



UNIVERSIDAD TECNICA
FEDERICO SANTA MARIA

Trabajo Práctico 1

INF356 - 2025-1 - 200

Computación Distribuida para Big Data

4 de mayo de 2025 - v1.0

Índice

1. Despliegue del cluster	4
1.1. Implementación del tutorial	4
1.2. Expansión del cluster	15
1.3. Cliente web	16
2. Instalación de Apache Hive	18
2.1. Procedimiento	18
2.2. Prueba	21
3. Exploración del HDFS	29
4. Uso del cluster	31
4.1. Importación	31
4.2. Parsing	33
4.3. Análisis	35

Índice de figuras

1. Captura de pantalla de la “AWS Management Console - EC2” que muestra las máquinas del cluster creado con el tutorial	14
2. Captura de pantalla de la “AWS Management Console - EC2” que muestra las máquinas del cluster expandido	16
3. Captura de pantalla del cliente web de Hadoop	17
4. Captura de pantalla del cliente web de HDFS	17
5. Instancias mejoradas para no sufrir	18

Índice de fragmentos de código

1. Conectarse por primera vez al Master usando la llave PEM con la IP publica	4
2. Actualizar repositorios locales y descargar actualizaciones	5
3. Descargar Java / OpenJDK 11 usando APT	5
4. Crear una llave publica en formato RSA y autorizarla	5
5. Descargar y descomprimir Hadoop en el Master	5
6. Editar las Variables de Entorno usando el usuario Root y el editor de texto Vim	5
7. Variables de Entorno que usaremos para el cluster	6

8.	hadoop-3.3.6/etc/hadoop/core-site.xml	7
9.	hadoop-3.3.6/etc/hadoop/hdfs-site.xml	7
10.	hadoop-3.3.6/etc/hadoop/mapred-site.xml	8
11.	hadoop-3.3.6/etc/hadoop/yarn-site.xml	9
12.	IP's privadas de los Workers en la carpeta Home del Master	10
13.	Borrar el archivo hadoop-3.3.6.tar.gz y comprimimos el directorio hadoop-3.3.6	10
14.	Script para configurar todos los Workers desde el Master	11
15.	Script para configurar todos los Workers desde el Master	12
16.	Formateamos el File System, lo iniciamos y llamamos al negociador de recursos YARN	13
17.	Corroborar estado de nuestro HDFS	13
18.	Estado del HDFS luego de haber ejecutado un Map Reduce	13
19.	Ejecutar el Map Reduce y leer su salida	14
20.	Los archivos que hay que modificar antes de ejecutar el script de instalacion nuevamente	15
21.	Usar WGET para descargar Apache Hive 4.0.1	18
22.	Descomprimos la carpeta de Apache Hive 4.0.1	18
23.	Archivo de configuracion de Apache Hive	19
24.	Levantar MetaStore y HiveServer2	19
25.	Usar Beeline como interfaz a Hive	20
26.	Ejemplo de uso de Apache Hive	21
27.	Salida de Apache Hive (Parte 1)	22
28.	Salida de Apache Hive (Parte 2)	23
29.	Salida de Apache Hive (Parte 3)	24
30.	Salida de Apache Hive (Parte 4)	25
31.	Salida de Apache Hive (Parte 5)	26
32.	Salida de Apache Hive (Parte 6)	27
33.	Salida de Apache Hive (Parte 7)	28
34.	Resultado esperado para ejemplo de uso de Apache Hive	28
35.	Ejemplo de uso de Apache Hive	29
36.	Muestra el archivo, los bloques que ocupa, el factor de replicacion y la ubicacion de cada bloque	29
37.	Instalar AWS CLI	30
38.	Importacion de los datos sucios	31
39.	Importacion de los datos sucios	32
40.	Limpieza de los datos sucios	33
41.	Importacion de los datos sucios	33
42.	Analisis de los datos limpios para entender que significa cada columna	34
43.	Observacion con nombres y tipos correctos	35
44.	Observacion con nombres y tipos correctos	36

1. Despliegue del cluster

1.1. Implementación del tutorial

Para la creación del Cluster en Amazon Web Services, fue necesario seguir los pasos indicados en el tutorial entregado en el Slide T01 - Hadoop Cluster - v1.0 entregado en Aula. El instructivo constaba de los siguientes pasos:

- Crear la cuenta de AWS
- Preparar el entorno del Cluster
- Configurar la máquina Master
- Configurar todos los Workers
- Inicializar el Cluster
- Ejecutar el script de Map Reduce con Hadoop

Crear la cuenta de AWS

Para la creación de la cuenta de AWS simplemente seguí las indicaciones en el correo de invitación indicado por el profesor. Posteriormente a eso, me fui directamente a la sección de Modules, en donde estaba el Módulo de Launch AWS Academy Learner Lab, que sirve para entrar a la Console Home de AWS.

Preparar el entorno del Cluster

Una vez en Console Home, creamos una instancia de tipo EC2, para la cual creamos los Key Pairs, el Security Group y la instancia Master.

Para los Key Pairs, creamos una llave de tipo RSA con formato de archivo PEM.

Para el Security Group, creamos tres reglas, una Inbound Rule de tipo SSH con tipo de fuente Anywhere-IPv4 para entrar por SSH, una Inbound Rule de tipo Custom TCP con rango de puerto 8088 y tipo de fuente Anywhere-IPv4 para la consola web de Hadoop, una Inbound Rule de tipo Custom TCP con rango de puerto 9870 y tipo de fuente Anywhere-IPv4 para la consola web de HDFS y finalmente, una Inbound Rule de tipo Custom TCP con rango de puerto de 0 a 65535 y tipo de fuente Anywhere-IPv4.

Para la instancia Master creamos un nodo con Ubuntu 24.04 de tipo t2.micro y con 16gb de almacenamiento. Esta instancia usará la Key y Security Group recién creados (al igual que las instancias Worker más adelante). Una vez configurados los parámetros, lanzamos la instancia y nos conectamos por SSH:

```
ssh ubuntu@34.207.118.42 -i .\Documents\llave.pem
```

Código 1: Conectarse por primera vez al Master usando la llave PEM con la IP pública

Configurar la maquina Master

El primer paso es logicamente actualizar los paquetes del sistema, para esto es necesario ejecutar el siguiente comando:

```
1 sudo apt update && sudo apt upgrade -y
```

Código 2: Actualizar repositorios locales y descargar actualizaciones

Una vez hecho esto, instalamos el paquete correspondiente a Java / OpenJDK 11:

```
1 sudo apt-get install openjdk-11-jdk -y
```

Código 3: Descargar Java / OpenJDK 11 usando APT

Una vez que hayamos resuelto todos los paquetes necesarios, tenemos que crear una llave publica con el siguiente comando:

```
1 ssh-keygen -t rsa
2 cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
```

Código 4: Crear una llave publica en formato RSA y autorizarla

Esta llave publica la guardaremos en la carpeta de Authorized Keys para SSH, lo cual permitira al Master conectarse a si mismo por LocalHost.

Descargamos Hadoop 3.3.6 desde la pagina oficial usando el comando WGET, el cual hace una request a la pagina para descargar el archivo comprimido en formato TAR GZ, el cual tendremos que descomprimir y modificar para indicarle parametros especiales. Ejecutamos entonces los siguientes comandos:

```
1 wget https://www.apache.org/dyn/closer.cgi/hadoop/common/hadoop-3.3.6/hadoop-3.3.6-src.tar.gz
2 tar -xvz hadoop-3.3.6.tar.gz
```

Código 5: Descargar y descomprimir Hadoop en el Master

Esto creara una carpeta llamada hadoop-3.3.6 en el directorio Home del Master, esta carpeta albergara los binarios yu scripts ejecutables para desplegar nuestro Distributed File System (DFS). Para poder ejecutar los comandos en dicha carpeta de manera directa, modificaremos las Variables de Entorno de la maquina Master, de esta manera podremos ejecutar los comandos en cualquier directorio. Dichas Variables tambien contienen rutas para la ejecucion de OpenJDK y la ubicacion de nuestra carpeta Hadoop, las cuales debemos obedecer para que funcione correctamente nuestro cluster. Las variables de entorno en Linux se guardan en el directorio /etc/environment y solo puede modificarse por el Super Usuario (Root). Personalmente, a mi me gusta usar Vim, pero para aquellas personas con gustos inferiores existen alternativas:

```
1 sudo vim /etc/environment
```

Código 6: Editar las Variables de Entorno usando el usuario Root y el editor de texto Vim

Las Variables debiesen quedar asi:

```
1 PATH="/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:..."
2 JAVA_HOME="/usr/lib/jvm/java-11-openjdk-amd64"
3 HADOOP_HOME="/home/ubuntu/hadoop-3.3.6"
4 HADOOP_COMMON_HOME="/home/ubuntu/hadoop-3.3.6"
```

Código 7: Variables de Entorno que usaremos para el cluster

Lo ultimo es modificar los archivos de configuracion de Hadoop en base a lo indicado por el profesor. Estos archivos luego seran enviados a los Workers cuando comprimamos la carpeta hadoop-3.3.6.

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4   <property>
5     <name>fs.defaultFS</name>
6     <value>hdfs://172.31.30.146</value>
7   </property>
8 </configuration>
```

Código 8: hadoop-3.3.6/etc/hadoop/core-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4   <property>
5     <name>dfs.block.size</name>
6     <value>33554432</value>
7   </property>
8   <property>
9     <name>dfs.replication</name>
10    <value>3</value>
11  </property>
12  <property>
13    <name>dfs.name.dir</name>
14    <value>file:///home/ubuntu/dfs/name</value>
15  </property>
16  <property>
17    <name>dfs.data.dir</name>
18    <value>file:///home/ubuntu/dfs/data</value>
19  </property>
20  <property>
21    <name>dfs.namenode.heartbeat.recheck-interval</name>
22    <value>10000</value>
23  </property>
24 </configuration>
```

Código 9: hadoop-3.3.6/etc/hadoop/hdfs-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4   <property>
5     <name>mapreduce.framework.name</name>
6     <value>yarn</value>
7   </property>
8   <property>
9     <name>yarn.app.mapreduce.am.env</name>
10    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
11  </property>
12  <property>
13    <name>mapreduce.map.env</name>
14    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
15  </property>
16  <property>
17    <name>mapreduce.map.java.opts</name>
18    <value>-Xmx1024m</value>
19  </property>
20  <property>
21    <name>mapreduce.map.memory.mb</name>
22    <value>1024</value>
23  </property>
24  <property>
25    <name>mapreduce.reduce.env</name>
26    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
27  </property>
28  <property>
29    <name>mapreduce.reduce.java.opts</name>
30    <value>-Xmx1024m</value>
31  </property>
32  <property>
33    <name>mapreduce.reduce.memory.mb</name>
34    <value>1024</value>
35  </property>
36  <property>
37    <name>mapreduce.job.maps</name>
38    <value>4</value>
39  </property>
40  <property>
41    <name>mapreduce.job.running.map.limit</name>
42    <value>4</value>
43  </property>
44  <property>
45    <name>mapreduce.job.reduces</name>
46    <value>4</value>
47  </property>
48  <property>
49    <name>mapreduce.job.running.reduce.limit</name>
50    <value>4</value>
51  </property>
52 </configuration>
```

Código 10: hadoop-3.3.6/etc/hadoop/mapred-site.xml


```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4   <property>
5     <name>yarn.resourcemanager.hostname</name>
6     <value>172.31.30.146</value>
7   </property>
8   <property>
9     <name>yarn.nodemanager.aux-services</name>
10    <value>mapreduce_shuffle</value>
11  </property>
12  <property>
13    <name>yarn.nodemanager.resource.memory-mb</name>
14    <value>1024</value>
15  </property>
16  <property>
17    <name>yarn.nodemanager.resource.cpu-vcores</name>
18    <value>1</value>
19  </property>
20  <property>
21    <name>yarn.app.mapreduce.am.resource.mb</name>
22    <value>1024</value>
23  </property>
24  <property>
25    <name>yarn.app.mapreduce.am.command-opts</name>
26    <value>-Xmx1024m</value>
27  </property>
28  <property>
29    <name>yarn.app.mapreduce.am.resource.cpu-vcores</name>
30    <value>1</value>
31  </property>
32 </configuration>
```

Código 11: hadoop-3.3.6/etc/hadoop/yarn-site.xml

Configurar todos los Workers

Los Workers también serán máquinas con Ubuntu 24.04 y de tipo t2.micro, usarán las mismas Key Pairs y el mismo Security Group. Lo único que cambiará entre ellos y el Master será el almacenamiento, las cuales constarán solo de 8gb en total. Una vez creados los Workers desde el interfaz de EC2 de AWS, tendremos que obtener sus IP's públicas y privadas. La primera para conectarnos por primera vez, la segunda para establecer una conexión permanente. Anotamos las IP's privadas de la siguiente manera:

```
1 172.31.30.0
2 172.31.17.157
3 172.31.24.51
4 172.31.27.79
```

Código 12: IP's privadas de los Workers en la carpeta Home del Master

Luego de tener esto, borramos el archivo comprimido de Hadoop que descargamos inicialmente y comprimimos el directorio que modificamos en el paso anterior para luego enviárselo a los Workers. Estos luego descomprimirán el archivo de manera local para así no abusar de nuestro ancho de banda limitado:

```
1 tar czf hadoop-3.3.6.tar.gz hadoop-3.3.6
```

Código 13: Borrar el archivo hadoop-3.3.6.tar.gz y comprimimos el directorio hadoop-3.3.6

Y como la tarea de configurar cada Worker manualmente es tedioso, usamos un script de Bash que realiza el trabajo por nosotros. El cual itera los siguientes pasos:

Para cada una de las instancias Worker proceder de la siguiente forma:

1. Instalar la llave pública de la máquina Master en el Worker
2. Crear una llave para el Worker
3. Instalar la llave del Worker en el Master
4. Instalar la llave del Worker en el Worker
5. Actualizar el sistema operativo en el Worker
6. Instalar Java en el Worker (openjdk-11-jdk)
7. Editar el archivo /etc/environment de la misma forma que para Master
8. Reiniciar el worker
9. Transferir el archivo de despliegue de Hadoop de la máquina Master al Worker
10. Expandir el archivo de despliegue en la máquina Worker

El script utilizado quedó de la siguiente fórmula:

```

1 #!/bin/bash
2
3 # === CONFIGURACIÓN ===
4
5 # IPs publicas de los Workers y el Master
6 WORKERS=("172.31.30.0" "172.31.17.157" "172.31.24.51" "172.31.27.79")
7 MASTER="172.31.30.146"
8
9 # User de el Master y los Workers
10 USER="ubuntu"
11 HOME="/home/$USER"
12
13 # Ubicacion de la llave privada
14 KEY="$HOME/.ssh/llave.pem"
15
16 # Variables de Entorno
17 PATH="/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games..."
18 JAVA_HOME="/usr/lib/jvm/java-11-openjdk-amd64"
19 HADOOP_HOME="/home/ubuntu/hadoop-3.3.6"
20 HADOOP_COMMON_HOME="/home/ubuntu/hadoop-3.3.6"
21
22 # Path del archivo comprimido de Hadoop
23 HADOOP_ARCHIVE="/home/ubuntu/hadoop-3.3.6.tar.gz"
24
25 # Ubicacion de las Variables de Entorno
26 ENV_VARS="/etc/environment"
27
28 for WORKER in "${WORKERS[@]}; do
29     echo ">>>Configurando al Worker de IP $WORKER<<<"
30
31     echo ""
32     echo "1) Instalando llave publica del Master en el Worker"
33     MASTER_KEY="$(cat ~/.ssh/id_rsa.pub)"
34     echo ">>ssh $USER@$WORKER -i $KEY echo $MASTER_KEY>> ~/.ssh/authorized_keys"
35     ssh $USER@$WORKER -i $KEY "echo $MASTER_KEY>> ~/.ssh/authorized_keys"
36
37     echo ""
38     echo "2) Creando llave publica para el Worker"
39     echo ">>ssh $USER@$WORKER -i $KEY ssh-keygen -t rsa"
40     ssh $USER@$WORKER -i $KEY "ssh-keygen -t rsa"
41     echo ""
42
43     echo ""
44     echo "3) Instalando llave del Worker en el Master"
45     echo ">>ssh $USER@$WORKER -i $KEY cat ~/.ssh/id_rsa.pub>> ~/.ssh/authorized_keys"
46     ssh $USER@$WORKER -i $KEY "cat ~/.ssh/id_rsa.pub" >> ~/.ssh/authorized_keys
47     echo ""
48
49     echo ""
50     echo "4) Instalando llave del Worker en el mismo Worker"
51     echo ">>ssh $USER@$WORKER cat ~/.ssh/id_rsa.pub>> ~/.ssh/authorized_keys"
52     ssh $USER@$WORKER "cat ~/.ssh/id_rsa.pub>> ~/.ssh/authorized_keys"
53     echo ""
54
55     echo ""
56     echo "5) Actualizando sistema operativo en el Worker"
57     echo ">>ssh $USER@$WORKER sudo apt update -qq && sudo apt upgrade -y -qq"
58     ssh $USER@$WORKER "sudo apt update -qq && sudo apt upgrade -y -qq"
59     echo ""
60     ...

```

Código 14: Script para configurar todos los Workers desde el Master

```

1  ...
2  echo ""
3  echo "6) Instalando Java en el Worker"
4  echo "└─> ssh $USER@$WORKER sudo apt-get install -y -qq openjdk-11-jdk"
5  ssh $USER@$WORKER "sudo apt-get install -y -qq openjdk-11-jdk"
6  echo ""
7
8  echo ""
9  # Esta seccion me dio problemas asi que este paso lo hice manual al final
10 echo "7) Copiando las Variables de Entorno en el Worker"
11     echo "└─> cat $HOME/environment"
12     cat $HOME/environment
13 echo "└─> scp $HOME/environment $USER@$WORKER:$HOME/"
14 scp $HOME/environment $USER@$WORKER:$HOME/
15 echo "ssh $USER@$WORKER cat $HOME/ambiente"
16 ssh $USER@$WORKER "cat $HOME/ambiente"
17 echo "└─> ssh $USER@$WORKER sudo cp $HOME/ambiente $ENV_VARS"
18 ssh $USER@$WORKER "sudo cp $HOME/ambiente $ENV_VARS"
19 echo "└─> ssh $USER@$WORKER cat $ENV_VARS"
20 ssh $USER@$WORKER "cat $ENV_VARS"
21
22 echo ""
23
24 echo ""
25 echo "8) Copiando archivo de despliegue Hadoop al Worker..."
26 echo "ssh $USER@$WORKER rm $HADOOP_ARCHIVE && ls -la"
27 ssh $USER@$WORKER "rm $HADOOP_ARCHIVE && ls -la"
28 echo "└─> scp $HADOOP_ARCHIVE $USER@$WORKER:$HOME/"
29 scp $HADOOP_ARCHIVE $USER@$WORKER:$HOME/
30 echo "└─> ssh $USER@$WORKER ls -la"
31 ssh $USER@$WORKER "ls -la"
32 echo ""
33
34 echo ""
35 echo "9) Descomprimiendo el archivo de despliegue en el Worker..."
36 echo "└─> ssh $USER@$WORKER tar -xzf $HADOOP_HOME.tar.gz -C $HOME/"
37 ssh $USER@$WORKER "tar -xzf $HADOOP_HOME.tar.gz -C $HOME/"
38 echo ""
39
40 echo ""
41 echo "10) Reiniciando Worker..."
42 echo "└─> ssh $USER@$WORKER sudo reboot now"
43 ssh $USER@$WORKER "sudo reboot now"
44 echo ""
45
46 echo "El Worker de IP $WORKER ha sido configurado correctamente! :)"
47 echo ""
48 done
49
50 echo "Todos los Workers han sido configurados! :D"

```

Código 15: Script para configurar todos los Workers desde el Master

Una vez configurado todo, reiniciamos el nodo Master y reingresamos cuando la maquina vuelva a estar activa.

Inicializar el Cluster

Si las maquinas estan correctamente configuradas, podremos levantar el File System distribuido sin problemas. Para esto es necesario formatear el disco e iniciar los servicios para la distribucion de este y para la gestion de recursos entre los Workers. Cabe notar que por disco me refiero a unidad de almacenamiento, como vimos en clases, un DFS no almacena datos de forma continua, si no que la reparte entre varios Nodos.

Dicho esto, ejecutamos la siguiente serie de comandos:

```
1 hdfs namenode -format
2 start-dfs.sh
3 start-yarn.sh
```

Código 16: Formateamos el File System, lo iniciamos y llamamos al negociador de recursos YARN

Luego, para corroborar de que todo esta en orden, ejecutamos el siguiente comando:

```
1 hdfs fsck / -files -blocks
```

Código 17: Corroborar estado de nuestro HDFS

La salida debiese ser algo del siguiente estilo:

```
1 ...
2 Status: HEALTHY
3 Number of data-nodes: 4
4 Number of racks: 1
5 Total dirs: 11
6 Total symlinks: 0
7
8 Replicated Blocks:
9 Total size: 428156 B
10 Total files: 9
11 Total blocks (validated): 8 (avg. block size 53519 B)
12 Minimally replicated blocks: 8 (100.0 %)
13 Over-replicated blocks: 0 (0.0 %)
14 Under-replicated blocks: 0 (0.0 %)
15 Mis-replicated blocks: 0 (0.0 %)
16 Default replication factor: 3
17 Average block replication: 3.0
18 Missing blocks: 0
19 Corrupt blocks: 0
20 Missing replicas: 0 (0.0 %)
21 Blocks queued for replication: 0
22 ...
23 The filesystem under path '/' is HEALTHY
```

Código 18: Estado del HDFS luego de haber ejecutado un Map Reduce

Notaremos que existe una seccion que muestra la cantidad de Data Nodes, si muestra 4 entonces configuramos correctamente nuestro cluster.

Instances (1/5) [Info](#)

Last updated less than a minute ago

[Connect](#) [Instance state](#) [Actions](#) [Launch instances](#)

Find Instance by attribute or tag (case-sensitive) [All states](#)

	Name	Instance ID	Instance state	Instance type	Status check	Alarm
<input checked="" type="checkbox"/>	worker-1	i-02726e3b5b124366c	Running	t2.micro	Initializing	View
<input type="checkbox"/>	worker-2	i-0e8beb4bab32bb293	Running	t2.micro	Initializing	View
<input type="checkbox"/>	worker-3	i-0ad704152037a39af	Running	t2.micro	Initializing	View
<input type="checkbox"/>	worker-4	i-0640fbfde68e3f126	Running	t2.micro	Initializing	View
<input type="checkbox"/>	master	i-0a241455b2a86b298	Running	t2.micro	2/2 checks passed	View

Figura 1: Captura de pantalla de la "AWS Management Console - EC2" que muestra las máquinas del cluster creado con el tutorial

Ejecutar el script de Map Reduce con Hadoop

Finalmente, ejecutamos nuestro algoritmo de Map Reduce. Para esto basta con crear un directorio /data, almacenar nuestro archivo a leer y finalmente ejecutar el ejecutable JAR contenido en nuestra carpeta de Hadoop. Una vez hecho esto, podemos analizar la salida del Map Reduce haciéndole un Cat y pasándolo por Pipe a More para desplegar la información por pantalla:

```
1 hdfs dfs -mkdir /data
2 hdfs dfs -put sw-script-e04.txt /data/
3 hadoop jar hadoop-3.3.6/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.3.6.jar \
4 wordcount /data/sw-script-e04.txt /data/sw-script-e04.wordcount
5 hadoop fs -cat /data/sw-script-e04.wordcount/part-r-00000 | more
```

Código 19: Ejecutar el Map Reduce y leer su salida

1.2. Expansión del cluster

Para expandir el cluster a 8 Workers, lo primero sería replicar las instancias desde AWS, estas debiesen tener los mismos atributos que los workers originales. Lo segundo sería actualizar la configuración desde el Master para incluir las nuevas IP's y habilitar el uso de 4 Nodos desde los XML de Hadoop. Lo tercero sería configurar las conexiones de SSH hacia los workers nuevos, para esto podemos reutilizar el script de Bash pero solo con las IP's privadas nuevas.

```
1 # hadoop-3.3.6/etc/hadoop/mapred-site.xml
2 <property>
3   <name>mapreduce.job.running.map.limit</name>
4   <value>8</value> # Modificar de 4 a 8
5 </property>
6 # ...
7 <property>
8   <name>mapreduce.job.reduces</name>
9   <value>8</value> # Modificar de 4 a 8
10 </property>
11 # ...
12 <property>
13   <name>mapreduce.job.running.reduce.limit</name>
14   <value>8</value> # Modificar de 4 a 8
15 </property>
16
17 # hadoop-3.3.6/etc/hadoop/workers
18 172.31.30.0
19 172.31.17.157
20 172.31.24.51
21 172.31.27.79
22 172.31.84.79
23 172.31.93.162
24 172.31.92.28
25 172.31.81.101
26
27 # /home/ubuntu/install.sh
28 # WORKERS=("172.31.30.0" "172.31.17.157" "172.31.24.51" "172.31.27.79")
29 WORKERS=("172.31.30.0" "172.31.17.157" "172.31.24.51" "172.31.27.79" \
30 "172.31.84.79" "172.31.93.162" "172.31.92.28" "172.31.81.101")
31 MASTER="172.31.30.146"
```

Código 20: Los archivos que hay que modificar antes de ejecutar el script de instalación nuevamente

The screenshot displays the AWS Management Console interface for EC2 instances. The top section shows a list of instances with columns for Name, Instance ID, Instance state, Instance type, Status check, Alarm status, Availability Zone, Public IPv4 DNS, Public IPv4 address, Elastic IP, IPv6 IP, Monitoring, and Security group name. The instances are named 'worker-5' through 'worker-8' and 'master'. The 'worker-5' instance is highlighted.

Below the list, the detailed view for the 'i-0796c3f9461f16202 (worker-8)' instance is shown. The 'Details' tab is active, displaying various instance attributes:

- Instance ID:** i-0796c3f9461f16202
- IPv4 address:** -
- Hostname type:** IP name: ip-172-31-51-101.ec2.internal
- Answer private resource DNS name:** IPv4 (A)
- Auto-assigned IP address:** 3.52.145.145 (Public IP)
- IAM Role:** -
- Instance type:** t2.micro
- VPC ID:** vpc-b449f1a5
- Subnet ID:** subnet-b0a48c20
- Instance ARN:** -
- Public IPv4 address:** 3.52.145.145
- Instance state:** Running
- Private IP DNS name (IPv4 only):** ip-172-31-51-101.ec2.internal
- Private IPv4 addresses:** 172.31.81.101
- Public IPv4 DNS:** ec2-3-52-145-145.compute-1.amazonaws.com
- Elastic IP addresses:** -
- AWS Compute Optimizer finding:** Open to AWS Compute Optimizer for recommendations
- Auto Scaling Group name:** -
- Monitoring:** Managed

Figura 2: Captura de pantalla de la "AWS Management Console - EC2" que muestra las máquinas del cluster expandido

1.3. Cliente web

Aquí podemos apreciar todas las instancias desde la interfaz de EC2 de AWS, las instancias expandidas tienen el nombre de Worker y un número que va desde 5 hasta 8. Sus direcciones IP privadas aparecen en el apartado de expansión y tienen las mismas características que los otros Nodos.

La figura 3 muestra la consola web de Hadoop y la figura 4 la consola web del HDFS.

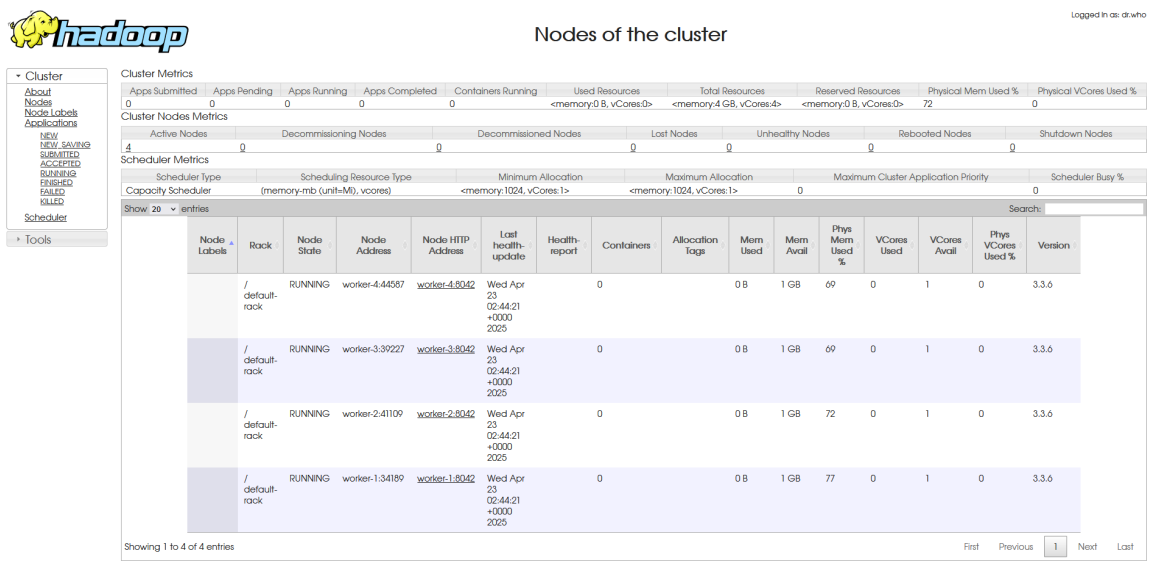


Figura 3: Captura de pantalla del cliente web de Hadoop

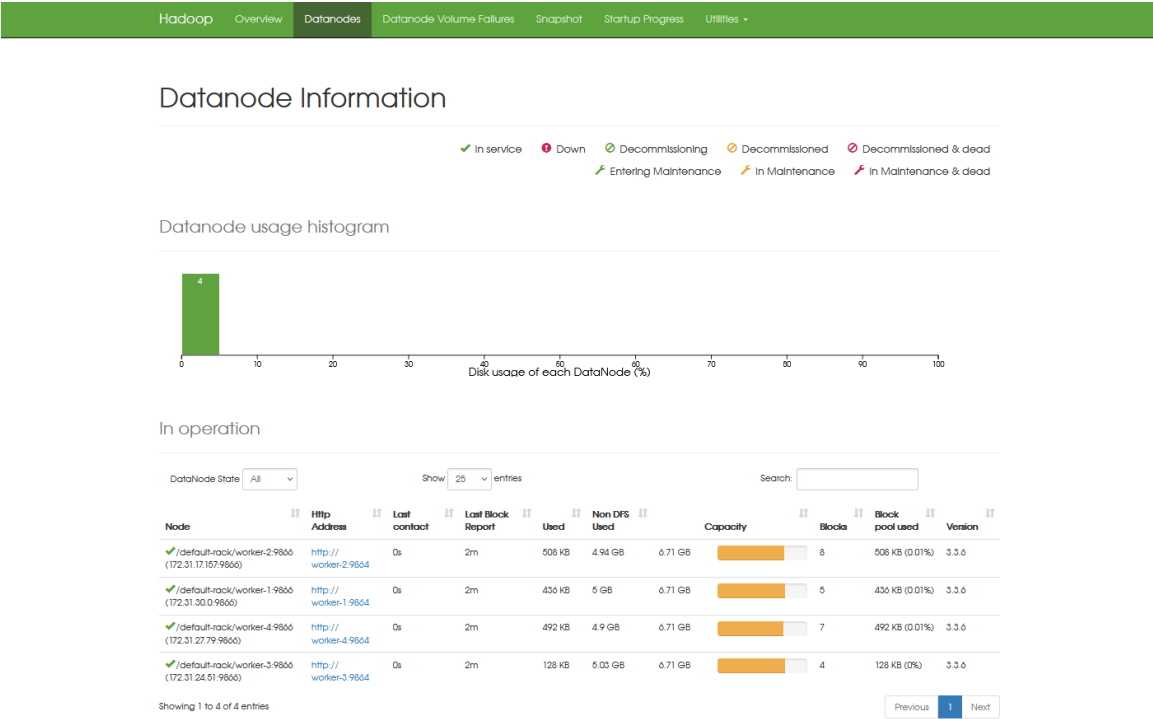


Figura 4: Captura de pantalla del cliente web de HDFS

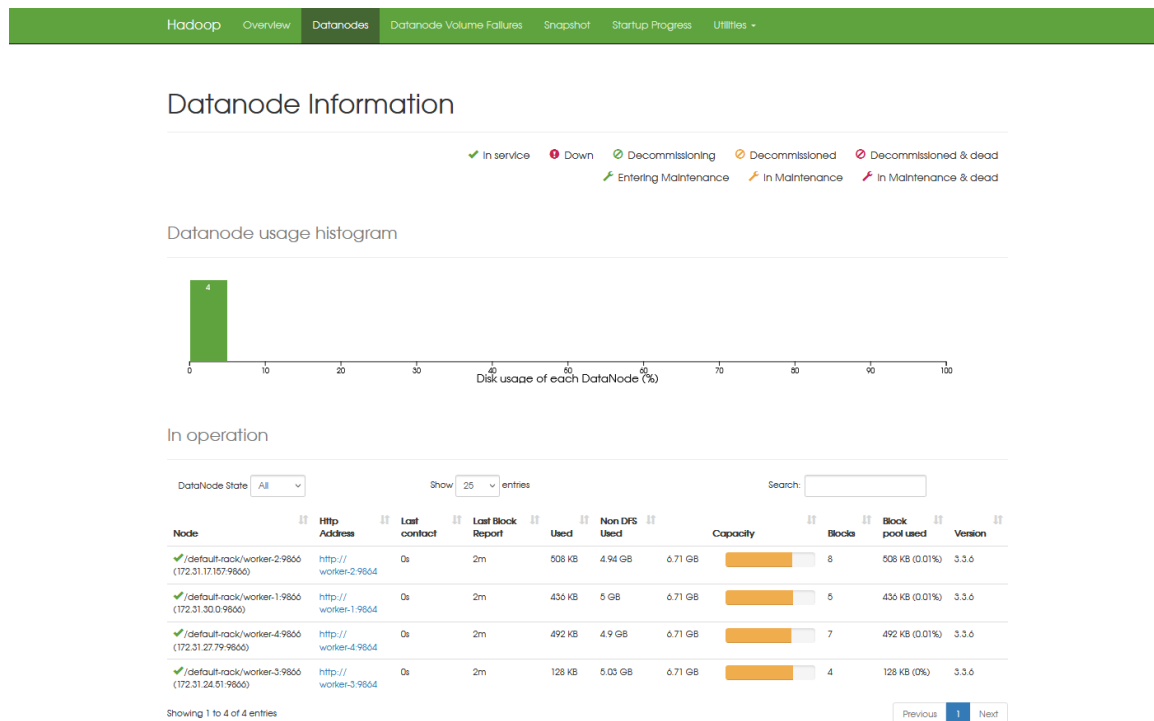


Figura 5: Instancias mejoradas para no sufrir

2. Instalación de Apache Hive

2.1. Procedimiento

Lo primero que debemos hacer para evitar cualquier posterior de cabeza grande, es ampliar las capacidades de las instancias que tenemos construidas en AWS. Las maquinas Worker pasaran a ser de tipo t3.small mientras que el Master a t3.medium. Esta mejora la realizamos desde la interfaz de EC2 de AWS.

Lo segundo que hacemos es ir a la pagina oficial de Apache Hive para descargar la version 4.0.1. Desde el nodo Master tendremos que usar el comando WGET para descargar el comprimido correspondiente:

```
1 wget https://d1cdn.apache.org/hive/hive-4.0.1/apache-hive-4.0.1-bin.tar.gz
```

Código 21: Usar WGET para descargar Apache Hive 4.0.1

Con el archivo ya descargado, lo descomprimos de la misma manera que con Hadoop:

```
1 tar -xvzf apache-hive-4.0.1-bin.tar.gz -C .
```

Código 22: Descomprimos la carpeta de Apache Hive 4.0.1

Luego, para poder utilizar el Apache Hive, debemos crear el archivo de configuracion necesario para que pueda funcionar por sobre nuestra instalacion de Hadoop, para esto creamos el siguiente archivo XML y lo guardamos en /home/ubuntu/apache-hive-4.0.1-bin/conf/

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4   <property>
5     <name>javax.jdo.option.ConnectionURL</name>
6     <value>jdbc:derby:metastore_db;create=true</value>
7   </property>
8   <property>
9     <name>javax.jdo.option.ConnectionDriverName</name>
10    <value>org.apache.derby.jdbc.EmbeddedDriver</value>
11  </property>
12  <property>
13    <name>datanucleus.schema.autoCreateAll</name>
14    <value>true</value>
15  </property>
16  <property>
17    <name>hive.metastore.schema.validation</name>
18    <value>false</value>
19  </property>
20  <property>
21    <name>hive.metastore.warehouse.dir</name>
22    <value>/user/hive/warehouse</value>
23  </property>
24  <property>
25    <name>hive.server2.thrift.port</name>
26    <value>10000</value>
27  </property>
28  <property>
29    <name>hive.server2.thrift.bind.host</name>
30    <value>0.0.0.0</value>
31  </property>
32  <property>
33    <name>hive.server2.authentication</name>
34    <value>NOSASL</value>
35  </property>
36  <property>
37    <name>hive.server2.enable.doAs</name>
38    <value>false</value>
39  </property>
40 </configuration>

```

Código 23: Archivo de configuracion de Apache Hive

Creado dicho archivo, reiniciamos las maquinas.

Con la maquinas ya reiniciadas, levantamos el servicio metastore y de hiveserver2 de la siguiente manera:

```

1 # Levanta ambos servicios en segundo plano y esconde su salida
2 nohup hive --service hiveserver2 > /dev/null 2>&1
3 nohup hive --service metastore > /dev/null 2>&1

```

Código 24: Levantar MetaStore y HiveServer2

Antes de seguir, es importante recordar que Apache Hive funciona por sobre Hadoop, por ende el DFS debe estar funcionando al igual que el negociador de recursos, así que los levantamos manualmente a través de

los scripts que vienen con Hadoop de forma normal, luego, cuando ambos servicios esten funcionando, corremos la shell de Apache Hive usando Beeline:

```
1 beeline -u "jdbc:hive2://localhost:10000/default;auth=noSasl" --verbose=true
```

Código 25: Usar Beeline como interfaz a Hive

2.2. Prueba

Se nos entrega la Query a continuacion para probar el funcionamiento correcto de Hive:

```
1 CREATE DATABASE sw_dialogs;
2
3 USE sw_dialogs;
4
5 CREATE TABLE sw04_dialogs (
6     line      INT,
7     character STRING,
8     dialog    STRING
9 )
10 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
11 WITH SERDEPROPERTIES (
12     "separatorChar" = " ",
13     "quoteChar" = "\""
14 )
15 STORED AS TEXTFILE
16 TBLPROPERTIES ("skip.header.line.count"="1");
17
18 LOAD DATA INPATH
19 '/data/sw-script-e04.txt'
20 INTO TABLE
21 sw04_dialogs;
22
23 SELECT character, COUNT(*) AS lines
24 FROM sw04_dialogs
25 GROUP BY character
26 ORDER BY lines DESC
27 LIMIT 12;
```

Código 26: Ejemplo de uso de Apache Hive

Para corroborar que Hive funciona correctamente, descargue la salida completa tras ejecutar la Query anteriormente mencionada. En las paginas siguientes esta la salida completa tras la ejecucion de dicho comando desde la Shell de Beeline.

```

1 ubuntu@ip-172-31-30-146:~$ beeline -u "jdbc:hive2://localhost:10000/default;auth=noSasl" --
   ↪ verbose=true
2 !connect jdbc:hive2://localhost:10000/default;auth=noSasl '' [passwd stripped]
3 Connecting to jdbc:hive2://localhost:10000/default;auth=noSasl
4 Connected to: Apache Hive (version 4.0.1)
5 Driver: Hive JDBC (version 4.0.1)
6 Transaction isolation: TRANSACTION_REPEATABLE_READ
7 Beeline version 4.0.1 by Apache Hive
8 0: jdbc:hive2://localhost:10000/default> show databases;
9 going to print operations logs
10 printed operations logs
11 going to print operations logs
12 printed operations logs
13 going to print operations logs
14 printed operations logs
15 going to print operations logs
16 INFO : Compiling command(queryId=ubuntu_20250501025436_a463ca5e-c723-4037-a672-4465
   ↪ e790fae9): show databases
17 INFO : Semantic Analysis Completed (retrial = false)
18 INFO : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:database_name, type:
   ↪ string, comment:from deserializer)], properties:null)
19 INFO : Completed compiling command(queryId=ubuntu_20250501025436_a463ca5e-c723-4037-a672
   ↪ -4465e790fae9); Time taken: 2.378 seconds
20 INFO : Concurrency mode is disabled, not creating a lock manager
21 INFO : Executing command(queryId=ubuntu_20250501025436_a463ca5e-c723-4037-a672-4465
   ↪ e790fae9): show databases
22 INFO : Starting task [Stage-0:DDL] in serial mode
23 INFO : Completed executing command(queryId=ubuntu_20250501025436_a463ca5e-c723-4037-a672
   ↪ -4465e790fae9); Time taken: 0.217 seconds
24 printed operations logs
25 Getting log thread is interrupted, since query is done!
26 +-----+
27 | database_name |
28 +-----+
29 | default       |
30 +-----+
31 1 row selected (3.366 seconds)
32 0: jdbc:hive2://localhost:10000/default> !run hive-script.hql
33 >>> CREATE DATABASE sw_dialogs;
34 going to print operations logs
35 printed operations logs
36 Getting log thread is interrupted, since query is done!
37 INFO : Compiling command(queryId=ubuntu_20250501025448_df7a8fbc-1987-45fb-ad48-
   ↪ d61a6bb732a5): CREATE DATABASE sw_dialogs
38 INFO : Semantic Analysis Completed (retrial = false)
39 INFO : Created Hive schema: Schema(fieldSchemas:null, properties:null)
40 INFO : Completed compiling command(queryId=ubuntu_20250501025448_df7a8fbc-1987-45fb-ad48-
   ↪ d61a6bb732a5); Time taken: 0.009 seconds
41 INFO : Concurrency mode is disabled, not creating a lock manager
42 INFO : Executing command(queryId=ubuntu_20250501025448_df7a8fbc-1987-45fb-ad48-
   ↪ d61a6bb732a5): CREATE DATABASE sw_dialogs
43 INFO : Starting task [Stage-0:DDL] in serial mode
44 INFO : Completed executing command(queryId=ubuntu_20250501025448_df7a8fbc-1987-45fb-ad48-
   ↪ d61a6bb732a5); Time taken: 0.224 seconds
45 No rows affected (0.253 seconds)
46 >>>
47 >>> USE sw_dialogs;
48 going to print operations logs
49 printed operations logs

```

Código 27: Salida de Apache Hive (Parte 1)

```

1 Getting log thread is interrupted, since query is done!
2 INFO : Compiling command(queryId=ubuntu_20250501025448_a7a58e45-3065-4b8d-987f-382
    ↪ f6be6ac0d): USE sw_dialogs
3 INFO : Semantic Analysis Completed (retrial = false)
4 INFO : Created Hive schema: Schema(fieldSchemas:null, properties:null)
5 INFO : Completed compiling command(queryId=ubuntu_20250501025448_a7a58e45-3065-4b8d-987f
    ↪ -382f6be6ac0d); Time taken: 0.026 seconds
6 INFO : Concurrency mode is disabled, not creating a lock manager
7 INFO : Executing command(queryId=ubuntu_20250501025448_a7a58e45-3065-4b8d-987f-382
    ↪ f6be6ac0d): USE sw_dialogs
8 INFO : Starting task [Stage-0:DDL] in serial mode
9 INFO : Completed executing command(queryId=ubuntu_20250501025448_a7a58e45-3065-4b8d-987f
    ↪ -382f6be6ac0d); Time taken: 0.041 seconds
10 No rows affected (0.108 seconds)
11 >>>
12 >>> CREATE TABLE sw04_dialogs (
13     line      INT,
14     character STRING,
15     dialog    STRING
16 )
17 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
18 WITH SERDEPROPERTIES (
19     "separatorChar" = " ",
20     "quoteChar" = "\""
21 )
22 STORED AS TEXTFILE
23 TBLPROPERTIES ("skip.header.line.count"="1");
24 going to print operations logs
25 printed operations logs
26 Getting log thread is interrupted, since query is done!
27 INFO : Compiling command(queryId=ubuntu_20250501025448_7e5cb1f6-8c26-4e7c-aa1e-7
    ↪ a905630ea48): CREATE TABLE sw04_dialogs (
28     line      INT,
29     character STRING,
30     dialog    STRING
31 )
32 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
33 WITH SERDEPROPERTIES (
34     "separatorChar" = " ",
35     "quoteChar" = "\""
36 )
37 STORED AS TEXTFILE
38 TBLPROPERTIES ("skip.header.line.count"="1")
39 INFO : Semantic Analysis Completed (retrial = false)
40 INFO : Created Hive schema: Schema(fieldSchemas:null, properties:null)
41 INFO : Completed compiling command(queryId=ubuntu_20250501025448_7e5cb1f6-8c26-4e7c-aa1e-7
    ↪ a905630ea48); Time taken: 0.102 seconds
42 INFO : Concurrency mode is disabled, not creating a lock manager
43 INFO : Executing command(queryId=ubuntu_20250501025448_7e5cb1f6-8c26-4e7c-aa1e-7
    ↪ a905630ea48): CREATE TABLE sw04_dialogs (
44     line      INT,
45     character STRING,
46     dialog    STRING
47 )
48 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
49 WITH SERDEPROPERTIES (
50     "separatorChar" = " ",
51     "quoteChar" = "\""
52 )
53 STORED AS TEXTFILE
54 TBLPROPERTIES ("skip.header.line.count"="1")
55 INFO : Starting task [Stage-0:DDL] in serial mode
56 INFO : Completed executing command(queryId=ubuntu_20250501025448_7e5cb1f6-8c26-4e7c-aa1e-7
    ↪ a905630ea48); Time taken: 0.724 seconds

```

```

1 No rows affected (0.866 seconds)
2 >>>
3 >>> LOAD DATA INPATH
4 '/data/sw-script-e04.txt'
5 INTO TABLE
6 sw04_dialogs;
7 going to print operations logs
8 printed operations logs
9 going to print operations logs
10 INFO : Compiling command(queryId=ubuntu_20250501025449_c327dfb8-5749-408f-b787-70
    ↳ bc8925a884): LOAD DATA INPATH
11 '/data/sw-script-e04.txt'
12 INTO TABLE
13 sw04_dialogs
14 INFO : Semantic Analysis Completed (retrial = false)
15 INFO : Created Hive schema: Schema(fieldSchemas:null, properties:null)
16 INFO : Completed compiling command(queryId=ubuntu_20250501025449_c327dfb8-5749-408f-b787
    ↳ -70bc8925a884); Time taken: 0.406 seconds
17 INFO : Concurrency mode is disabled, not creating a lock manager
18 INFO : Executing command(queryId=ubuntu_20250501025449_c327dfb8-5749-408f-b787-70
    ↳ bc8925a884): LOAD DATA INPATH
19 '/data/sw-script-e04.txt'
20 INTO TABLE
21 sw04_dialogs
22 INFO : Starting task [Stage-0:MOVE] in serial mode
23 INFO : Loading data to table sw_dialogs.sw04_dialogs from hdfs://172.31.30.146/data/sw-
    ↳ script-e04.txt
24 INFO : Starting task [Stage-1:STATS] in serial mode
25 INFO : Executing stats task
26 INFO : Table sw_dialogs.sw04_dialogs stats: [numFiles=2, totalSize=156556,
    ↳ numFilesErasureCoded=0]
27 INFO : Completed executing command(queryId=ubuntu_20250501025449_c327dfb8-5749-408f-b787
    ↳ -70bc8925a884); Time taken: 0.866 seconds
28 printed operations logs
29 Getting log thread is interrupted, since query is done!
30 No rows affected (1.308 seconds)
31 >>>
32 >>> SELECT character, COUNT(*) AS lines
33 FROM sw04_dialogs
34 GROUP BY character
35 ORDER BY lines DESC
36 LIMIT 12;
37 going to print operations logs
38 printed operations logs
39 going to print operations logs
40 printed operations logs
41 going to print operations logs
42 printed operations logs
43 going to print operations logs
44 printed operations logs
45 going to print operations logs
46 INFO : Compiling command(queryId=ubuntu_20250501025450_36332ece-1f44-45b7-998f-69
    ↳ c953bdb7af): SELECT character, COUNT(*) AS lines
47 FROM sw04_dialogs
48 GROUP BY character
49 ORDER BY lines DESC
50 LIMIT 12

```

Código 29: Salida de Apache Hive (Parte 3)


```

1 INFO : No Stats for sw_dialogs@sw04_dialogs, Columns: character
2 INFO : Semantic Analysis Completed (retrial = false)
3 INFO : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:character, type:string,
    ↳ comment:null), FieldSchema(name:lines, type:bigint, comment:null)], properties:null)
4 INFO : Completed compiling command(queryId=ubuntu_20250501025450_36332ece-1f44-45b7-998f
    ↳ -69c953bdb7af); Time taken: 3.888 seconds
5 INFO : Concurrency mode is disabled, not creating a lock manager
6 INFO : Executing command(queryId=ubuntu_20250501025450_36332ece-1f44-45b7-998f-69
    ↳ c953bdb7af): SELECT character, COUNT(*) AS lines
7 FROM sw04_dialogs
8 GROUP BY character
9 ORDER BY lines DESC
10 LIMIT 12
11 WARN : Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions.
    ↳ Consider using a different execution engine (i.e. tez) or using Hive 1.X releases.
12 INFO : Query ID = ubuntu_20250501025450_36332ece-1f44-45b7-998f-69c953bdb7af
13 INFO : Total jobs = 2
14 INFO : Launching Job 1 out of 2
15 INFO : Starting task [Stage-1:MAPRED] in serial mode
16 INFO : Number of reduce tasks not specified. Estimated from input data size: 1
17 INFO : In order to change the average load for a reducer (in bytes):
18 INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
19 INFO : In order to limit the maximum number of reducers:
20 INFO :   set hive.exec.reducers.max=<number>
21 INFO : In order to set a constant number of reducers:
22 INFO :   set mapreduce.job.reduces=<number>
23 printed operations logs
24 going to print operations logs
25 printed operations logs
26 going to print operations logs
27 printed operations logs
28 going to print operations logs
29 printed operations logs
30 going to print operations logs
31 INFO : number of splits:2
32 INFO : Submitting tokens for job: job_1746065552179_0001
33 INFO : Executing with tokens: []
34 printed operations logs
35 going to print operations logs
36 INFO : The url to track the job: http://master:8088/proxy/application_1746065552179_0001/
37 INFO : Starting Job = job_1746065552179_0001, Tracking URL = http://master:8088/proxy/
    ↳ application_1746065552179_0001/
38 INFO : Kill Command = /home/ubuntu/hadoop-3.3.6/bin/mapred job -kill
    ↳ job_1746065552179_0001
39 printed operations logs
40 going to print operations logs
41 printed operations logs
42 going to print operations logs
43 printed operations logs
44 going to print operations logs
45 printed operations logs
46 going to print operations logs
47 printed operations logs
48 going to print operations logs

```

Código 30: Salida de Apache Hive (Parte 4)

```
1 printed operations logs
2 going to print operations logs
3 printed operations logs
4 going to print operations logs
5 printed operations logs
6 going to print operations logs
7 printed operations logs
8 going to print operations logs
9 printed operations logs
10 going to print operations logs
11 printed operations logs
12 going to print operations logs
13 INFO : Hadoop job information for Stage-1: number of mappers: 2; number of reducers: 1
14 INFO : 2025-05-01 02:55:12,158 Stage-1 map = 0%, reduce = 0%
15 printed operations logs
16 going to print operations logs
17 printed operations logs
18 going to print operations logs
19 printed operations logs
20 going to print operations logs
21 printed operations logs
22 going to print operations logs
23 printed operations logs
24 going to print operations logs
25 printed operations logs
26 going to print operations logs
27 printed operations logs
28 going to print operations logs
29 INFO : 2025-05-01 02:55:21,611 Stage-1 map = 50%, reduce = 0%, Cumulative CPU 2.27 sec
30 INFO : 2025-05-01 02:55:22,648 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 4.43 sec
31 printed operations logs
32 going to print operations logs
33 printed operations logs
34 going to print operations logs
35 printed operations logs
36 going to print operations logs
37 printed operations logs
38 going to print operations logs
39 printed operations logs
40 going to print operations logs
41 INFO : 2025-05-01 02:55:33,002 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 6.68 sec
42 printed operations logs
43 going to print operations logs
44 INFO : MapReduce Total cumulative CPU time: 6 seconds 680 msec
45 INFO : Ended Job = job_1746065552179_0001
46 INFO : Launching Job 2 out of 2
47 INFO : Starting task [Stage-2:MAPRED] in serial mode
48 INFO : Number of reduce tasks determined at compile time: 1
49 INFO : In order to change the average load for a reducer (in bytes):
50 INFO : set hive.exec.reducers.bytes.per.reducer=<number>
```

Código 31: Salida de Apache Hive (Parte 5)

```

1 INFO : In order to limit the maximum number of reducers:
2 INFO :   set hive.exec.reducers.max=<number>
3 INFO : In order to set a constant number of reducers:
4 INFO :   set mapreduce.job.reduces=<number>
5 INFO : number of splits:1
6 INFO : Submitting tokens for job: job_1746065552179_0002
7 INFO : Executing with tokens: []
8 INFO : The url to track the job: http://master:8088/proxy/application_1746065552179_0002/
9 INFO : Starting Job = job_1746065552179_0002, Tracking URL = http://master:8088/proxy/
    ↪ application_1746065552179_0002/
10 INFO : Kill Command = /home/ubuntu/hadoop-3.3.6/bin/mapred job -kill
    ↪ job_1746065552179_0002
11 printed operations logs
12 going to print operations logs
13 printed operations logs
14 going to print operations logs
15 printed operations logs
16 going to print operations logs
17 printed operations logs
18 going to print operations logs
19 printed operations logs
20 going to print operations logs
21 printed operations logs
22 going to print operations logs
23 INFO : Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
24 INFO : 2025-05-01 02:55:50,592 Stage-2 map = 0%, reduce = 0%
25 printed operations logs
26 going to print operations logs
27 printed operations logs
28 going to print operations logs
29 printed operations logs
30 going to print operations logs
31 INFO : 2025-05-01 02:56:00,914 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 2.66 sec
32 printed operations logs
33 going to print operations logs
34 printed operations logs
35 going to print operations logs
36 printed operations logs
37 going to print operations logs
38 INFO : 2025-05-01 02:56:12,332 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 5.7 sec
39 INFO : MapReduce Total cumulative CPU time: 5 seconds 700 msec
40 INFO : Ended Job = job_1746065552179_0002
41 INFO : MapReduce Jobs Launched:
42 INFO : Stage-Stage-1: Map: 2 Reduce: 1 Cumulative CPU: 6.68 sec HDFS Read: 179440
    ↪ HDFS Write: 1807 HDFS EC Read: 0 SUCCESS
43 INFO : Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 5.7 sec HDFS Read: 10647 HDFS
    ↪ Write: 348 HDFS EC Read: 0 SUCCESS
44 INFO : Total MapReduce CPU Time Spent: 12 seconds 380 msec
45 INFO : Completed executing command(queryId=ubuntu_20250501025450_36332ece-1f44-45b7-998f
    ↪ -69c953bdb7af); Time taken: 78.686 seconds
46 printed operations logs
47 Getting log thread is interrupted, since query is done!

```

Código 32: Salida de Apache Hive (Parte 6)

```

1 +-----+-----+
2 | character | lines |
3 +-----+-----+
4 | LUKE      | 508   |
5 | HAN       | 306   |
6 | THREEPIO  | 238   |
7 | BEN       | 164   |
8 | LEIA      | 114   |
9 | VADER     | 82    |
10 | RED LEADER | 74    |
11 | BIGGS     | 68    |
12 | TARKIN    | 56    |
13 | OWEN      | 50    |
14 | TROOPER   | 38    |
15 | WEDGE     | 28    |
16 +-----+-----+
17 12 rows selected (82.764 seconds)

```

Código 33: Salida de Apache Hive (Parte 7)

Mientras que la salida de referencia era:

```

1 +-----+-----+
2 | character | lines |
3 +-----+-----+
4 | LUKE      | 524   |
5 | HAN       | 135   |
6 | THREEPIO  | 119   |
7 | BEN       | 82    |
8 | LEIA      | 57    |
9 | VADER     | 41    |
10 | RED LEADER | 73    |
11 | BIGGS     | 34    |
12 | TARKIN    | 28    |
13 | OWEN      | 25    |
14 | TROOPER   | 19    |
15 | WEDGE     | 14    |
16 +-----+-----+
17 12 rows selected (67.678 seconds)

```

Código 34: Resultado esperado para ejemplo de uso de Apache Hive

3. Exploración del HDFS

Para desarrollar el Script que muestre algo similar a lo que sale a continuacion, me aproveche de la salida que tuve. En mi caso, se dio perfectamente que Hadoop guardo el archivo en un unico bloque, por ende hice un Script en Bash que funcionara para mi situacion especifica. De no ser así el caso, hubiese guardado cada IP por bloque en archivos separados y luego guardarlos en un arreglo usano un pipe desde LS.

```

1 File:           /mydata/myfile.txt
2 Blocks:         2
3 Avg. Replication: 3.0
4
5 0 - blk_1073742511_1694 - 172.31.40.100 172.31.32.111 172.31.35.247
6 1 - blk_1073742511_1695 - 172.31.32.111 172.31.210.111 172.31.35.247

```

Código 35: Ejemplo de uso de Apache Hive

El Script que hice fue el siguiente:

```

1 #!/bin/bash
2 HDFS_FILE="/user/hive/warehouse/sw_dialogs.db/sw04_dialogs/sw-script-e04.txt"
3
4 hdfs fsck $HDFS_FILE -files -blocks -locations > fsck_output
5
6 grep -o '[0-9]*\.[0-9]*\.[0-9]*\.[0-9]*' fsck_output.txt > IP_with_Star_Wars.txt
7 grep -o 'blk_[0-9]*_[0-9]*' fsck_output.txt > Blocks_with_Star_Wars.txt
8 grep "Average_block_replication:" fsck_output.txt | grep -o "[0-9\.]*" > Rep_Fact_Star_Wars
   ↪ .txt
9
10 BLOCKS_QTY="$(wc -l Blocks_with_Star_Wars.txt | grep -o '[0-9]*')"
11 AVG_REP="$(cat Rep_Fact_Star_Wars.txt)"
12
13 COUNTER=0
14
15 echo "File: $HDFS_FILE"
16 echo "Blocks: $BLOCKS_QTY"
17 echo "Avg. Replication: $AVG_REP"
18
19 cat Blocks_with_Star_Wars.txt | while read line
20 do
21     echo "$COUNTER-$line-$(cat Blocks_with_Star_Wars.txt)"
22     ((counter++))
23 done

```

Código 36: Muestra el archivo, los bloques que ocupa, el factor de replicacion y la ubicacion de cada bloque

Para descargar el archivo solicitado, segui las instrucciones para instalar AWS CLI indicadas en la documentación de AWS (<https://docs.aws.amazon.com/cli/latest/userguide/getting-started-install.html>) y luego descargue el archivo usando el ultimo comando:

```
1 # Instalacion de AWS CLI
2
3 # Descargar archivo comprimido
4 curl "https://awscli.amazonaws.com/awscli-exe-linux-x86_64.zip" -o "awscliv2.zip"
5 # Descomprimir
6 unzip awscliv2.zip
7 # Instalar
8 sudo ./aws/install
9
10 # Descargar el CSV indicado
11 aws s3 cp s3://utfsm-inf356-dataset/vlt_observations_000.csv vlt_observations_000.csv --no-
    ↪ sign-request
```

Código 37: Instalar AWS CLI

4. Uso del cluster

4.1. Importación

Primero, para determinar el orden correcto de los datos, los insertare de forma directa a la Base de Datos desde Hive, para esto hice una Query generada por un Script en donde cada parametro recibe el nombre de parametro_n, donde n es el numero de la columna, el cual se obtiene contando la cantidad de comas en la primera fila del archivo CSV. De forma preliminar todas las columnas seran de tipo String, ya que no tenemos forma de saber el tipo de dato de antemano.

```

1 # Campos que tendra el archivo. Como desconocemos el contenido, le asignamos parametros
  ↳ genericos
2
3 # Exportamos cada linea del Script de manera dinamica
4 echo "CREATE_DATABASE_dirty_database;" > dirty_data_script.hql
5 echo "USE_dirty_database;" >> dirty_data_script.hql
6 echo "CREATE_TABLE_dirty_table(" >> dirty_data_script.hql
7
8 # Obtenemos la cantidad de parametros extrayendo una fila aleatoria del archivo
9 # al que posteriormente se le remplazaron las comas por saltos de linea
10 grep -Eo "," vlt_observations_000_single_row.csv > comas
11 N_OF_PARAMS=$(wc -l < comas)
12
13 # Contador para nombrar cada parametro
14 COUNTER=0
15
16 # Guarda cada parametro generico como un String, posteriormente los
17 # segregaremos en base a las salidas que nos dan las Querys de esta
18 # Base de Datos Generica
19 while [ "$COUNTER" -le "$N_OF_PARAMS" ]; do
20     echo "UUUUUUparametro_$COUNTER_string," >> dirty_data_script.hql
21     COUNTER=$((COUNTER+1))
22 done
23
24 # Quitar la ultima Coma que queda molestando
25 truncate -s-2 dirty_data_script.hql
26
27 echo "U)" >> dirty_data_script.hql
28 echo "COMMENT_'Datos_Genericos'" >> dirty_data_script.hql
29 echo "ROW_FORMAT_DELIMITED" >> dirty_data_script.hql
30 echo "FIELDS_TERMINATED_BY',';" >> dirty_data_script.hql
31
32 # Finalmente cargamos el CSV en nuestro HDFS para rellenar los datos
33 echo "LOAD_DATA_INPATH_'/data/vlt_observations_000.csv'_INTO_TABLE_vlt_obs_generic;" >>
  ↳ dirty_data_script.hql
34
35 # Subir el archivo a nuestro DFS
36 hdfs dfs -put vlt_observations_000.csv /data/

```

Código 38: Importacion de los datos sucios

```
1 CREATE DATABASE dirty_database;
2 USE dirty_database;
3 CREATE TABLE dirty_table (
4     parametro_0 string,
5     parametro_1 string,
6     parametro_2 string,
7     parametro_3 string,
8     parametro_4 string,
9     parametro_5 string,
10    parametro_6 string,
11    parametro_7 string,
12    parametro_8 string,
13    parametro_9 string,
14    parametro_10 string,
15    parametro_11 string,
16    parametro_12 string,
17    parametro_13 string,
18    parametro_14 string,
19    parametro_15 string,
20    parametro_16 string,
21    parametro_17 string,
22    parametro_18 string,
23    parametro_19 string,
24    parametro_20 string,
25    parametro_21 string,
26    parametro_22 string,
27    parametro_23 string,
28    parametro_24 string,
29    parametro_25 string,
30    parametro_26 string,
31    parametro_27 string,
32    parametro_28 string,
33    parametro_29 string,
34    parametro_30 string,
35    parametro_31 string,
36    parametro_32 string
37 )
38 COMMENT 'Datos_Genericos'
39 ROW FORMAT DELIMITED
40 FIELDS TERMINATED BY ',';
41 LOAD DATA INPATH '/data/vlt_observations_000.csv' INTO TABLE dirty_table;
```

Código 39: Importacion de los datos sucios

Una vez con la Base de Datos generica creada, realizamos una Query para corroborar que los datos se insertaron correctamente.

4.2. Parsing

Al revizar los datos, notamos inmediatamente que el Script toma las comas que estan dentro de comillas dobles. Para solucionar esto cree un Script en Python que se encarga de borrar todas las comas que se encuentran dentro de Strings con comillas dobles.

```

1 import csv
2 import re
3
4 # Lee el archivo original y guarda la correccion en otro archivo
5 with open("vlt_observations_000.csv") as fin, open("vlt_observations_000_clean.csv", "w")
   ↪ as fout:
6     for line in fin:
7         # Reemplaza comas dentro de paréntesis por otro carácter temporal
8         line_fixed = re.sub(r'\([^(]*\)\'', lambda m: m.group(0).replace(",", " "), line)
9         # Muestra por pantalla la linea a corregir
10        print(f"Linea a corregir: {line_fixed}")
11        # Escribe en el archivo la linea corregida
12        fout.write(line_fixed)

```

Código 40: Limpieza de los datos sucios

Con los datos arreglados, realizamos otra tabla ya con los datos arreglados. Con esta tabla nueva creada, realizamos una Query para cada parametro en donde se devuelven los 100 valores unicos por cada columna, de esta manera podemos saber que significa cada columna del CSV original para poder crear una ultima tabla con cada variable con un nombre y tipo mas apropiado.

```

1 CREATE DATABASE dirty_database;
2 USE dirty_database;
3 CREATE TABLE clean_table (
4     parametro_0 string,
5     parametro_1 string,
6     parametro_2 string,
7     ...
8     parametro_30 string,
9     parametro_31 string,
10    parametro_32 string
11 )
12 COMMENT 'Datos Genericos'
13 ROW FORMAT DELIMITED
14 FIELDS TERMINATED BY ',';
15 LOAD DATA INPATH '/data/vlt_observations_000_clean.csv' INTO TABLE clean_table;

```

Código 41: Importacion de los datos sucios

```
1 USE dirty_database;
2 SELECT DISTINCT parametro_0 FROM clean_table LIMIT 100;
3 SELECT DISTINCT parametro_1 FROM clean_table LIMIT 100;
4 SELECT DISTINCT parametro_2 FROM clean_table LIMIT 100;
5 SELECT DISTINCT parametro_3 FROM clean_table LIMIT 100;
6 SELECT DISTINCT parametro_4 FROM clean_table LIMIT 100;
7 SELECT DISTINCT parametro_5 FROM clean_table LIMIT 100;
8 SELECT DISTINCT parametro_6 FROM clean_table LIMIT 100;
9 SELECT DISTINCT parametro_7 FROM clean_table LIMIT 100;
10 SELECT DISTINCT parametro_8 FROM clean_table LIMIT 100;
11 SELECT DISTINCT parametro_9 FROM clean_table LIMIT 100;
12 SELECT DISTINCT parametro_10 FROM clean_table LIMIT 100;
13 SELECT DISTINCT parametro_11 FROM clean_table LIMIT 100;
14 SELECT DISTINCT parametro_12 FROM clean_table LIMIT 100;
15 SELECT DISTINCT parametro_13 FROM clean_table LIMIT 100;
16 SELECT DISTINCT parametro_14 FROM clean_table LIMIT 100;
17 SELECT DISTINCT parametro_15 FROM clean_table LIMIT 100;
18 SELECT DISTINCT parametro_16 FROM clean_table LIMIT 100;
19 SELECT DISTINCT parametro_17 FROM clean_table LIMIT 100;
20 SELECT DISTINCT parametro_18 FROM clean_table LIMIT 100;
21 SELECT DISTINCT parametro_19 FROM clean_table LIMIT 100;
22 SELECT DISTINCT parametro_20 FROM clean_table LIMIT 100;
23 SELECT DISTINCT parametro_21 FROM clean_table LIMIT 100;
24 SELECT DISTINCT parametro_22 FROM clean_table LIMIT 100;
25 SELECT DISTINCT parametro_23 FROM clean_table LIMIT 100;
26 SELECT DISTINCT parametro_24 FROM clean_table LIMIT 100;
27 SELECT DISTINCT parametro_25 FROM clean_table LIMIT 100;
28 SELECT DISTINCT parametro_26 FROM clean_table LIMIT 100;
29 SELECT DISTINCT parametro_27 FROM clean_table LIMIT 100;
30 SELECT DISTINCT parametro_28 FROM clean_table LIMIT 100;
31 SELECT DISTINCT parametro_29 FROM clean_table LIMIT 100;
32 SELECT DISTINCT parametro_30 FROM clean_table LIMIT 100;
33 SELECT DISTINCT parametro_31 FROM clean_table LIMIT 100;
34 SELECT DISTINCT parametro_32 FROM clean_table LIMIT 100;
```

Código 42: Analisis de los datos limpios para entender que significa cada columna

Como ya hemos automatizado todo a este punto, le pasamos las tablas a nuestra LLM de preferencia y le preguntamos a que corresponden los datos devueltos por la Query.

4.3. Análisis

En base a el analisis entregado por la Inteligencia Artificial, ademas de indagacion y criterio propio, la Query que agrupa los valores de cada columna con su nombre y tipo adecuado queda de la siguiente forma:

```

1  -- Creacion de una Base de Daots con los registros limpios
2  CREATE DATABASE base_de_datos_observaciones;
3  USE base_de_datos_observaciones;
4
5  -- Datos con su significado y tipo correspondiente
6  CREATE TABLE observaciones(
7      tipo_de_calibracion string,
8      tiempo_en_ascencion_recta time,
9      tiempo_en_declinacion time,
10     fecha_y_hora_de_observacion datetime,
11     codigo_asignado_es0 string,
12     investigadores string,
13     tipo_de_observador string,
14     titulo_de_proyecto string,
15     tipo_de_programa string,
16     tipo_de_instrumento string,
17     tipo_de_archivo string,
18     tipo_de_calibracion string,
19     modo_de_observacion string,
20     nombre_unico string,
21     nombre_fits string,
22     fecha_de_observacion date,
23     objeto_observado string,
24     codigo_observatorio string,
25     pipeline_de_reduccion string,
26     tiempo_de_finalizacion time
27     redshift float
28     longitud_onda_central float,
29     longitud_onda_final float,
30     modo_de_adquisicion string,
31     dispersor_usado string,
32     configuracion_espectral string,
33     geometria_de_rendija string,
34     fecha_julian_date double,
35     airmass float,
36     seeing float,
37     null_0 string,
38     null_1 string,
39     null_2 string
40 )
41 -- Darles el formato de lectura separado por comas
42 ROW FORMAT DELIMITED
43 FIELDS TERMINATED BY ',';
44 -- Leer registros desde el archivo subido al DFS
45 LOAD DATA INPATH '/data/vlt_observations_000_clean.csv' INTO TABLE observaciones;

```

Código 43: Observacion con nombres y tipos correctos

```
1 ...
2
3 -- Agrupar en base a los modos de observacion mas comunes
4 SELECT modo_de_observacion, COUNT(*) AS cantidad
5 FROM observaciones
6 GROUP BY modo_de_observacion;
7
8 -- Tiempo promedio en declinacion por cada tipo de instrumento
9 SELECT tipo_de_instrumento, AVG(tiempo_en_declinacion) AS promedio_declinacion
10 FROM observaciones
11 GROUP BY tipo_de_instrumento;
12
13 -- Seeing promedio por hora de observacion
14 SELECT HOUR(fecha_y_hora_de_observacion) AS hora, AVG(seeing) AS promedio_seeing
15 FROM observaciones
16 GROUP BY HOUR(fecha_y_hora_de_observacion)
17 ORDER BY hora;
```

Código 44: Observacion con nombres y tipos correctos