Assignment 1

BUAN 6346

Big Data Analytics

Dr. Jerry F. Perez



Yug Jyotirmay Singh YJS220000

Table of Contents

Chapter 2 - Setup

- 2.1 Getting Started
- 2.2 Run the setup script

Chapter 3 - Access HDFS with Command Line and Hue

- 3.1 Explore the HDFS Command Line Interface
- 3.2 Upload Files to HDFS
- 3.3 View HDFS files
- 3.4 Use the Hue File Browser to browse, view and manage files

Chapter 4 - Run a Yard Job

- 4.1 Explore the YARD Cluster
- 4.2 Submit and Application to the YARN Cluster
- 4.3 View the Application in the Hue Job Browser
- 4.4 View the Application in the YARN UI

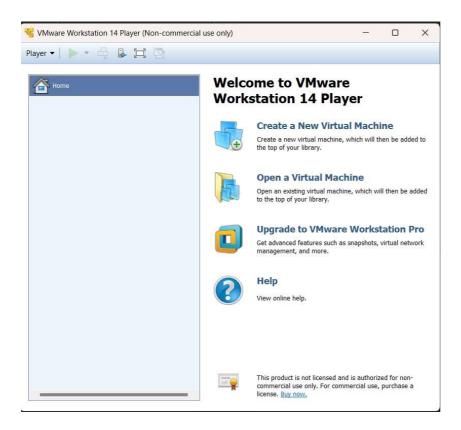
<u>Chapter 5 - Import Data from MySQL Using Sqoop</u>

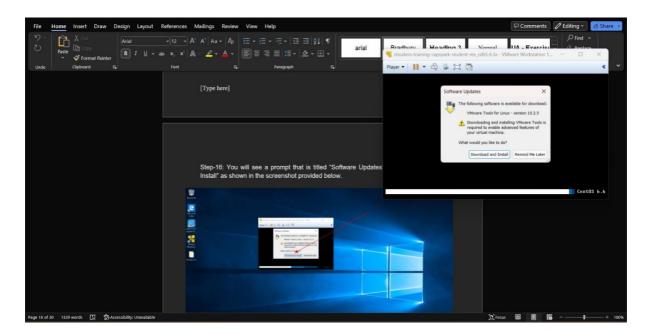
- 5.1 Import the Accounts table from MySQL
- 5.2 View the Imported Data
- 5.3 Import Incremental Updates to Accounts
- 5.4 Import webpage data using an alternate find delimiter

Chapter 2

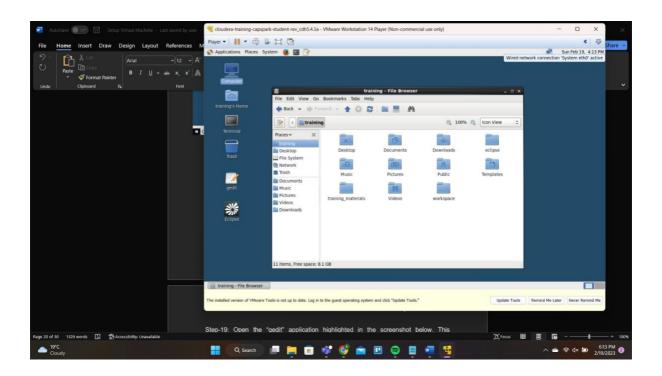
2.1 Getting Started

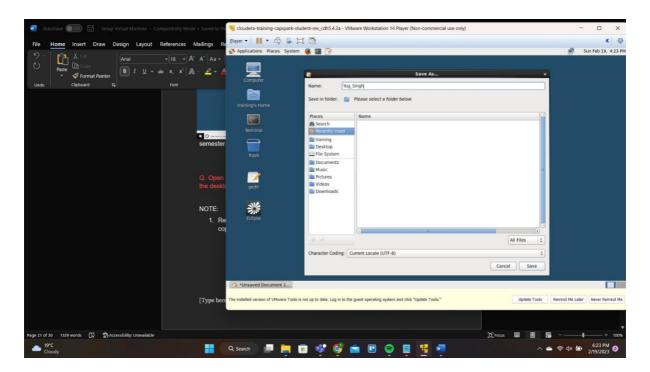
We are now setting up the VM and observing the working of the VM.

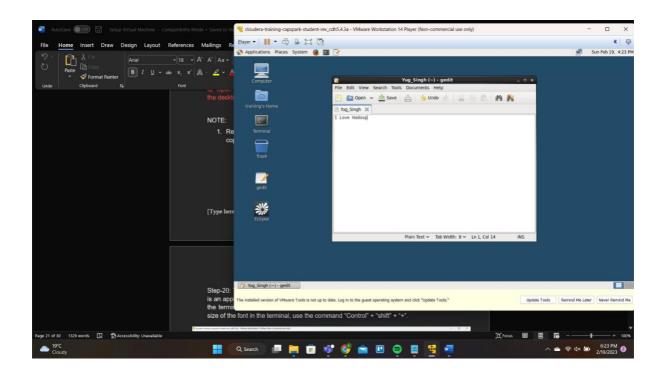


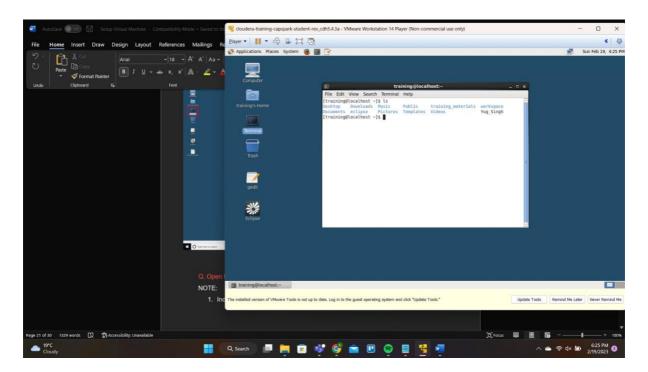


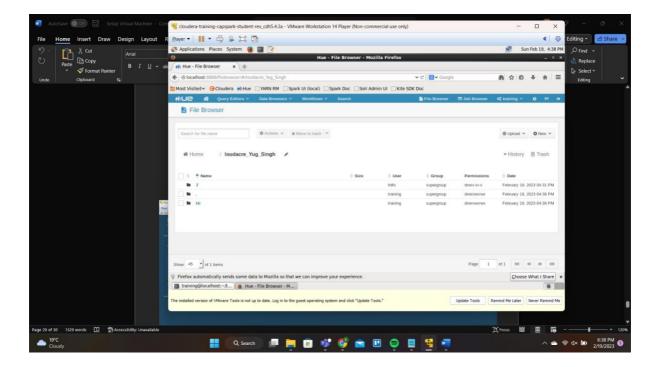
We will upload the files to HDFS before working with them.







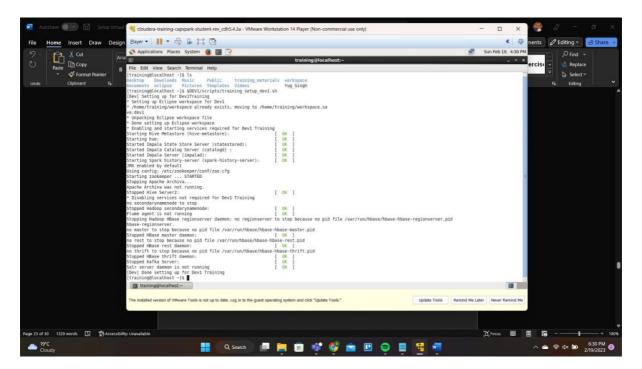


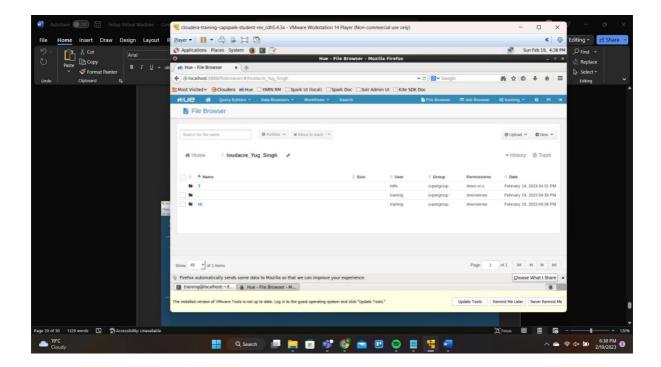


2.2 Run the setup script

The initial task to do here is to set up the environment, and we run the setup script.

\$DEV1/scripts/training_setup_dev1.sh





Chapter 3

3.1 Explore the HDFS Command Line Interface

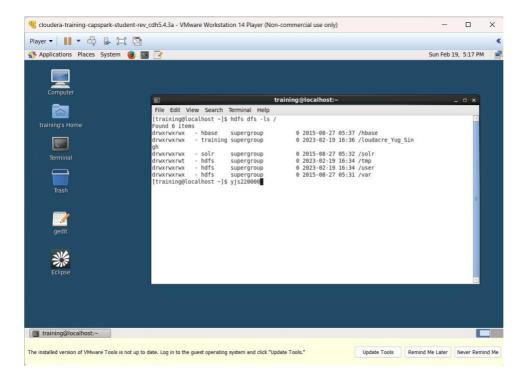
The easiest way to interact with HDFS is using the *hdfs* command.

To execute file system commands within HDFS, use the *hdfs dfs* command.

Use the following steps to explore HDFS CLI:

1.Execute the following command in terminal

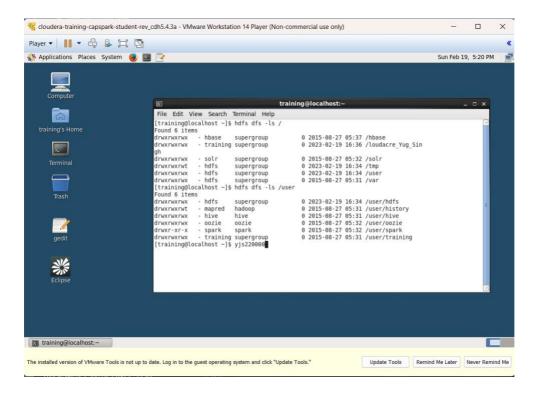
\$hdfs dfs -ls /



We observe that the directory has 5 items.

2. Viewing the contents of the /user directory by running.

\$ hdfs dfs -ls /user



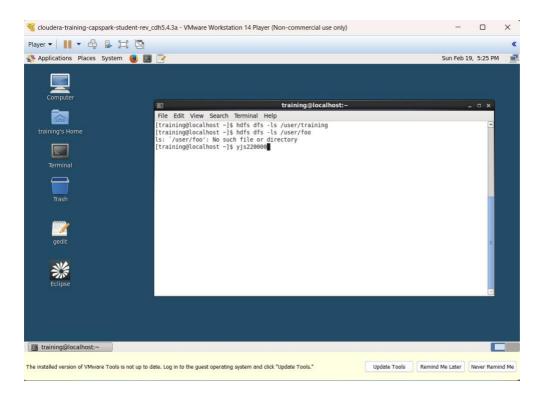
We observe that the user directory has 5 times.

3. Viewing the contents of the home directory by running:

\$ hdfs dfs -ls /user/training

I also tried to look into a non-existent folder which displayed an error message.

\$ hdfs dfs -ls /user/foo



3.2 Upload Files to HDFS

Besides browsing the existing filesystem, another essential thing you can do with the HDFS command-line interface is to upload new data into HDFS.

Use the following steps to achieve this:

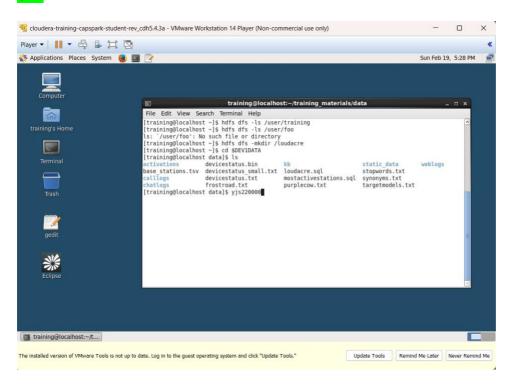
1. Start by creating a new top-level directory for homework assignments, will use this directory throughout the rest of the course.

\$ hdfs dfs -mkdir /loudacre

2. Start by creating a new top-level directory for homework assignments, will use this directory throughout the rest of the course.

\$cd \$DEV1DATA

\$ ls

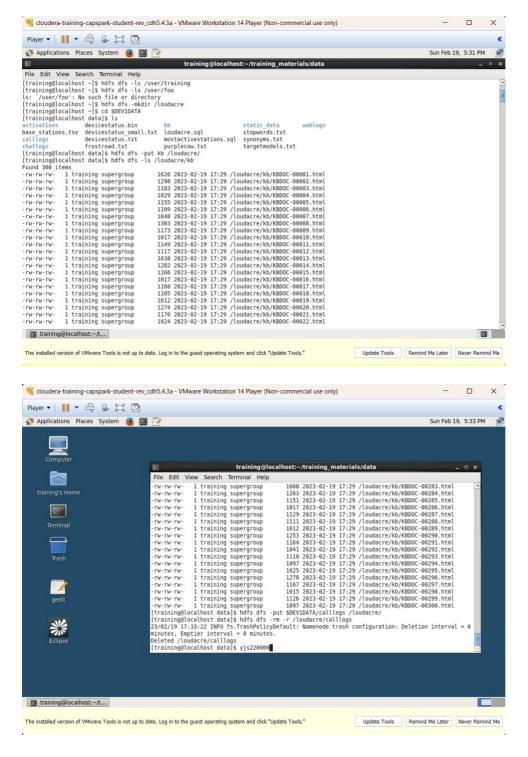


3. Insert this directory into HDFS. This copies the local kb directory and its contents into a remote HDFS directory named /loudacre/kb

\$ hdfs dfs -put kb /loudacre/

4. Listing the files in the new HDFS directory.

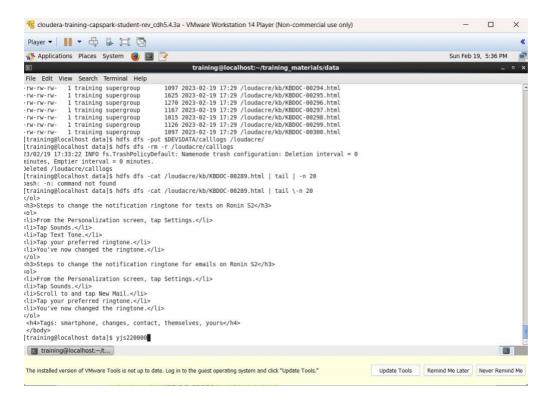
\$ hdfs dfs -ls /loudacre/kb



5. Uploading a directory, then removing it, as it is not needed for future exercises.

\$ hdfs dfs -put \$DEV1DATA/calllogs /loudacre/

\$ hdfs dfs -rm -r /loudacre/calllogs



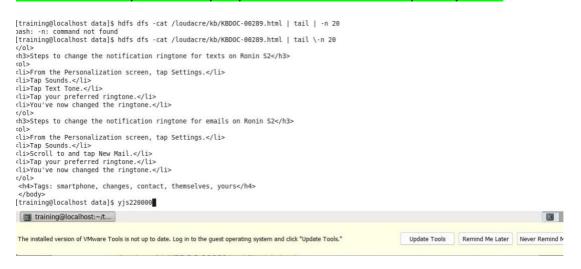
3.3 View HDFS files

Now, let's view some of the data we just copied into HDFS.

Use the following steps to achieve this:

1. The code is:

\$ hdfs dfs -cat /loudacre/kb/KBDOC-00289.html | tail \ -n 20



NOTE: This prints the last 20 lines of the article to your terminal. This command is handy for viewing HDFS data. An individual file is often huge, making it inconvenient to view the entire file in the terminal. For this reason,

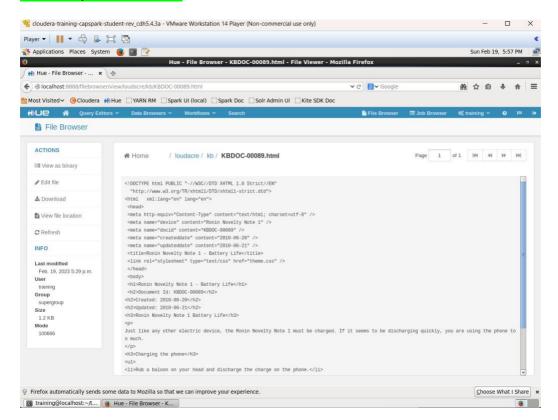
it's often a good idea to pipe the output of the fs -cat command into the head, tail, more, or less.

2. You can use the hdfs dfs -get command to retrieve a local copy of a file or directory from HDFS. This command takes two arguments: an HDFS path and a local path. It copies the HDFS contents into the local filesystem:

\$ hdfs dfs -get \

/loudacre/kb/KBDOC-00289.html ~/article.html

\$ less ~/article.html



3. Several other operations are available with the hdfs dfs command to perform the most common filesystem manipulations: my, cp, mkdir, etc.

The following command provides help:

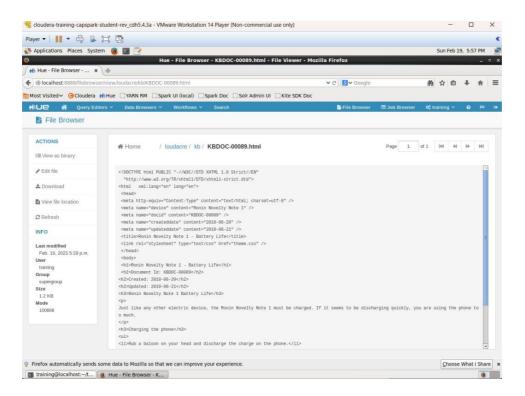
\$hdfs dfs

3.4 Use the Hue File Browser to browse, view and manage files

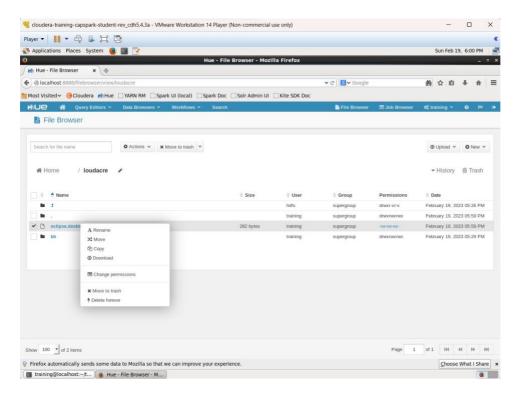
Start Firefox on the VM by clicking its icon on the main menu panel at the top of the screen. Click the Hue bookmark, or visit http://localhost:8888

Because this is the first time anyone has logged into Hue on this server, you will be prompted to create a new user account. Enter username training and password training, then click Create Account. (If started, you may click "Remember Password")

- Note: You may see a misconfiguration warning when you first log in to Hue. This is because not all the services Hue depends on are running on the course VM. You can disregard the message.
- -Hue has many useful features, many of which will be covered later in the course. To access HDFS, click File Browser in the Hue menu bar.



-We can see that the file is modified at a timestamp that matches to our command time.



- The file actions that we can see are shown as above.

Note: Using the upload option, we can upload the files into the file system, and I uploaded the eclipse.desktop file using the upload option.

Chapter 4

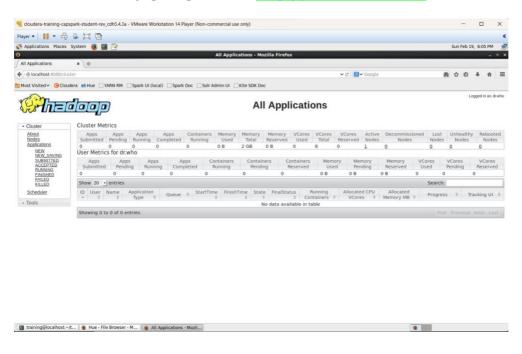
The focus of this exercise is not on what the application does, but on how YARN distributes tasks in a job across a cluster, and how to monitor an application and view its log files.

Files and Data Used in this Homework

Exercise directory:

\$DEV1/exercises/yarn Dataset (HDFS): /loudacre/kb

- 4.1 Explore the YARD Cluster
- 1.Visit the YARN Resource Manager (RM) UI in Firefox using the provided bookmark, or by going to URL http://localhosv/8088



-No jobs are currently running so the current view shows the cluster "at rest".

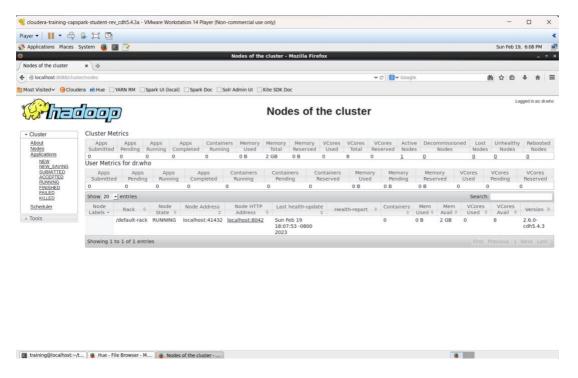
Who is Dr. Who?

You may notice that YARN says you are logged in as dr.who. This is what is displayed when user authentication is disabled for the cluster, as it is on

the training VM. If user authentication were enabled, you would have to log in as a valid user to view the YARN UI, and your actual user name would be displayed, together with user metrics such as how many applications you had run, how much system resources your applications used and so on.

We see that total memory is 2GB, vcores 8 and active nodes 1.
-Vcores are virtual cores, it is a share of host CPU that the YARN Node Manager allocates to available resources

2. Click on the Nodes link in the Cluster menu on the left. The bottom section will display a list of worker nodes in the cluster. The pseudo-distributed cluster used for training has only a single node, which is running on the local machine. In the real world, this list would show multiple worker nodes.



3.Click on the Node HTTP Address to open the Node Manager UI on that specific node. This displays statistics about the selected node, including amount of available memory, currently running applications (none, currently) and so on.



- 4.To return to the Resource Manager, expand Resource Manager on the left and press RM Home.
- 4.2 Submit an Application to the YARN Cluster
- 1. In a terminal window, change to the exercise directory:

\$ cd \$DEV1/exercises/yarn

2.Run the example wordcount.py program on the YARN cluster to count the frequency of words in the knowledge dataset:

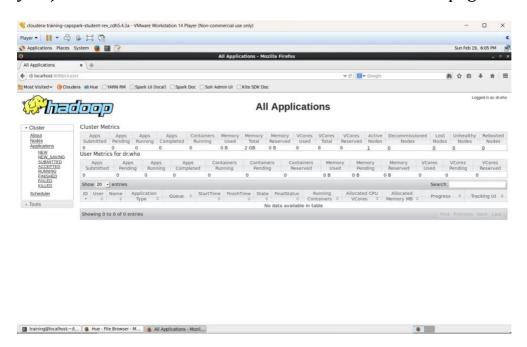
\$ spark-submit --master yarn-cluster \ wordcount.py /loudacre/kb/*

```
[training@localhost yarn]$ spark-submit --master yarn-cluster \ wordcount.py /loudacre/kb/* Exception in thread "main" java.net.URISyntaxException: Illegal character in path at index 0: wordco
           at java.net.URI$Parser.fail(URI.java:2829)
at java.net.URI$Parser.checkChars(URI.java:3002)
           at java.net.URI$Parser.parseHierarchical(URI.java:3086)
           at java.net.URI$Parser.parse(URI.java:3044)
           at java.net.URI.<init>(URI.java:595)
           at org.apache.spark.util.Utils$.resolveURI(Utils.scala:1663)
           at org.apache.spark.deploy.SparkSubmitArguments.parses1(SparkSubmitArguments.scala:432) at org.apache.spark.deploy.SparkSubmitArguments.parse0pts(SparkSubmitArguments.scala:288)
           at org.apache.spark.deploy.SparkSubmitArguments.<init>(SparkSubmitArguments.scala:87) at org.apache.spark.deploy.SparkSubmit$.main(SparkSubmit.scala:105)
           at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/zookeeper/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/Stat
r.class1
SLF4J: Found binding in [jar:file:/usr/lib/flume-ng/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/Stati
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
[training@localhost yarn]$ date
Mon Feb 20 09:58:12 PST 2023
[training@localhost yarn]$ yjs220000
```

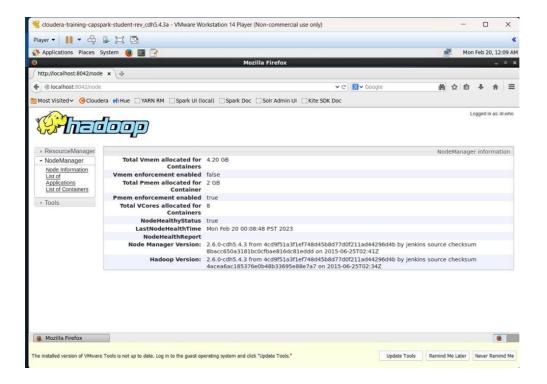
- 4.3 View the Application in the Hue Job Browser
- Go to Hue in Firefox, and select the Job Browser. (Depending on the width of your browser, you may see the whole label, or just the icon.)
- The Job Browser displays a list of currently running and recently completed applications. (If you don't see the application you just started, wait a few seconds, the page will automatically reload; it can take some time for the application to be accepted and start running.)
- This page allows you to click the application ID to see details of the running application, or to kill a running job. (Don't do that now though!)
- 4.4 View the Application in the YARN UI

Note: To get a more detailed view of the cluster, use the YARN UI.

1.Reload the YARN RM page in Firefox. You will now see the application you just started in the bottom section of the RM home page.



2. Select Nodes on the left menu, and Select the Node HTTP Address to open the Node Manager UI.



- 3. Now that an application is running, you can click on List of Applications to see the application you submitted. If it is still running, you will be able to see the containers as shown below.
- It doesn't appear when the whole application is executed, just run the command at 4.2.2 and refresh the browser.



Chapter 5

In this exercise, we will import tables from MySQL into HDFS with Sqoop

Files and Data Used in this Homework

Exercise directory: \$DEV1/exercises/sqoop

MySQL Database: loudacre

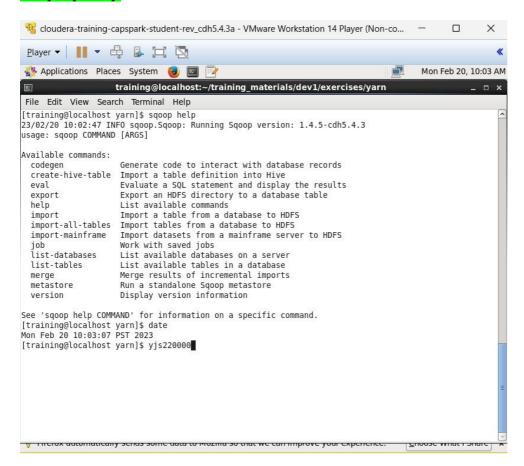
 ${\it MySQL}$ Tables: accounts, webpage

5.1 Import the Accounts table from MySQL

We can use Sqoop to look at the table layout in MySQL. With Sqoop, we can also import the table from MySQL to HDFS.

1. In a terminal window, Run the *sqoop help* command to familiarize ourselves with the options in Sqoop:

\$ sqoop help



2. List the tables in the loudacre database:

\$ sqoop list-tables \

--connect jdbc:mysql://localhost/loudacre \

--username training --password training

```
Try --help for usage instructions.
[training@localhost yarn]$ sqoop list-tables \
> --connect jdbc:mysql://localhost/loudacre \
> --username training --password training
23/02/20 10:05:55 INFO sqoop.Sqoop: Running Sqoop version: 1.4.5-cdh5.4.3
23/02/20 10:05:55 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Con
sider using -P instead.
23/02/20 lo:05:56 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
accountdevice
accounts
basestations
calllog
customerservicerep
device
knowledgebase
mostactivestations
toparticles
webpage
websitehit
[training@localhost yarn]$ date
Mon Feb 20 10:06:01 PST 2023
[training@localhost yarn]$ yjs220000
  Therox automatically serias some data to Prozina so that we cult improve your experience.
```

- We see the tables are in the loudacre database:

accountdevice

accounts

basestations

calllog

customerservicerep

device

knowledgebase

mostactivestations

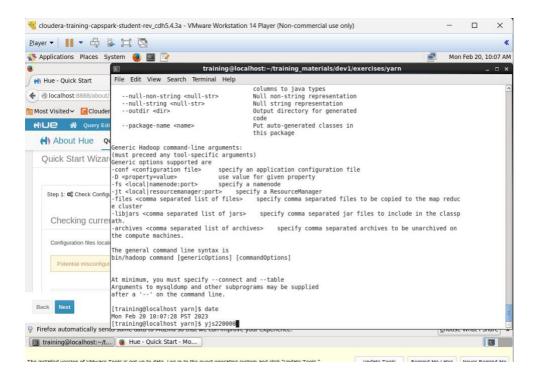
toparticles

webpage

websitehit

3. Run the sqoop import command to see its options:

\$ sqoop import --help

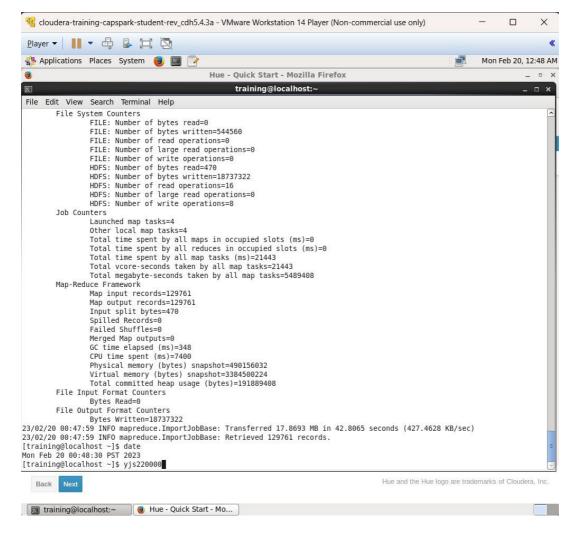


4.Use Sqoop to import the *accounts* table in the *loudacre* database and save it in HDFS under */loudacre*:

\$ sqoop import \

- --connect jdbc:mysql://localhost/loudacre \
- --username training --password training \
- --table accounts \
- --target-dir /loudacre/accounts \
- --null-non-string '\\N'

Note: The *--null-non-string* option tells Sqoop to represent null values as \N, which makes the imported data compatible with Hive and Impala.



5.2 View the Imported Data

Sqoop imports the contents of the specified tables to HDFS. You can use the hdfs command line or the Hue File Browser to view the files and their contents.

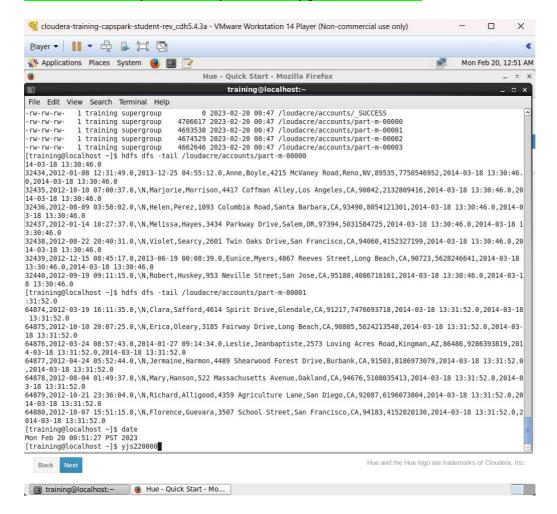
1. List the contents of the *accounts* directory

\$ hdfs dfs -ls /loudacre/accounts

```
Bytes Written=18737322
23/02/20 00:47:59 INFO mapreduce.ImportJobBase: Transferred 17.8693 MB in 42.8065 seconds (427.4628 KB/sec)
23/02/20 00:47:59 INFO mapreduce.ImportJobBase: Retrieved 129761 records.
[training@localhost ~]$ date
Mon Feb 20 00:48:30 PST 2023
[training@localhost ~]$ hdfs dfs -ls /loudacre/accounts
Found 5 items
                                                 0 2023-02-20 00:47 /loudacre/accounts/ SUCCESS
              1 training supergroup
- rw- rw- rw-
              1 training supergroup
                                          4706617 2023-02-20 00:47 /loudacre/accounts/part-m-00000
                                          4693530 2023-02-20 00:47 /loudacre/accounts/part-m-00001
4674529 2023-02-20 00:47 /loudacre/accounts/part-m-00002
- rw- rw- rw-
              1 training supergroup
             1 training supergroup
1 training supergroup
- rw-rw-rw-
                                           4662646 2023-02-20 00:47 /loudacre/accounts/part-m-00003
[training@localhost ~]$ yjs220000
```

2. Use either the Hue File Browser or the HDFS *tail* command to view the last part of the file for each of the MapReduce partition files, e.g.:

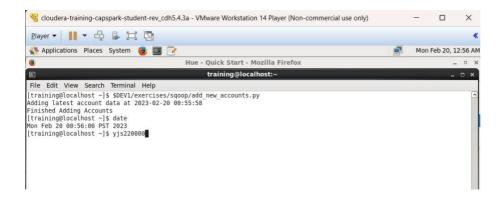
- \$ hdfs dfs -tail /loudcare/accounts/part-m-00000
- \$ hdfs dfs -tail /loudcare/accounts/part-m-00001
- \$ hdfs dfs -tail /loudcare/accounts/part-m-00002
- \$ hdfs dfs -tail /loudcare/accounts/part-m-00003



- I attached the screenshot of last command that is of part 3 file, and I see the greatest account number being 129761.
- 5.3 Import Incremental Updates to Accounts

As Loudacre adds new accounts, the account data in HDFS must be updated as accounts are created. You can use Sqoop to append these new records.

- 1. Run the add_new_accounts.py script to add the latest accounts to MySQL.
- \$\$DEV1/exercises/sqoop/add_new_accounts.py



2. Incrementally import and append the newly added accounts to the accounts directory. Use Sqoop to import on the last value on the acct_num column largest account ID:

\$ sqoop import \

- --connect jdbc:mysql://localhost/loudacre \
- --username training
- --password training \
- --incremental append \
- --null-non-string '\\N' \
- --table accounts \
- --target-dir /loudacre/accounts \
- --check-column acct_num \
- --last-value 129761

```
File Edit View Search Terminal Help

File System Counters

FILE: Number of bytes read=0
FILE: Number of bytes written=543672
FILE: Number of pread operations=0
FILE: Number of read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes written=3203
HDFS: Number of large read operations=0
HDFS: Number of large read operations=0
HDFS: Number of write operations=8
Job Counters

Launched map tasks=4
Total time spent by all maps in occupied slots (ms)=0
Total time spent by all maps in occupied slots (ms)=0
Total time spent by all map tasks (ms)=14091
Total vore-seconds taken by all map tasks=14091
Total vore-seconds taken by all map tasks=3607296
Map-Reduce Framework
Map input records=53
Map output records=53
Input split bytes=480
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
Go time elapsed (ms)=233
CPU time spent (ms)=2030
Physical memory (bytes) snapshot=481181696
Virtual memory (bytes) snapshot=3378053120
Total committed heap usage (bytes)=191889408
File Input Format Counters
Bytes Written=3203
23/02/20 e1:15:19 INFO mapreduce.ImportJobBase: Transferred 3.1279 KB in 31.0022 seconds (103.3154 bytes/sec)
23/02/20 e1:15:19 INFO mapreduce.ImportJobBase: Retrieved 53 records.

[training@localhost -]$ date
More About the spent of bytes yielded and the specific of the spent of t
```

3. List the contents of the *accounts* directory to verify the Sqoop import:

\$ hdfs dfs -ls /loudacre/accounts

```
[training@localhost ~]$ hdfs dfs -ls /loudacre/webpage
Found 5 items
- rw- rw- rw-
               1 training supergroup
                                                      0 2023-02-20 01:15 /loudacre/webpage/ SUCCESS
                                                   751 2023-02-20 01:15 /loudacre/webpage/part-m-00000
874 2023-02-20 01:15 /loudacre/webpage/part-m-00001
804 2023-02-20 01:15 /loudacre/webpage/part-m-00002
               1 training supergroup
- rw- rw- rw-
               1 training supergroup
              1 training supergroup
- rw- rw- rw-
              1 training supergroup
                                                   774 2023-02-20 01:15 /loudacre/webpage/part-m-00003
[training@localhost ~]$ date
Mon Feb 20 01:17:23 PST 2023
[training@localhost ~]$ yjs220000
```

4. You should see three new files. Use Hadoop's *cat* command to view the entire contents of these files.

\$ hdfs dfs -cat /loudacre/accounts/part-m-0000[456]

Note: I tried running this command line, but several times I was shown No such file or directory.

- 5.4 Import webpage data using an alternate find delimiter
- 1. We also want to import the webpage table to HDFS. But first look at a few records in that table using the sqoop *eval* command.

\$ sqoop eval \

--query "SELECT * FROM webpage LIMIT 10" \

--connect jdbc:mysql://localhost/loudacre \

--username training --password training

- We notice that the values in the last column contain commas. By default, sqoop uses commas as field separators, but because the data itself uses commas, we can't do that this time.
- 2. Use sqoop to import the webpage table, but use the tab character (\t) instead of the default (comma) as the field terminator.

\$ sqoop import \

- --connect jdbc:mysql://localhost/loudacre \
- --username training --password training ackslash
- --table webpage \
- --target-dir /loudacre/webpage \
- --fields-terminated-by "\t"

```
File Edit View Search Terminal Help

File System Counters

FILE: Number of bytes read=0
FILE: Number of bytes written=543672
FILE: Number of pread operations=0
FILE: Number of read operations=0
FILE: Number of write operations=0
FILE: Number of write operations=0
HDFS: Number of bytes written=3203
HDFS: Number of bytes read=480
HDFS: Number of large read operations=0
HDFS: Number of large read operations=0
HDFS: Number of write operations=8
Job Counters

Launched map tasks=4
Other local map tasks=4
Total time spent by all maps in occupied slots (ms)=0
Total time spent by all map tasks (ms)=14091
Total tore-seconds taken by all map tasks=14091
Total tore-seconds taken by all map tasks=14091
Total magabyte-seconds taken by all map tasks=3607296
Map-Reduce Framework
Map input records=53
Map output records=53
Input split bytes=480
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=233
CPU time spent (ms)=2939
Physical memory (bytes) snapshot=3378053120
Total committed heap usage (bytes)=191889408
File Input Format Counters
Bytes Written=3203
23/02/20 el:15:19 INFO mapreduce.ImportJobBase: Transferred 3.1279 KB in 31.0022 seconds (103.3154 bytes/sec)
23/02/20 el:15:19 INFO mapreduce.ImportJobBase: Transferred 53 records.

Itraining@localhost -] 5 date
Mon Feb 20 el:15:46 PST 2023
Itraining@localhost -] 5 date
More Feb 20 el:15:46 PST 2023
Itraining@localhost -] 5 date
More Feb 20 el:15:46 PST 2023
Itraining@localhost -] 5 yiz220000
```

3. Using Hue or the hdfs command line, review the data files imported to the */loudacre/webpage* directory. Take note of the structure of the data; you will use this data in the next exercises.