

## **Felipe Guth – COMP41110 – Student Number: 14210231**

### **Project-3 (Hadoop/Map/Reduce- log analysis) 1- Student**

In this project, you will gain hands-on experience with the Map/Reduce programming model building and running an application using Hadoop. A common task in digital forensics is the analysis of log files. Each line (entry) of a web server's log file normally contains important information such as: IP address, data and time of request, request line, HTTP status, etc. Write a Map Reduce program that takes a log file as an input and then extracts the following information from the log files:

- a) Total number of connections to the server (i.e. total numbers of entries).
- b) List of distinct IPs.
- c) Number of entries for each IP.
- d) List of distinct IPs, number of entries for each IP in a given periods (e.g. from 01/02/2015 to 28/02/2015)

You can use any Hadoop platforms to test your program (AWS Map/Reduce/Hadoop, Google Hadoop Cloud, etc.).

### **Model Design**

The next figure shows how the solution for log analysis was designed.

**Log File:** Is the input file of a server that follow a determine pattern of logs of web servers. For this application it was used a sample log downloaded from the Apache HTTP server project. It contains a list of request to a server from different IPs and shows data as IP, date time, link, type of operation, etc.

**Map Reducer:** Map/Reduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

### **Map tasks:**

**Map Input:** The input data of the map tasks are a number of n lines of the log file that contains data about the requests to a web server.

**Map Output:** The output of the map tasks are the emission of {i,k} pairs of key, value. Where the keys are the IPs gathered in each line of the log file and the value 1 that is used for counting.

### **Reducer tasks:**

**Reducer Input:** The {i,k} key, values pairs produces by the mappers are the input for the reducer tasks. Each reduce task processes a determine number of keys and produce a counting.

**Reducer Output:** The output of each reducer tasks is an output file that contain the keys (IPs) and values (count) processed by each reducer.

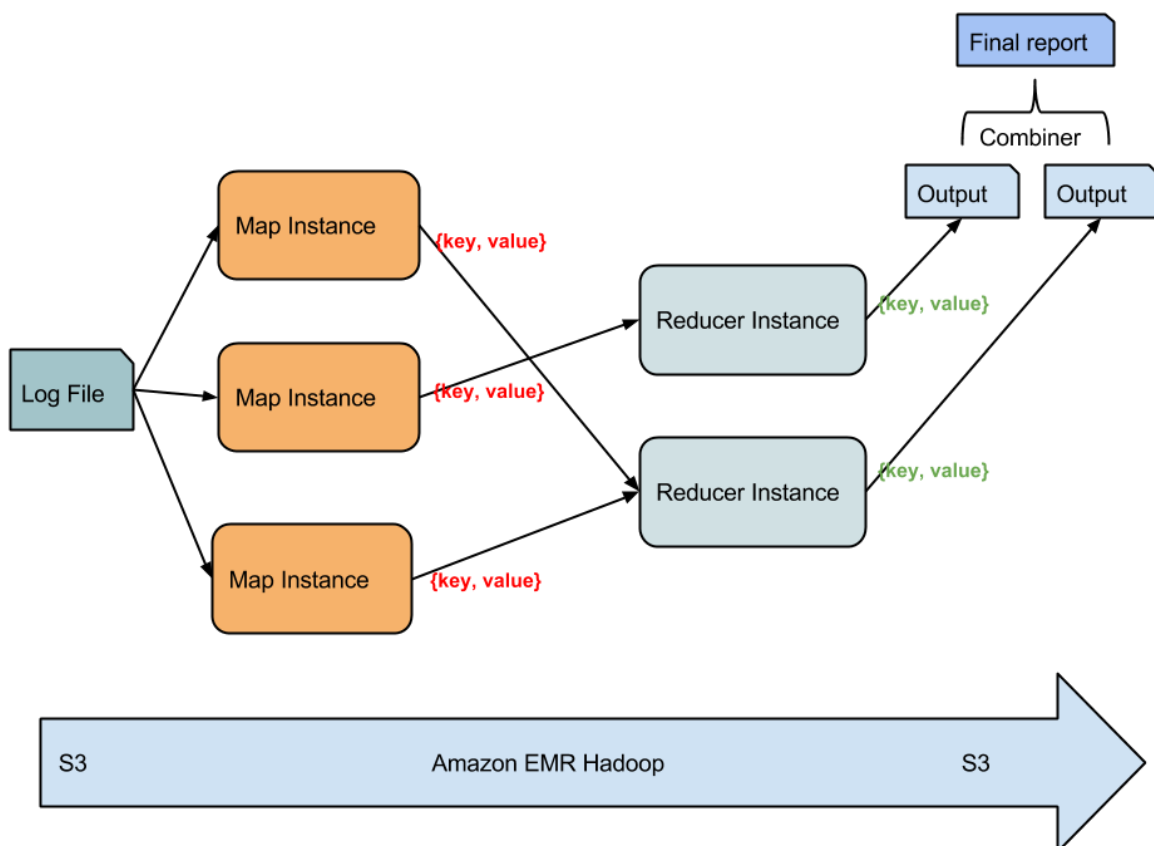
#### Combiner:

In order to generate a final report accordingly to the project propose, a combiner was implemented to read the output of each reducer and count the number of access from each IP to produce the total number of connections to the server.

#### Filtering:

With the purpose of filter the data given a data range it was implemented in the map function a routine that takes into account the dates specified by the user. As in each Map Reduce application on Hadoop, the user must specify the archives for the mapper and the reducer, it turns out that the most time effective solution is to include the range dates directly in the map script. In that way, many different Map Reduce applications can be configured with different copies of the map script filtering the data in distinct date ranges.

In local tests it is possible to develop a script to receive the dates as parameters by the system call, i.e., `./map.py 05/03.2004 06/03/2004`. However in the map reduce Hadoop environment this seems not possible until now. A question was created in the Stack Overflow regarding the passage of parameters in Hadoop, but there is no answer until this date (<http://stackoverflow.com/questions/29549740/how-to-pass-arguments-to-streaming-job-on-amazon-emr>). For this reason, the date parameters were fixed in the code of the map function.



## **System Architecture**

In order to test the Map Reduce application for log analysis the system was deployed in the Amazon Infrastructure. Amazon provides different version of Hadoop that use server and storage integrated resources, also provided by Amazon.

### **Hadoop:**

Hadoop Map/Reduce is a software framework for easily writing applications which process vast amounts of data (multi-terabyte data-sets) in-parallel on large clusters (thousands of nodes) of commodity hardware in a reliable, fault-tolerant manner.

A Map/Reduce job usually splits the input data-set into independent chunks which are processed by the map tasks in a completely parallel manner. The framework sorts the outputs of the maps, which are then input to the reduce tasks. Typically both the input and the output of the job are stored in a file-system. The framework takes care of scheduling tasks, monitoring them and re-executes the failed tasks.

Typically the compute nodes and the storage nodes are the same, that is, the Map/Reduce framework and the Hadoop Distributed File System (see HDFS Architecture Guide) are running on the same set of nodes. This configuration allows the framework to effectively schedule tasks on the nodes where data is already present, resulting in very high aggregate bandwidth across the cluster.

### **Amazon Elastic Map Reduce (EMR):**

Amazon Elastic Map/Reduce (Amazon EMR) is a web service that makes it easy to quickly and cost-effectively process vast amounts of data. Amazon EMR is based on Hadoop, a Java-based programming framework that supports the processing of large data sets in a distributed computing environment. Amazon EMR processes data across a Hadoop cluster of virtual servers on the Amazon Elastic Compute Cloud (EC2). The elastic in EMR's name refers to its dynamic resizing ability, which allows it to ramp up or reduce resource use depending on the demand at any given time.

Amazon EMR is used for data analysis in log analysis, web indexing, data warehousing, machine learning, financial analysis, scientific simulation, bioinformatics and more. (Amazon.com)

### **Amazon Simple Storage Service (S3):**

Amazon Simple Storage Service (Amazon S3), provides developers and IT teams with secure, durable, highly-scalable object storage. Amazon S3 is easy to use, with a simple web services interface to store and retrieve any amount of data from anywhere on the web. With Amazon S3, you pay only for the storage you actually use. (Amazon.com).

The S3 is integrated with the Amazon solution (EMR), where the scripts of the map and reducer and the input log are linked in the Map Reduce application. The EMR also generates the reducer output files in a pre-defined S3 bucket. After that, these files are downloaded from the Amazon S3 cloud and are processed by the combiner script to produce the final report.

## Local Testing

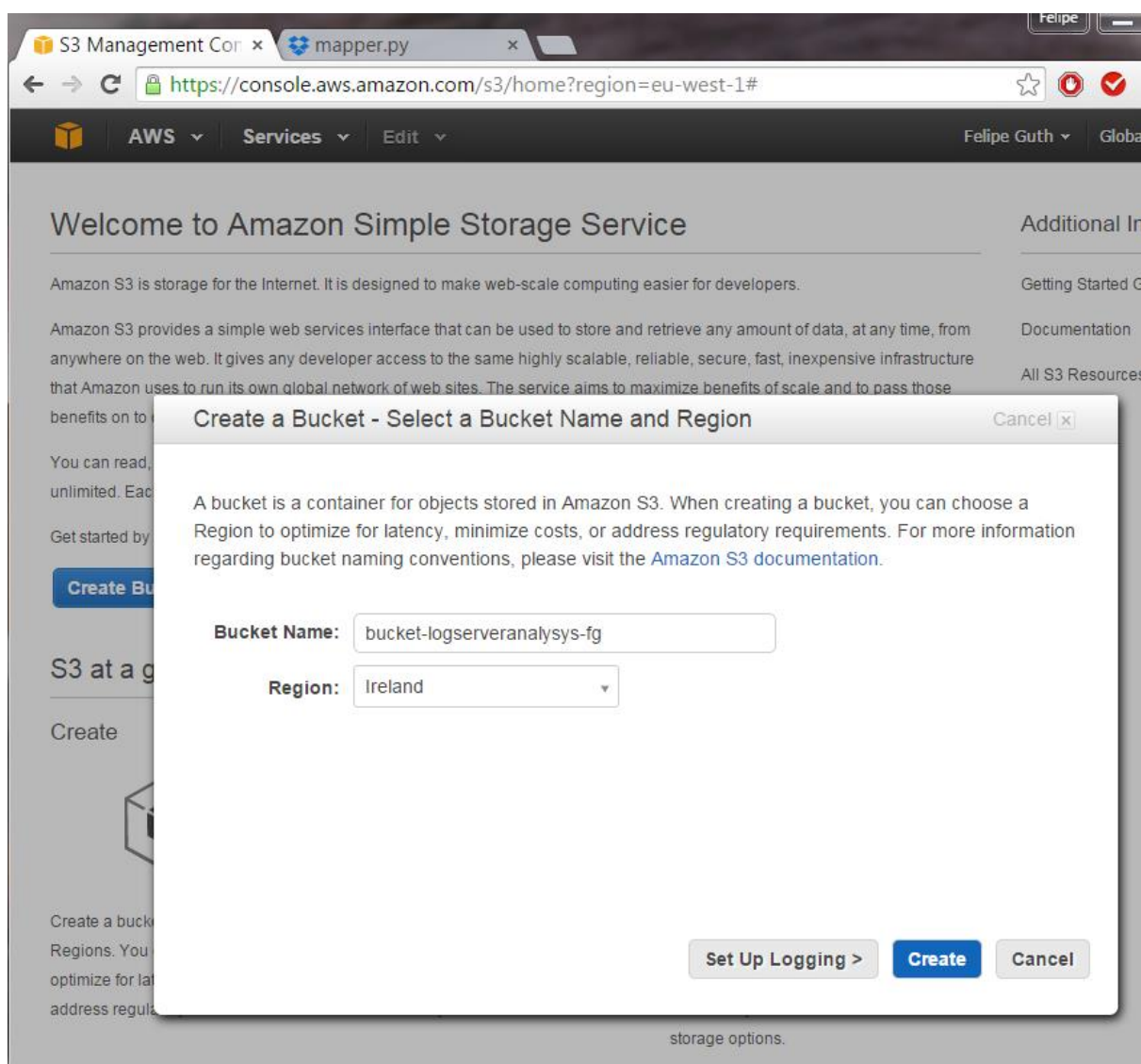
Beforehand the deployment of the system in the cloud infrastructure, local tests were performed. The scripts were development in the Python language. To test it is necessary to identify the scripts are executable archives. In Linux: `chmod +x filename`. After that the test is straightforward:

```
cat access_log | ./map.py | ./reducer.py
```

## Deployment of the System in the Amazon Cloud

### Creation of a storage bucket in the Amazon S3

The first step was to create a bucket to store the input files, scripts and output.



After the creation of the bucket, folders were created inside the bucket for the input log file and the script programs, afterwards the scripts and the input log file were upload to the Amazon cloud.

## Hadoop Cluster Creation with Amazon EMR

The following screenshots show the creation of the EMR cluster in the Amazon cloud platform.

Elastic MapReduce

Create Cluster

Cluster Configuration

Configure sample application

Cluster name

loganalysis

Termination protection

☒ Yes

☐ No

Prevents accidental termination of the cluster: to shut down the cluster, you must turn off termination protection. [Learn more](#)

Logging

☒ Enabled

Copy the cluster's log files automatically to S3. [Learn more](#)

Log folder S3 location

s3://aws-logs-050067417136-eu-west-1/elasticmapreduce/

s3://<bucket-name>/<folder>/

Debugging

☐ Enabled

Index logs to enable console debugging functionality (requires logging). [Learn more](#)

Tags

Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are propagated to the underlying EC2 instances. [Learn more](#) about tagging your Amazon EMR clusters.

Key	Value (optional)
<div>Add a key to create a tag</div>	

Software Configuration

Hadoop distribution

☒ Amazon

Use Amazon's Hadoop distribution. [Learn more](#)

AMI version

3.6.0

Determines the base configuration of the instances in your cluster, including the Hadoop version. [Learn more](#)

☐ MapR

Use MapR's Hadoop distribution. [Learn more](#)

Applications to be installed	Version			
Hive	0.13.1			
Pig	0.12.0			
Hue	3.7.1			

Additional applications

Select an application

Configure and add

## File System Configuration

**i** The [EMR File System \(EMRFS\)](#) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores data on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable [S3 server-side encryption](#) or [S3 client-side encryption](#) and [consistent view](#) for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS.

EMRFS S3 Encryption

Choose encryption method for objects written to or read from S3 using EMRFS. Please note that this will not encrypt files written to HDFS. [Learn more](#)

Consistent view ☐ Enabled

Monitors list and read-after-write (for new puts) consistency for files in S3. [Learn more](#)

## Hardware Configuration

**i** Specify the [networking](#) and [hardware](#) configuration for your cluster. If you need more than 20 EC2 instances, [complete this form](#). [Request Spot instances](#) (unused EC2 capacity) to save money.

Network

Use a Virtual Private Cloud (VPC) to process sensitive data or connect to a private network. [Create a VPC](#)

EC2 Subnet

[Create a Subnet](#)

Type	Name	EC2 instance type	Count	Request spot	Bid price			
Master	<input type="text" value="Master instance group - 1"/>	<input type="text" value="m1.medium"/>	<input type="text" value="1"/>	<input type="checkbox"/>				<b>?</b>
Core	<input type="text" value="Core instance group - 2"/>	<input type="text" value="c1.medium"/>	<input type="text" value="2"/>	<input type="checkbox"/>				<b>?</b>
Task	<input type="text" value="Task instance group - 3"/>	<input type="text" value="c1.medium"/>	<input type="text" value="0"/>	<input type="checkbox"/>				<b>x ?</b>

[Add task instance group](#)

## Security and Access

EC2 key pair

Use an existing EC2 key pair to SSH into the master node of the Amazon EMR cluster. [Learn more](#)

IAM user access ☒ All other IAM users  
☐ No other IAM users

Control the visibility of this cluster to other IAM users. [Learn more](#)

### IAM Roles

**i** An IAM role for the EMR service and an EC2 instance profile for instances in an EMR cluster are recommended. You can create and assign these roles to limit the permissions of the EMR service and applications running on a cluster. [Learn more](#)

Roles configuration ☐ Proceed without roles

No roles will be used. Selecting "Default" is recommended. [Learn more](#)

☒ Default  
[View policies for default roles](#)

Use default roles for your cluster. If not present, they will be automatically created for you. [Learn more](#)

☐ Custom

Select custom roles to tailor permissions for your cluster. [Learn more](#)

The Hardware configuration should be selected according to the needs of the user. A range of different hardware configurations are available. Best resources result in a more expensive bill.

### ▼ EC2 Security Groups

**i** An EC2 security group acts as a virtual firewall for your cluster nodes to control inbound and outbound traffic. There are two types of security groups you can configure, [EMR managed security groups](#) and [additional security groups](#). EMR will [automatically update](#) the rules in the EMR managed security groups in order to launch a cluster. [Learn more](#).

Type	EMR managed security groups <small>EMR will automatically update the selected group</small>	Additional security groups <small>EMR will not modify the selected groups</small>
Master	Default: sg-a1d5abc4 (ElasticMapReduce-master) ▼	No security groups selected
Core & Task	Default: sg-a0d5abc5 (ElasticMapReduce-slave) ▼	No security groups selected

[Create a security group](#)

### Bootstrap Actions

**i** Bootstrap actions are scripts that are executed during setup before Hadoop starts on every cluster node. You can use them to install additional software and customize your applications. [Learn more](#)

Bootstrap action type	Name	S3 location	Optional arguments		
Add bootstrap action		Select a bootstrap action ▼			
Configure and add					

### Steps

**i** A step is a unit of work you submit to the cluster. A step might contain one or more Hadoop jobs, or contain instructions to install or configure an application. You can submit up to 256 steps to a cluster. [Learn more](#)

Name	Action on failure	JAR location	Arguments		
Add step		<div>Streaming program ▼</div> <div>Configure and add</div>			
Auto-terminate		Automatically terminate cluster after the last step is completed.			
<input checked="" type="radio"/> Yes					
<input type="radio"/> No		Keep cluster running until you terminate it.			

**i** No EC2 key pair has been selected, so you will not be able to SSH to this cluster or connect to HUE (unless you are using a VPN). [Learn how to create an EC2 Key Pair](#).

Cancel

Create cluster

The configuration of the Map/Reduce tasks occurs in the Steps section of the cluster setup. As Python scripts were used, in the Add step option, the Streaming program was selected, after this one should click in the “Configure and add” button.

**Linking EMR with S3: Select the scripts of the map and reducer and defining input and output directories**

The “configure and add” button results in the following screen:

Add Step

Step type

Streaming program

Name\*

LogAnalysis

Mapper\*

s3://bucket-logserveranalysis-fg/programs/mapperLog.py

S3 location of the map function or the name of the Hadoop streaming command to run.

Reducer\*

s3://bucket-logserveranalysis-fg/programs/reducer.py

S3 location of the reduce function or the name of the Hadoop streaming command to run.

Input S3 location\*

s3://bucket-logserveranalysis-fg/input/

s3://<bucket-name>/<folder>/

Output S3 location\*

s3://bucket-logserveranalysis-fg/output2/

s3://<bucket-name>/<folder>/

Arguments

Action on failure

Continue

What to do if the step fails.

Cancel

Save

In this step, the user should select the map and reduce programs of the streaming job to be executed in the Hadoop cluster. The input and output directory must be selected as well. The output directory must be an inexistent directory that will be automatically created.

After click in save, one can see the configuration of the Map Reduce steps as shown in the next screenshot. The “create cluster” button initialize the resources and perform the configured steps.

## Steps

**i** A step is a unit of work you submit to the cluster. A step might contain one or more Hadoop jobs, or contain instructions to install or configure an application. You can submit up to 256 steps to a cluster. [Learn more](#)

Name	Action on failure	JAR location	Arguments
LogAnalysis	Continue	/home/hadoop/contrib/streaming/hadoop-streaming.jar	-files s3://bucket-logserveranalysis-fg/programs/mapperLog.py,s3://bucket-logserveranalysis-fg/programs/reducer.py -mapper mapperLog.py -reducer reducer.py -input s3://bucket-logserveranalysis-fg/input/ -output s3://bucket-logserveranalysis-fg/output2/

Add step 

Select a step

Configure and add

Auto-terminate ☒ Yes

☐ No

Automatically terminate cluster after the last step is completed.

Keep cluster running until you terminate it.

**i** 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

Create cluster



## Cluster Running

The user can monitor the cluster by the Amazon Cluster Details option of the EMR.

The screenshot displays the AWS Elastic MapReduce console interface. At the top, there's a navigation bar with 'Elastic MapReduce', 'Cluster List', and 'Cluster Details'. Below this, a cluster named 'logserverfinal33' is shown in a 'Running' state. The console provides a comprehensive overview of the cluster's configuration, including connections, summary, configuration details (AMI version, Hadoop distribution, applications, log URI), network and hardware (availability zone, subnet, master and core instances), and security and access (key name, EC2 instance profile, EMR role, security groups). A 'Steps' section is expanded, showing a table of four completed steps: 'LogAnalysis', 'Setup pig', 'Setup hive', and 'Setup hadoop debugging'. Each step entry includes its ID, name, status, start time, elapsed time, and a link to view logs. The bottom of the console shows 'Bootstrap Actions' and a footer with copyright information and a feedback button.

## Task Detail

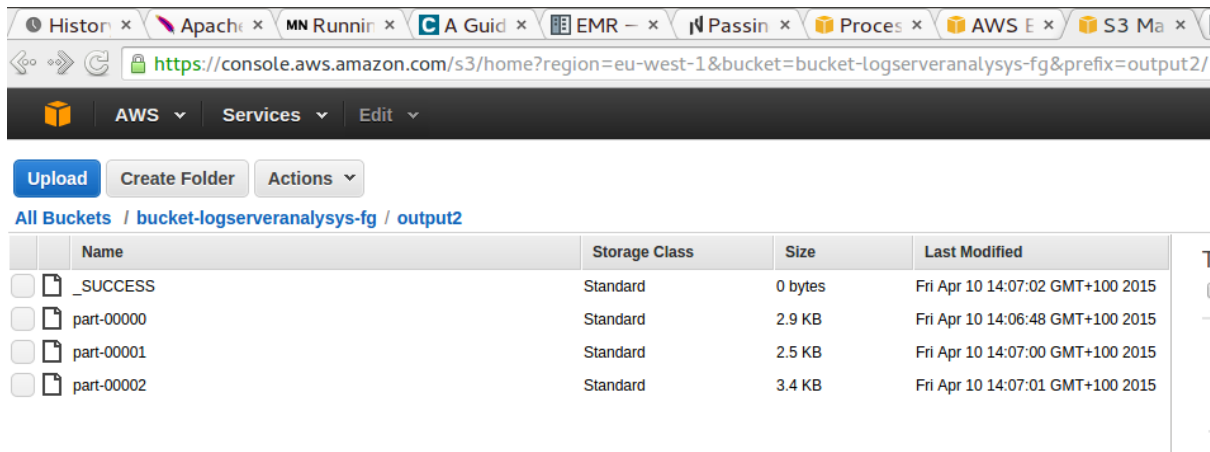
Also, it is possible to see the logs and details of the steps of the Map Reduce tasks as shown in the screenshot.

This screenshot shows the 'Task Detail' view within the AWS Elastic MapReduce console. It focuses on a specific step, 's-LZ4JRS55B05B', and displays a table of tasks associated with it. The tasks are categorized by type (REDUCE or MAP) and their state (COMPLETED). The table includes columns for the task ID, type, state, start time, and a link to view attempts. The tasks listed are all in a 'COMPLETED' state, with start times ranging from 2015-04-10 14:06:11 to 2015-04-10 14:06:39 UTC+1.

Task	Type	State	Start time (UTC+1)	Actions
r_000002	REDUCE	COMPLETED	2015-04-10 14:06:39 (UTC+1)	<a href="#">View attempts</a>
r_000001	REDUCE	COMPLETED	2015-04-10 14:06:38 (UTC+1)	<a href="#">View attempts</a>
r_000000	REDUCE	COMPLETED	2015-04-10 14:06:36 (UTC+1)	<a href="#">View attempts</a>
m_000003	MAP	COMPLETED	2015-04-10 14:06:24 (UTC+1)	<a href="#">View attempts</a>
m_000002	MAP	COMPLETED	2015-04-10 14:06:12 (UTC+1)	<a href="#">View attempts</a>
m_000001	MAP	COMPLETED	2015-04-10 14:06:12 (UTC+1)	<a href="#">View attempts</a>
m_000000	MAP	COMPLETED	2015-04-10 14:06:11 (UTC+1)	<a href="#">View attempts</a>

## S3 Result

After the completion of the tasks, the user can check the output files in the configured directory, as shown in the screenshot. In this case, these outputs are the result of the log analysis performed by the map and reducer tasks. The results are distributed in files according to the number of reducers that were allocated by the Hadoop cluster.



## Combining the results and showing the report

After downloading the output files from the Amazon S3, the files are processed by a combiner in order to show the result in an organized way, as demonstrated in the next screenshot. In the combiner script it is necessary to configure the directory where the files generated by the reducer are located using the variable "dirOutput".

In this screenshot the logs were filter by date from: 05/03/2004 date to: 05/03/2004. In this range there are just 5 registers. The directory where the output files are stored can be configured in the combiner script. This was used to present the result in a feasible way in this report. The combiner can be run as:

```
./combiner.py
```

```
fclub@fgub-DELL:~/Documents/AptanaStudio3/IpAnalisys/sources$ ./combiner.py
-----NUMBER OF CONNECTIONS TO THE SERVER:      5
-----NUMBER OF ENTRIES FOR EACH IP:
64.242.88.10      3
prxint-sxb3.e-i.net      2
-----LIST OF DISTINCT IPs:
64.242.88.10
prxint-sxb3.e-i.net
fclub@fgub-DELL:~/Documents/AptanaStudio3/IpAnalisys/sources$
```