

CS 696 Intro to Big Data: Tools and Methods
Fall Semester, 2019
Doc 17 Spark on AWS
Mar 28, 2019

Copyright ©, All rights reserved. 2019 SDSU & Roger Whitney,
5500 Campanile Drive, San Diego, CA 92182-7700 USA.
OpenContent (<http://www.opencontent.org/opl.shtml>) license
defines the copyright on this document.

Running Program on AWS EMR

Make sure program runs locally

Create program file containing code

Create s3 bucket(s) for

- program file

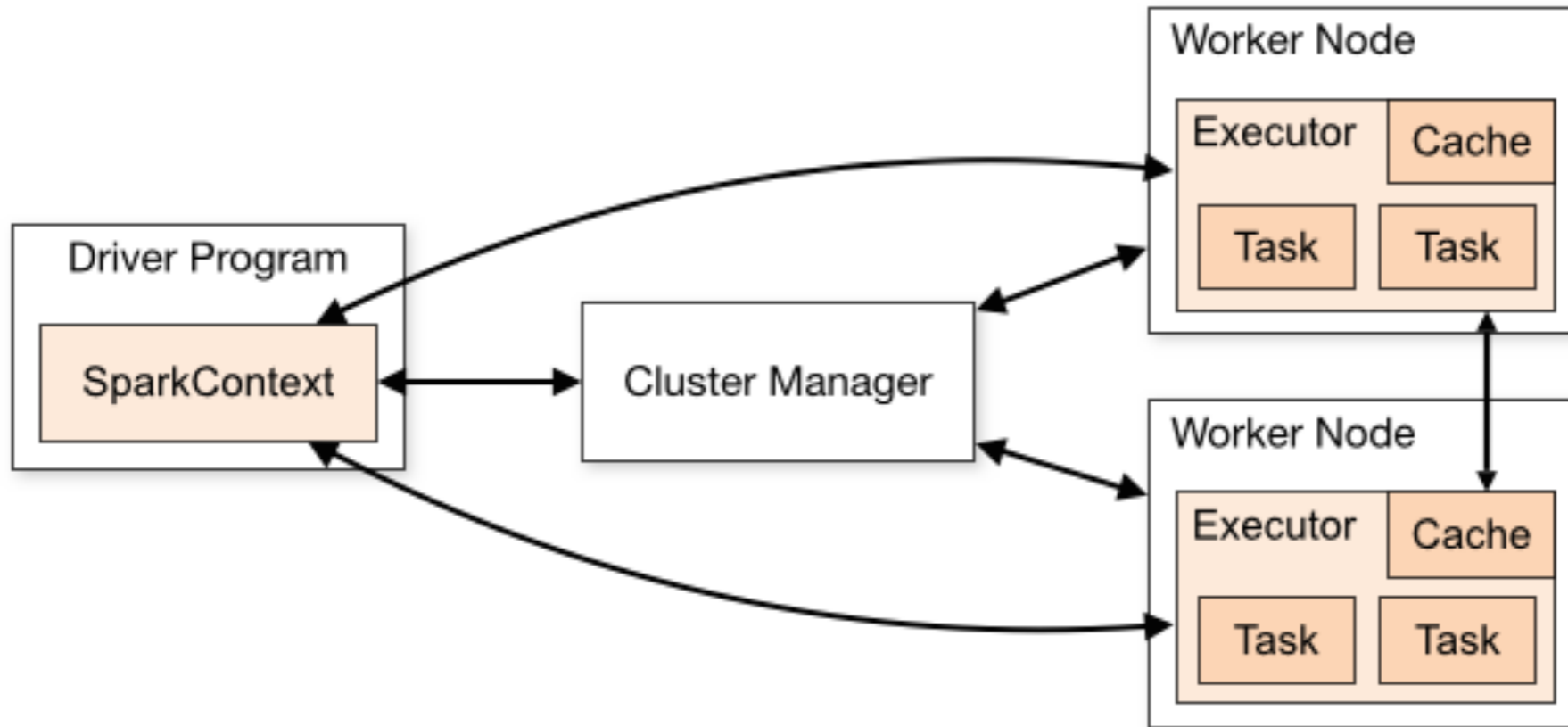
- logs

- input

- output

Upload program & data files to s3

Spark Components



Terms

Application

User program built on Spark

Driver program + executors

Driver program

The process running the main() function of the application and creating the SparkContext

Cluster manager

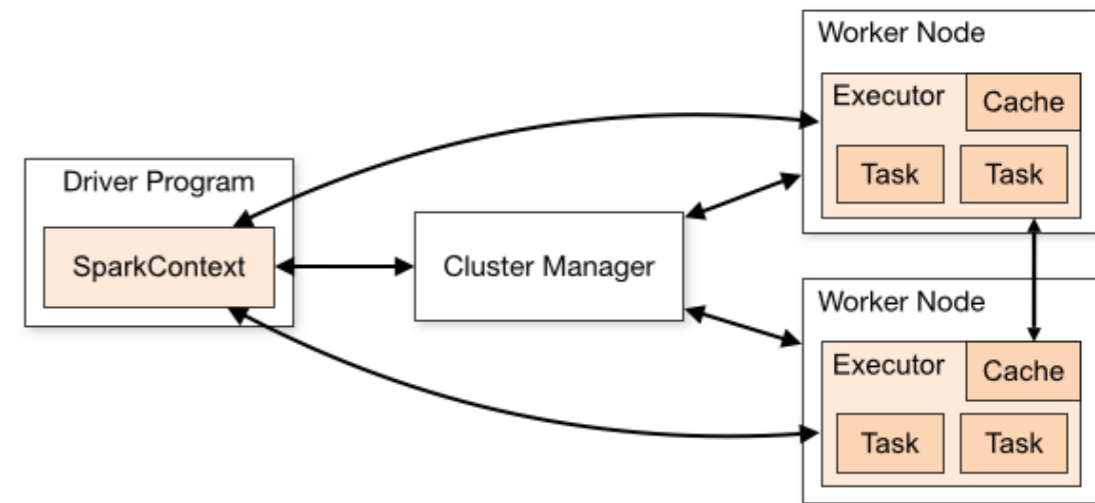
External service for acquiring resources on the cluster

Deploy mode

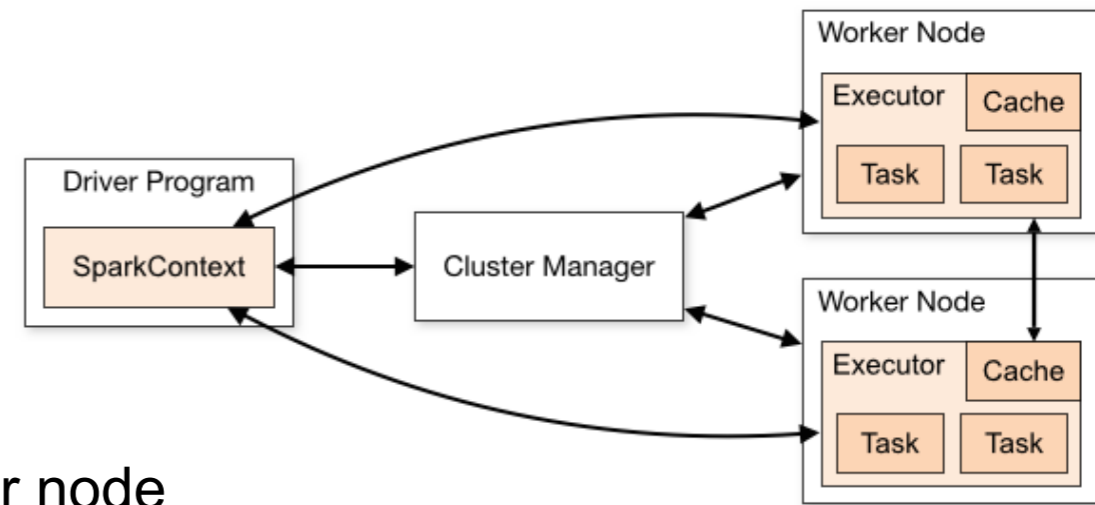
Where the driver process runs

"cluster" - the driver inside of the cluster

"client" - the driver outside of the cluster



Terms



Executor

A process launched for an application on a worker node

Runs tasks and keeps data in memory or disk storage across them.

Each application has its own executors

Task

A unit of work that will be sent to one executor

Job

A parallel computation consisting of multiple tasks

Gets spawned in response to a Spark action (e.g. save, collect)

Stage

Job divided into smaller tasks called stages

Depend on each other

Test Program 1 - Pi

```
from random import random
from operator import add
```

```
from pyspark import SparkContext
```

```
if __name__ == "__main__":
    sc = SparkContext(appName="PythonPi")
    partitions = 3
    n = 100000 * partitions
```

```
def f(_):
    x = random() * 2 - 1
    y = random() * 2 - 1
    return 1 if x ** 2 + y ** 2 < 1 else 0
```












```
count = sc.parallelize(range(1, n + 1), partitions).map(f).reduce(add)
print("Pi is roughly %f" % (4.0 * count / n))
```

```
sc.stop()
```

Designed to have no
Command line dependancies
No input or output files

My S3 Buckets



▶	 aws-emr-resources-834365227482-us-west-2	•	
▶	 aws-logs-834365227482-us-west-2	•	
▼	 rw-696-flight	•	3/27/19
	▶  2015-summary.json	21 KB	3/26/19
	▶  flight.py	2 KB	3/26/19
	▶  pi.py	49...tes	3/27/19
▶	 rw-cs696-data	•	
▶	 rw-hadoop-jars	•	
▶	 rw-hadoop-logs	•	
▼	 rw-output-data	•	3/27/19
▶	 rw-wc-input-data	•	
▶	 rw-wc-output-data	•	

Spark on AWS - EMR Console

Create cluster

View details

Clone

Terminate

Filter: All clusters 5 clusters (all loaded)

	Name	ID	Status	Creation time (UTC-7)	Elapsed time	Normalized instance h
<input type="checkbox"/>	▶ My cluster	j-10UWFTGNQOV0T	Terminated User request	2019-03-26 16:44 (UTC-7)	18 minutes	16
<input type="checkbox"/>	▶ NotebookCluster	j-3TT5GQUPDXJE	Terminated User request	2019-03-26 14:49 (UTC-7)	6 minutes	8
<input type="checkbox"/>	▶ NotebookCluster	j-37KDDYAF5P3Y3	Terminated User request	2019-03-26 14:47 (UTC-7)	8 minutes	8
<input type="checkbox"/>	▶ My cluster	j-32U2OW9DKADLG	Terminated User request	2019-03-26 10:36 (UTC-7)	15 minutes	24
<input type="checkbox"/>	▶ FlightExample	j-32NKDKNYTF1JE	Terminated User request	2019-03-26 10:39 (UTC-7)	1 minute	0

Using Quick Options

General Configuration

Cluster name

Logging ⓘ

S3 folder

Launch mode Cluster ⓘ Step execution ⓘ

Software configuration

Release ⓘ

Applications

- Core Hadoop: Hadoop 2.8.5 with Ganglia 3.7.2, Hive 2.3.4, Hue 4.3.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.1
- HBase: HBase 1.4.9 with Ganglia 3.7.2, Hadoop 2.8.5, Hive 2.3.4, Hue 4.3.0, Phoenix 4.14.1, and ZooKeeper 3.4.13
- Presto: Presto 0.214 with Hadoop 2.8.5 HDFS and Hive 2.3.4 Metastore
- Spark: Spark 2.4.0 on Hadoop 2.8.5 YARN with Ganglia 3.7.2 and Zeppelin 0.8.1 ←
- Use AWS Glue Data Catalog for table metadata ⓘ

Hardware configuration

Instance type ⓘ ←

Number of instances (1 master and 2 core nodes)

Security and access

EC2 key pair ⓘ Learn how to create an EC2 key pair ←

Permissions Default Custom

Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

EMR role [EMR_DefaultRole](#) ⓘ

EC2 instance profile [EMR_EC2_DefaultRole](#) ⓘ

←

Use Advanced Options

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name

Logging ⓘ

S3 folder ⓘ

Launch mode **Cluster** ⓘ **Step execution** ⓘ

Software configuration

Release ⓘ

- Applications**
- Core Hadoop: Hadoop 2.8.5 with Ganglia 3.7.2, Hive 2.3.4, Hue 4.3.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.1
 - HBase: HBase 1.4.9 with Ganglia 3.7.2, Hadoop 2.8.5, Hive 2.3.4, Hue 4.3.0, Phoenix 4.14.1, and ZooKeeper 3.4.13
 - Presto: Presto 0.214 with Hadoop 2.8.5 HDFS and Hive 2.3.4 Metastore
 - Spark: Spark 2.4.0 on Hadoop 2.8.5 YARN with Ganglia 3.7.2 and Zeppelin 0.8.1

Use AWS Glue Data Catalog for table metadata ⓘ

Advanced Options

Software Configuration

Release  

- | | | |
|--------------------------------------------------|-----------------------------------------|--------------------------------------------|
| <input checked="" type="checkbox"/> Hadoop 2.8.5 | <input type="checkbox"/> Zeppelin 0.8.1 | <input type="checkbox"/> Livy 0.5.0 |
| <input type="checkbox"/> JupyterHub 0.9.4 | <input type="checkbox"/> Tez 0.9.1 | <input type="checkbox"/> Flink 1.7.1 |
| <input type="checkbox"/> Ganglia 3.7.2 | <input type="checkbox"/> HBase 1.4.9 | <input type="checkbox"/> Pig 0.17.0 |
| <input type="checkbox"/> Hive 2.3.4 | <input type="checkbox"/> Presto 0.215 | <input type="checkbox"/> ZooKeeper 3.4.13 |
| <input type="checkbox"/> MXNet 1.3.1 | <input type="checkbox"/> Sqoop 1.4.7 | <input type="checkbox"/> Mahout 0.13.0 |
| <input type="checkbox"/> Hue 4.3.0 | <input type="checkbox"/> Phoenix 4.14.1 | <input type="checkbox"/> Oozie 5.1.0 |
| <input checked="" type="checkbox"/> Spark 2.4.0 | <input type="checkbox"/> HCatalog 2.3.4 | <input type="checkbox"/> TensorFlow 1.12.0 |

AWS Glue Data Catalog settings (optional)


Use for Spark table metadata 

Edit software settings

Enter configuration Load JSON from S3

`classification=config-file-name,properties=[myKey1=myValue1,myKey2=myValue2]`

Add steps (optional)

Step type 

Auto-terminate cluster after the last step is completed

Advanced Options

Hardware Configuration ⓘ

If you need more than 20 EC2 instances, see [this topic](#).

- Instance group configuration
- Uniform instance groups**
Specify a single instance type and purchasing option for each node type.
 - Instance fleets**
Specify target capacity and how Amazon EMR fulfills it for each node type. Mix instance types and purchasing options. [Learn more](#)

Network [Create a VPC](#) ⓘ

EC2 Subnet

Root device EBS volume size GiB ⓘ

Choose the instance type, number of instances, and a purchasing option. You can choose to use On-Demand instances, Spot instances, or both. The instance type and purchasing option apply to all EC2 instances in each instance group, and you can only specify these options for an instance group when you create it. [Learn more about instance purchasing options](#)

Node type	Instance type	Instance count	Purchasing option	Auto Scaling
Master Master - 1	m3.xlarge 8 vCore, 15 GiB memory, 80 SSD GB storage EBS Storage: none Add configuration settings	1 Instances	<input checked="" type="radio"/> On-demand ⓘ <input type="radio"/> Spot ⓘ <input type="text" value="Use on-demand as max price"/>	Not available for Master
Core Core - 2	m3.xlarge 8 vCore, 15 GiB memory, 80 SSD GB storage EBS Storage: none Add configuration settings	2 Instances	<input checked="" type="radio"/> On-demand ⓘ <input type="radio"/> Spot ⓘ <input type="text" value="Use on-demand as max price"/>	Not enabled
Task ✕ Task - 3	m3.xlarge 8 vCore, 15 GiB memory, 80 SSD GB storage EBS Storage: none Add configuration settings	3 Instances	<input checked="" type="radio"/> On-demand ⓘ <input type="radio"/> Spot ⓘ <input type="text" value="Use on-demand as max price"/>	Not enabled

[+ Add task instance group](#)

Cancel

Previous

Next

Advanced Options

General Options

Cluster name ←

Logging ⓘ

S3 folder ←

Debugging ⓘ

Termination protection ⓘ

Tags ⓘ

Key	Value (optional)
<input type="text" value="Add a key to create a tag"/>	<input type="text"/>

Additional Options

EMRFS consistent view ⓘ

Custom AMI ID ⓘ

▶ Bootstrap Actions



Cancel


Previous


Next

Advanced Options

Security Options

EC2 key pair Proceed without an EC2 key pair   ←

Cluster visible to all IAM users in account 

 This feature will be deprecated soon.

Permissions

Default Custom

Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.


EMR role [EMR_DefaultRole](#)  

EC2 instance profile [EMR_EC2_DefaultRole](#)  

Auto Scaling role [EMR_AutoScaling_DefaultRole](#)  

▶ Authentication and encryption

▶ EC2 security groups

 No EC2 key pair has been selected, so you will not be able to SSH to this cluster. [Learn how to create an EC2 Key Pair.](#)

Cancel

Previous

Create cluster



Cluster Created - Either Quick or Advanced

Amazon EMR

- Clusters
- Security configurations
- VPC subnets
- Events
- Notebooks
- Help
- What's new

Clone Terminate AWS CLI export

Cluster: CourseSlides **Starting**

Summary Application history Monitoring Hardware Configurations Events Steps Bootstrap actions

Connections: --

Master public DNS: --

Tags: -- [View All / Edit](#)

Summary

ID: j-U3D2ZYQ3E304

Creation date: 2019-03-28 10:01 (UTC-7)

Elapsed time: 30 seconds

Auto-terminate: No

Termination Protection: [Change](#)

Configuration details

Release label: emr-5.22.0

Hadoop distribution: Amazon 2.8.5

Applications: Spark 2.4.0

Log URI: s3://aws-logs-834365227482-us-west-2/elasticmapreduce/

EMRFS consistent view: Disabled

Custom AMI ID: --

Network and hardware

Availability zone: us-west-2a

Subnet ID: [subnet-0f55196b](#)

Master: Provisioning 1 m3.xlarge

Core: Provisioning 2 m3.xlarge

Task: --

Adding a Step

The screenshot shows the 'Steps' page in the AWS EMR console for a cluster named 'CourseSlides' which is in a 'Starting' state. At the top, there are buttons for 'Clone', 'Terminate', and 'AWS CLI export'. Below these are navigation tabs: 'Summary', 'Application history', 'Monitoring', 'Hardware', 'Configurations', 'Events', 'Steps', and 'Bootstrap actions'. The 'Steps' tab is selected, and an arrow points to it from above. Under the 'Steps' tab, there are buttons for 'Add step', 'Clone step', and 'Cancel step'. Below these buttons, there are links for 'View all interactive jobs' and 'View all jobs'. A filter bar shows 'Filter: All steps' and '1 step (all loaded)'. Below the filter is a table with the following data:

	ID	Name	Status	Start time (UTC-7)	Elapsed time	Log files	Actions
<input type="radio"/>	s-KSWMDHULHSD8	Setup hadoop debugging	Pending		--	View logs	View jobs

Add step ✕

Step type

Name

Deploy mode Run your driver on a slave node (cluster mode) or on the master node as an external client (client mode).

Spark-submit options Specify other options for spark-submit.

Application location* Path to a JAR with your application and dependencies (client deploy mode only supports a local path).

Arguments Specify optional arguments for your application.

Action on failure What to do if the step fails.

Cancel Add

Modes - client or cluster

Either works

client mode gives access to standard out

Clone Terminate AWS CLI export

Cluster: CourseSlides **Starting** Configuring cluster software

Summary Application history Monitoring Hardware Configurations Events Steps Bootstrap actions

Add step Clone step Cancel step

Steps

[View all interactive jobs](#) | [View all jobs](#)

Filter: All steps Filter steps ... 2 steps (all loaded)

	ID	Name	Status	Start time (UTC-7) ▼	Elapsed time	Log files	Actions
<input type="radio"/>	s-3Q0ZVOAZV3VR	Spark application	Pending		–	View logs	View jobs
<input type="radio"/>	s-KSWMDHULHSD8	Setup hadoop debugging	Pending		–	View logs	View jobs



Cluster: CourseSlides **Waiting** Cluster ready after last step completed.

Summary Application history Monitoring Hardware Configurations Events Steps Bootstrap actions

Add step Clone step Cancel step

Steps

[View all interactive jobs](#) | [View](#)

Filter: All steps Filter steps ... 2 steps (all loaded)

ID	Name	Status	Start time (UTC-7)	Elapsed time	Log files	Actions
s-3Q0ZVOAZV3VR	Spark application	Completed	2019-03-28 10:05 (UTC-7)	36 seconds	View logs	View jobs
s-KSWMDHULHSD8	Setup hadoop debugging	Completed	2019-03-28 10:05 (UTC-7)	2 seconds	View logs	View jobs

Add step Clone step Cancel step

Steps

[View all interactive jobs](#) | [View](#)

Filter: All steps Filter steps ... 2 steps (all loaded)

ID	Name	Status	Start time (UTC-7)	Elapsed time	Log files	Actions
s-3Q0ZVOAZV3VR	Spark application	Completed	2019-03-28 10:05 (UTC-7)	36 seconds	controller syslog* stderr stdout 	View jobs
s-KSWMDHULHSD8	Setup hadoop debugging	Completed	2019-03-28 10:05 (UTC-7)	2 seconds	View logs	View jobs

```
Pi is roughly 3.138120
```

Example 2

```
def flight(input, output):
    import pyspark.sql.functions as F
    from pyspark.sql import SparkSession
    spark = SparkSession.builder \
        .appName("Flight") \
        .getOrCreate()









    flight_df = spark.read.json(input)

    grouped_df = flight_df.groupBy('DEST_COUNTRY_NAME').agg(F.sum('count'))
    grouped_df.write.format('csv').save(output)

def files_from_args():
    import argparse
    parser = argparse.ArgumentParser()
    parser.add_argument('-i', '--input', default='input')
    parser.add_argument('-o', '--output', default='output')
    args = parser.parse_args()
    return (args.input, args.output)

if __name__ == "__main__":
    inputfile, outputfile = files_from_args()
    flight(inputfile, outputfile)
```

S3 Buckets

▶  aws-emr-resources-834365227482-us-west-2	•		
▶  aws-logs-834365227482-us-west-2	•		
▼  rw-696-flight	•		
 2015-summary.json		21 KB	3/26/19
 flight.py		2 KB	3/26/19
 pi.py		49...tes	3/27/19
▶  rw-cs696-data	•		
▶  rw-hadoop-jars	•		

Added Step

Add step ✕

Step type

Name

Deploy mode Run your driver on a slave node (cluster mode) or on the master node as an external client (client mode).

Spark-submit options Specify other options for spark-submit.

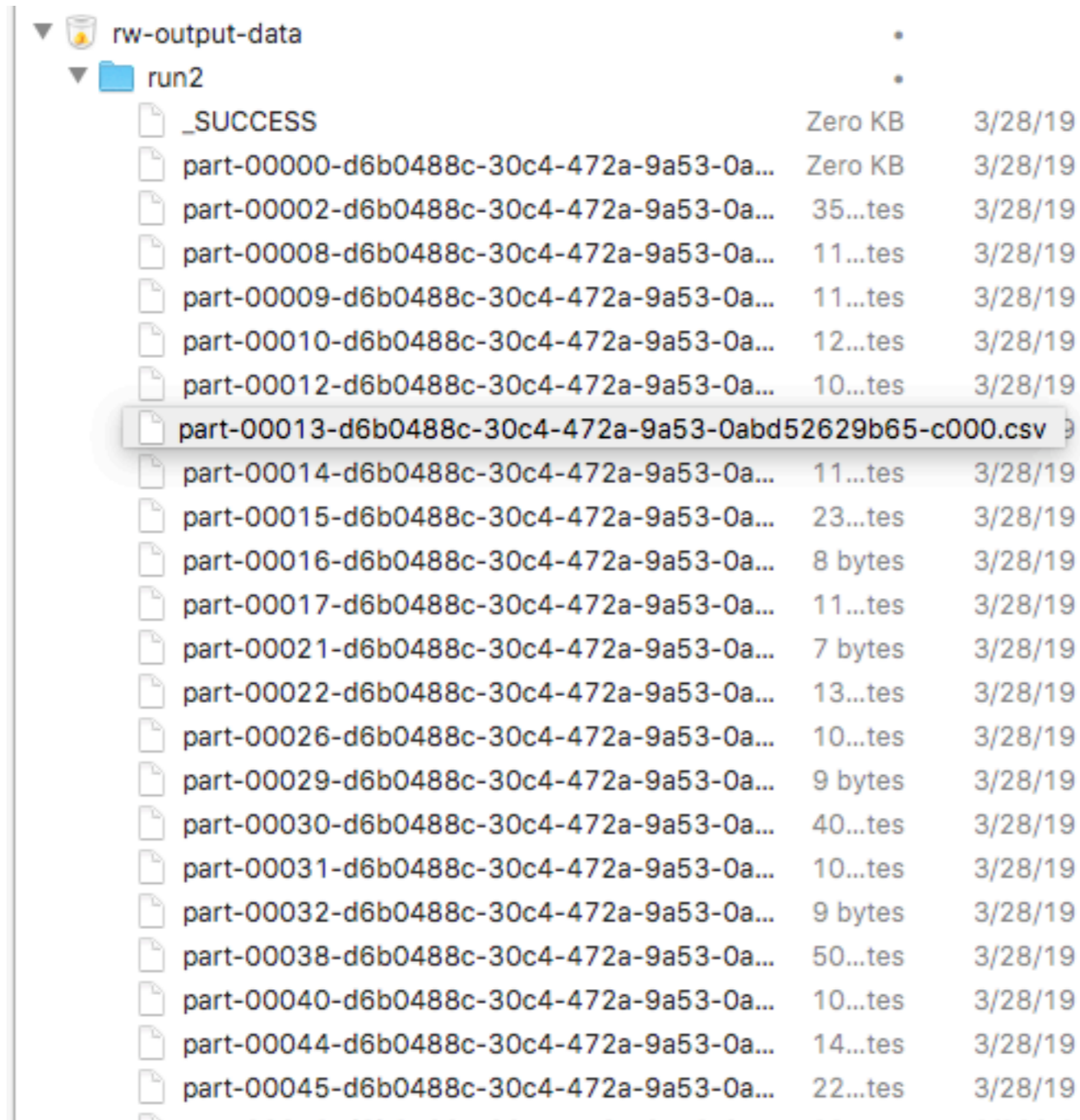
Application location* Path to a JAR with your application and dependencies (client deploy mode only supports a local path).

Arguments Specify optional arguments for your application.

Action on failure What to do if the step fails.

Cancel Add

S3 output



The screenshot shows a file explorer view of an S3 bucket. The root directory is 'rw-output-data', which contains a sub-directory 'run2'. The 'run2' directory lists 20 files. The file 'part-00013-d6b0488c-30c4-472a-9a53-0abd52629b65-c000.csv' is highlighted with a mouse cursor.

File Name	Size	Modified
rw-output-data		
run2		
_SUCCESS	Zero KB	3/28/19
part-00000-d6b0488c-30c4-472a-9a53-0a...	Zero KB	3/28/19
part-00002-d6b0488c-30c4-472a-9a53-0a...	35...tes	3/28/19
part-00008-d6b0488c-30c4-472a-9a53-0a...	11...tes	3/28/19
part-00009-d6b0488c-30c4-472a-9a53-0a...	11...tes	3/28/19
part-00010-d6b0488c-30c4-472a-9a53-0a...	12...tes	3/28/19
part-00012-d6b0488c-30c4-472a-9a53-0a...	10...tes	3/28/19
part-00013-d6b0488c-30c4-472a-9a53-0abd52629b65-c000.csv		
part-00014-d6b0488c-30c4-472a-9a53-0a...	11...tes	3/28/19
part-00015-d6b0488c-30c4-472a-9a53-0a...	23...tes	3/28/19
part-00016-d6b0488c-30c4-472a-9a53-0a...	8 bytes	3/28/19
part-00017-d6b0488c-30c4-472a-9a53-0a...	11...tes	3/28/19
part-00021-d6b0488c-30c4-472a-9a53-0a...	7 bytes	3/28/19
part-00022-d6b0488c-30c4-472a-9a53-0a...	13...tes	3/28/19
part-00026-d6b0488c-30c4-472a-9a53-0a...	10...tes	3/28/19
part-00029-d6b0488c-30c4-472a-9a53-0a...	9 bytes	3/28/19
part-00030-d6b0488c-30c4-472a-9a53-0a...	40...tes	3/28/19
part-00031-d6b0488c-30c4-472a-9a53-0a...	10...tes	3/28/19
part-00032-d6b0488c-30c4-472a-9a53-0a...	9 bytes	3/28/19
part-00038-d6b0488c-30c4-472a-9a53-0a...	50...tes	3/28/19
part-00040-d6b0488c-30c4-472a-9a53-0a...	10...tes	3/28/19
part-00044-d6b0488c-30c4-472a-9a53-0a...	14...tes	3/28/19
part-00045-d6b0488c-30c4-472a-9a53-0a...	22...tes	3/28/19

Warning on AWS

It can take 5-10 minutes to start cluster

Logs do not show your logging statements

When you configure Steps incorrectly they fail
Error messages are not very helpful

SSH to your Master Node

Create Amazon EC2 Key pair

Instructions

<http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ec2-key-pairs.html#having-ec2-create-your-key-pair>

Open EC2 Dashboard - Select Key Pairs

SSH to your Master Node

In Create Cluster - Quick Options

- HIVE 2.3.0 Metastore
- Spark: Spark 2.2.0 on Hadoop 2.7.3 YARN with Ganglia 3.7.2 and Zeppelin 0.7.2
- Use AWS Glue Data Catalog for table metadata ⓘ

Hardware configuration

- Instance type ⌵
- Number of instances (1 master and 2 core nodes)

Security and access

- EC2 key pair ⌵ ⓘ [Learn how to create an EC2 key pair.](#)
- Permissions Default Custom
Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.
- EMR role [EMR_DefaultRole](#) ⓘ
- EC2 instance profile [EMR_EC2_DefaultRole](#) ⓘ

Cancel

Create cluster

SSH to your Master Node

[Clone](#) [Terminate](#) [AWS CLI export](#)

Cluster: LargeRun **Waiting** Cluster ready after last step completed.

[Summary](#) [Application history](#) [Monitoring](#) [Hardware](#) [Events](#) [Steps](#) [Configurations](#) [Bootstrap actions](#)

Connections: [Enable Web Connection](#) – Zeppelin, Spark History Server, Ganglia, Resource Manager ... (View All)
Master public DNS: ec2-34-211-49-97.us-west-2.compute.amazonaws.com [SSH](#)
Tags: -- [View All / Edit](#)

Click for Instructions

Command-line Tools

Flintrock

Open-source command-line tool for launching Apache Spark clusters

<https://github.com/nchammas/flintrock>

aws cli

Amazon's command line tool

<https://aws.amazon.com/cli/>

Generating the Command Line

[Clone](#) [Terminate](#) [AWS CLI export](#)

Cluster: CourseSlides **Starting** Configuring cluster software

[Summary](#) [Application history](#) [Monitoring](#) [Hardware](#) [Configurations](#) [Events](#) [Steps](#) [Bootstrap actions](#)

[Add step](#) [Clone step](#) [Cancel step](#)

Steps

[View all interactive jobs](#) | [View all jobs](#)

Filter: Filter steps ... 2 steps (all loaded) [Refresh](#)

	ID	Name	Status	Start time (UTC-7)	Elapsed time	Log files 🔗	Actions
<input type="radio"/>	s-3Q0ZVOAZV3VR	Spark application	Pending		–	View logs	View jobs
<input type="radio"/>	s-KSWMDHULHSD8	Setup hadoop debugging	Pending		–	View logs	View jobs

AWS CLI export

```
aws emr create-cluster --termination-protected --applications Name=Hadoop Name=Spark --ec2-attributes
'{"InstanceProfile":"EMR_EC2_DefaultRole","SubnetId":"subnet-0f55196b","EmrManagedSlaveSecurityGroup":"sg-
65bffa1c","EmrManagedMasterSecurityGroup":"sg-62bffa1b"}' --release-label emr-5.22.0 --log-uri 's3n://aws-
logs-834365227482-us-west-2/elasticmapreduce/' --steps '[{"Args":["spark-submit","--deploy-
mode","client","s3://rw-696-flight/pi.py"],"Type":"CUSTOM_JAR","ActionOnFailure":"CONTINUE","Jar":"command-
runner.jar","Properties":"","Name":"Spark application"}]' --instance-groups
'[{ "InstanceCount":1,"InstanceGroupType":"MASTER","InstanceType":"m3.xlarge","Name":"Master - 1"},
{ "InstanceCount":2,"InstanceGroupType":"CORE","InstanceType":"m3.xlarge","Name":"Core - 2"}]' --auto-scaling-
role EMR_AutoScaling_DefaultRole --ebs-root-volume-size 10 --service-role EMR_DefaultRole --enable-
debugging --name 'CourseSlides' --scale-down-behavior TERMINATE_AT_TASK_COMPLETION --region us-west-2
```

Hadoop Ecosystem

Hadoop

HDFS

MapReduce

YARN

Tez

Pig

Hive

Hbase

Sqoop

Oozie

Falcon

Spark

ZooKeeper

Mahout

Phoenix

BigTop

+ others

Apache Pig

Programming Map-Reduce can be low level

Apache Pig - high-level platform for creating programs for Hadoop

Pig Latin

```
input_lines = LOAD '/tmp/my-copy-of-all-pages-on-internet' AS (line:chararray);
words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;
filtered_words = FILTER words BY word MATCHES '\\w+';
word_groups = GROUP filtered_words BY word;
word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS
                count, group AS word;

ordered_word_count = ORDER word_count BY count DESC;
STORE ordered_word_count INTO '/tmp/number-of-words-on-internet';
```


Apache Hive

SQL is common way to interact with data

Hive provides SQL like query language for HDFS, Amazon S3 data

HiveQL - converted into MapReduce

```
DROP TABLE IF EXISTS docs;
CREATE TABLE docs (line STRING);
LOAD DATA INPATH 'input_file' OVERWRITE INTO TABLE docs;
CREATE TABLE word_counts AS
SELECT word, count(1) AS count FROM
(SELECT explode(split(line, '\s')) AS word FROM docs) temp
GROUP BY word
ORDER BY word;
```

Apache HBase

BigTable for Hadoop

Non-relational distributed database

Fault-tolerant way of storing large quantities of sparse data

Apache Sqoop

People have data in non-hadoop databases

Sqoop

Transferring data between relational databases & Hadoop

Apache Phoenix

But SQL is common

Phoenix

Massively parallel relational database for Hadoop

Uses HBase to store data

Apache Spark

Hadoop has latency issues - reads data from disk

MapReduce is not conducive to solving all problems

Spark

- Uses distributed shared memory: Resilient distributed dataset (RDD)

- Iterative algorithms

- Implemented in Scala

Spark Core

Spark SQL

- Dataframes & SQL

Spark Streaming

Spark MLlib

- Machine learning

Apache Mahout

Hadoop does not have machine learning libraries

Mahout

Environment for quickly creating scalable machine learning applications

Samsara - R-line syntax & environment

Apache Flink, Apache Storm

Hadoop does batch jobs

Spark streaming has delays

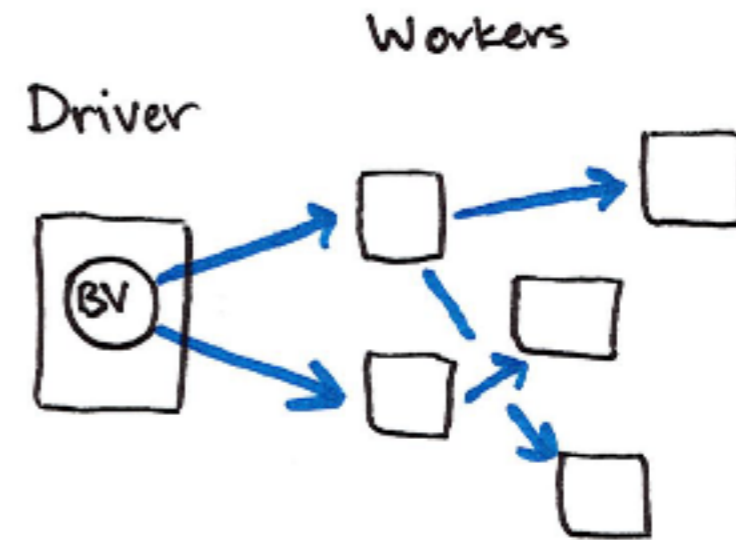
Flink & Storm

Each can have high throughput and low latency streaming

Distributed Variables

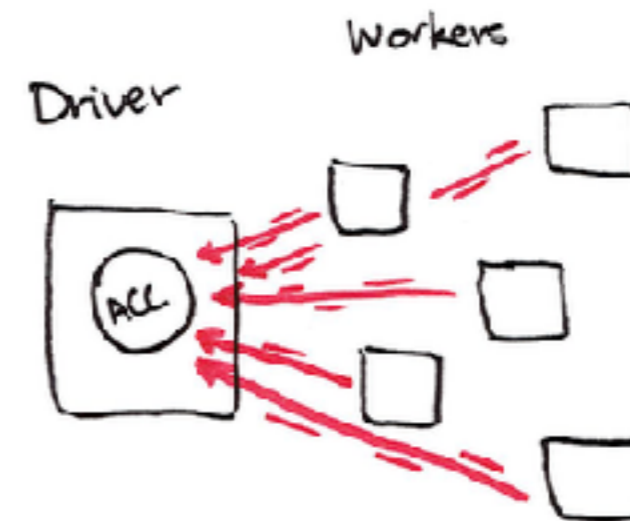
Broadcast

Read-only data shared among workers



Accumulator

Write only by workers
Read only on master



Broadcast Example

```
from pyspark.sql import SparkSession
```

```
spark = SparkSession\  
    .builder\  
    .appName("variables")\  
    .getOrCreate()
```

```
courseSize = 45
```

```
courseSizeBroadcast = spark.sparkContext.broadcast(courseSize)
```

```
courseSizeBroadcast.value
```

```
data = spark.sparkContext.parallelize((1,2,3,4,5,6,7,8), 2)
```

```
data.map(lambda x: x + courseSizeBroadcast.value).collect()
```

Using ComplexType

```
sampleMap = { 'a': 10, 'bat': 1 }
```

```
sampleBroadCast = spark.sparkContext.broadcast(sampleMap)
```

```
sampleBroadCast.value
```

```

import org.apache.spark.sql.SparkSession
val blockSize = "4096"
val spark = SparkSession.builder().
  appName("Broadcast Test").
  config("spark.broadcast.blockSize", blockSize).
  getOrCreate()

val sc = spark.sparkContext
val slices = 2
val num = 10000000

val arr1 = (0 until num).toArray

for (i <- 0 until 3) {
  println("Iteration " + i)
  println("=====")
  val startTime = System.nanoTime
  val barr1 = sc.broadcast(arr1)
  val observedSizes = sc.parallelize(1 to 10, slices).map(_ => barr1.value.length)
  observedSizes.collect().foreach(i => println(i))
  println("Iteration %d took %.0f milliseconds".format(i, (System.nanoTime - startTime) / 1E6))
}

```

Accumulator Example

```
from pyspark.sql import SparkSession
```

```
spark = SparkSession\  
    .builder\  
    .appName("variables")\  
    .getOrCreate()
```

```
counter = spark.sparkContext.accumulator(0)
```

```
def count(item):  
    global counter  
    print("item: ", item.id)  
    counter.add(1)
```

```
df = spark.range(16)  
smaller = df.coalesce(4)  
smaller.foreach(count)  
counter.value
```

Output

16

Accumulator

add()

value

Numbers only

Can create custom accumulators

Machine Learning in Spark

MLlib

RDD-based

`org.apache.spark.mllib`

Maintenance mode

DataFrame based (Spark ML)

`org.apache.spark.ml`

Pipelines

Inspired by Python scikit-learn

Classification

Regression

Clustering

Collaborative Filtering

Dimension reduction

Linear Algebra

Statistics

MLlib: Main Guide

- [Basic statistics](#)
- [Data sources](#)
- [Pipelines](#)
- [Extracting, transforming and selecting features](#)
- [Classification and Regression](#)
- [Clustering](#)
- [Collaborative filtering](#)
- [Frequent Pattern Mining](#)
- [Model selection and tuning](#)
- [Advanced topics](#)

Machine Learning Library (MLlib) Guide

MLlib is Spark's machine learning (ML) library. Its goal is to make practical machine learning scalable provides tools such as:

- **ML Algorithms:** common learning algorithms such as classification, regression, clustering, and co
- **Featurization:** feature extraction, transformation, dimensionality reduction, and selection
- **Pipelines:** tools for constructing, evaluating, and tuning ML Pipelines
- **Persistence:** saving and load algorithms, models, and Pipelines
- **Utilities:** linear algebra, statistics, data handling, etc.

Announcement: DataFrame-based API is primary

The MLlib RDD-based API is now in maintenance mode.

As of Spark 2.0, the [RDD](#)-based APIs in the `spark.mllib` package have entered maintenance mode.

Python Examples

`$SPARK_INSTALL_DIR/examples/src/main/python/ml`

aft_survival_regression.py
als_example.py
binarizer_example.py
bisecting_k_means_example.py
bucketed_random_projection_lsh_example.py
bucketizer_example.py
chi_square_test_example.py
chisq_selector_example.py
correlation_example.py
count_vectorizer_example.py
cross_validator.py
dataframe_example.py
dct_example.py
decision_tree_classification_example.py
decision_tree_regression_example.py
elementwise_product_example.py
estimator_transformer_param_example.py
feature_hasher_example.py
fpgrowth_example.py
gaussian_mixture_example.py
generalized_linear_regression_example.py
gradient_boosted_tree_classifier_example.py
gradient_boosted_tree_regressor_example.py
imputer_example.py

index_to_string_example.py
isotonic_regression_example.py
kmeans_example.py
lda_example.py
linear_regression_with_elastic_net.py
linearsvc.py
logistic_regression_summary_example.py
logistic_regression_with_elastic_net.py
max_abs_scaler_example.py
min_hash_lsh_example.py
min_max_scaler_example.py
multiclass_logistic_regression_with_elastic_net.py
multilayer_perceptron_classification.py
n_gram_example.py
naive_bayes_example.py
normalizer_example.py
one_vs_rest_example.py
onehot_encoder_estimator_example.py
pca_example.py
pipeline_example.py
polynomial_expansion_example.py
prefixspan_example.py
quantile_discretizer_example.py
random_forest_classifier_example.py