bayesian Classification
  - K-means clustering
  - …

https://www.oreilly.com/ideas/investigating-sparks-performance

# TeraSort (2014)

|  | **Hadoop World Record** | **Spark 100 TB** | **Spark 1 PB** |
|---|---|---|---|
| Data Size | 102.5 TB | 100 TB | 1000 TB |
| Elapsed Time | 72 mins | 23 mins | 234 mins |
| # Nodes | 2100 | 206 | 190 |
| # Cores | 50400 | 6592 | 6080 |
| # Reducers | 10,000 | 29,000 | 250,000 |
| Rate | 1.42 TB/min | 4.27 TB/min | 4.27 TB/min |
| Rate/node | 0.67 GB/min | 20.7 GB/min | 22.5 GB/min |
| Sort Benchmark Daytona Rules | Yes | Yes | No |
| Environment | dedicated data center | EC2 (i2.8xlarge) | EC2 (i2.8xlarge) |

https://gigaom.com/2014/10/10/databricks-demolishes-big-data-benchmark-to-prove-spark-is-fast-on-disk-too/

# Bibliografía: tutoriales

- Página Web oficial:
  - http://spark.apache.org/

- Introducción a cómo funciona Spark:
  - http://spark.apache.org/docs/latest/quick-start.html

- Tutorial de cómo instalar y usar Spark:
  - http://spark.apache.org/docs/latest/index.html
  - http://spark.apache.org/docs/latest/configuration.html

# Bibliografía: libro

- Learning Spark, Advanced Analytics with Spark:
  - http://shop.oreilly.com/product/0636920028512.do
  - http://shop.oreilly.com/product/0636920035091.do

# Agradecimientos

- Por último pero no por ello menos importante, agradecer al personal del
  Laboratorio del Departamento de Informática
  todos los comentarios y sugerencias para esta presentación.

# Diseño de Sistemas Distribuidos

Máster en Ciencia y Tecnología Informática
Curso 2018-2019

# Sistemas escalables en entornos distribuidos.
# Introducción a Spark

Alejandro Calderón Mateos  &  Jaime Pons Bailly-Bailliere

acaldero@inf.uc3m.es          jaime@lab.inf.uc3m.es