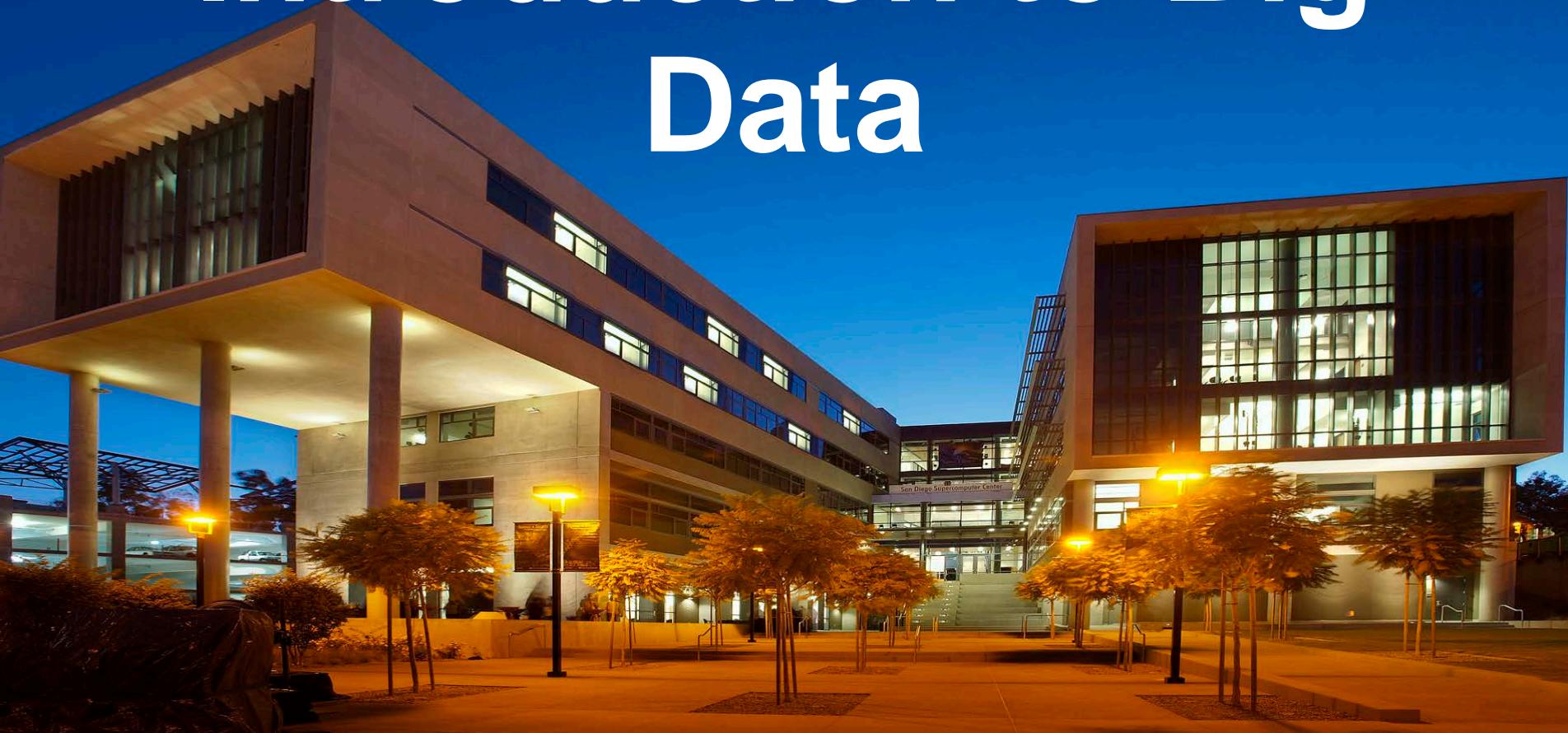
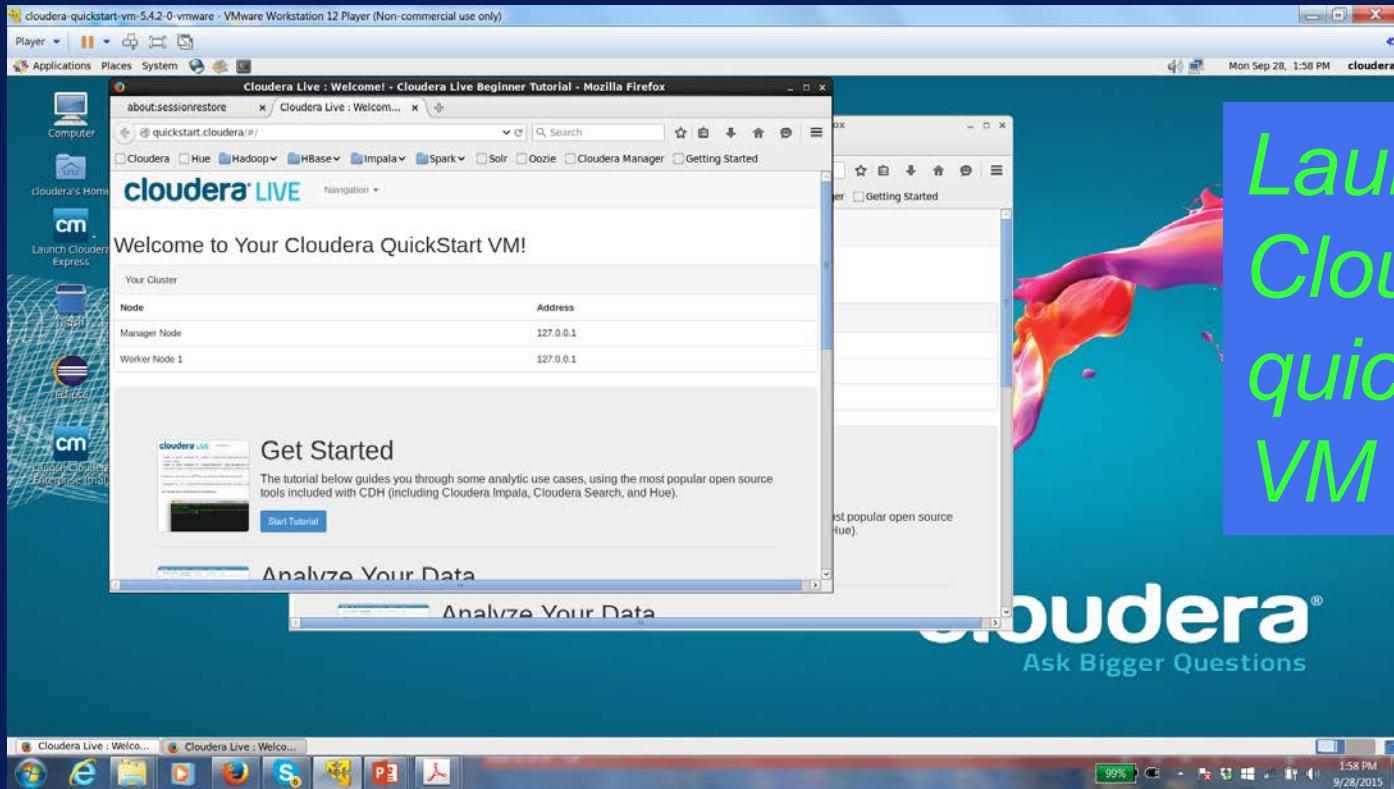


# Introduction to Big Data



# Lecture #1

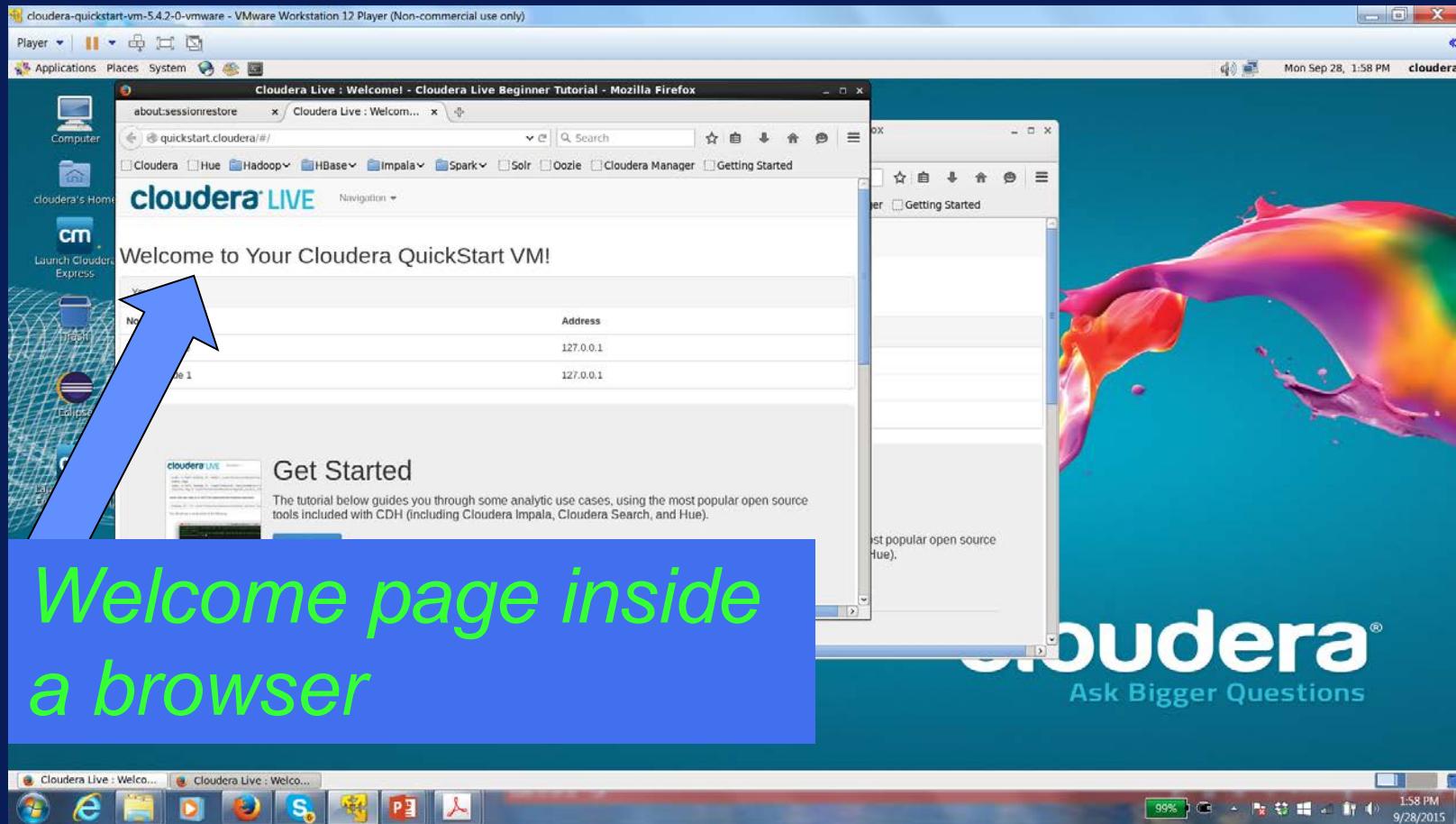
# Cloudera VM Tour

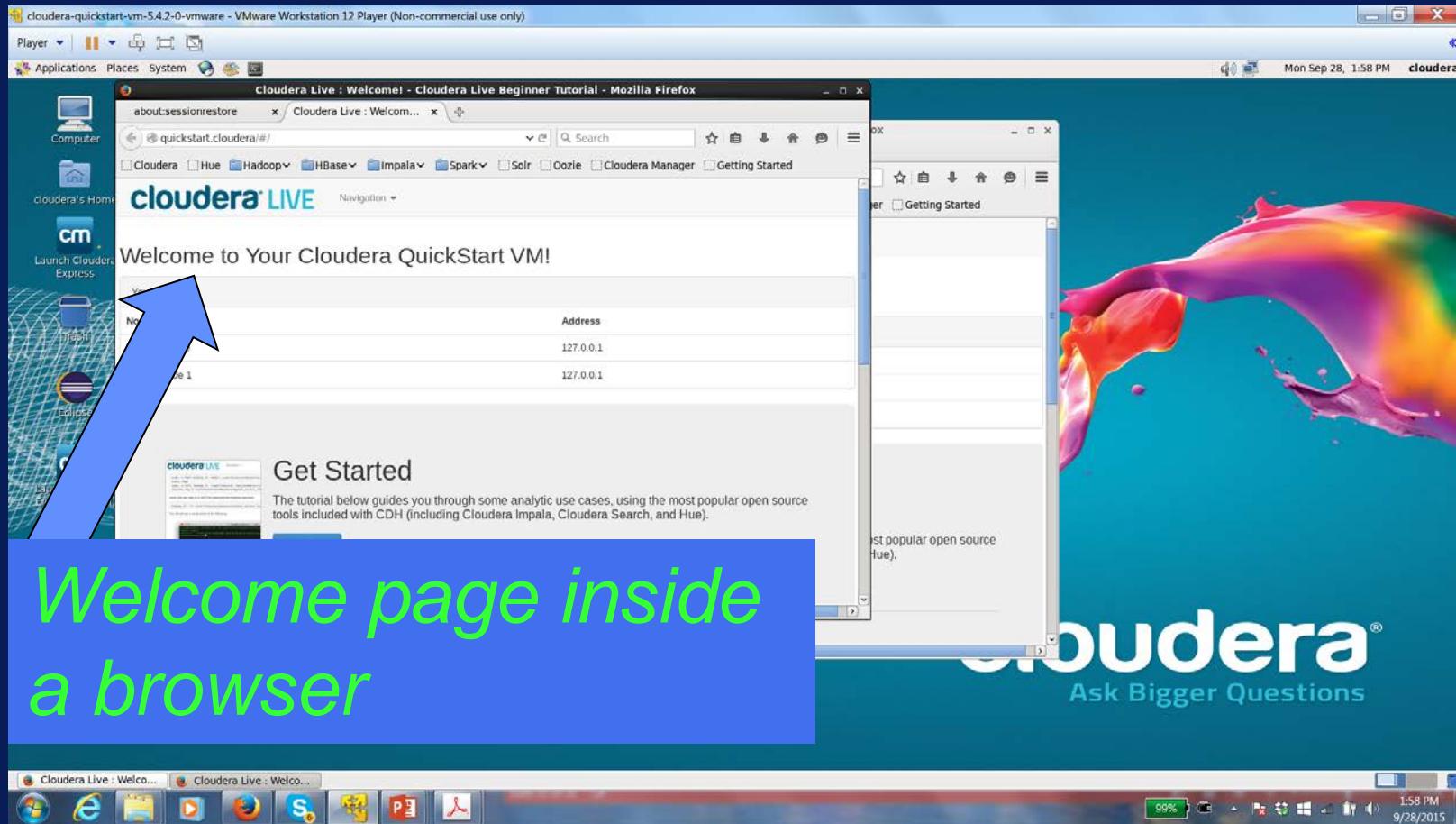


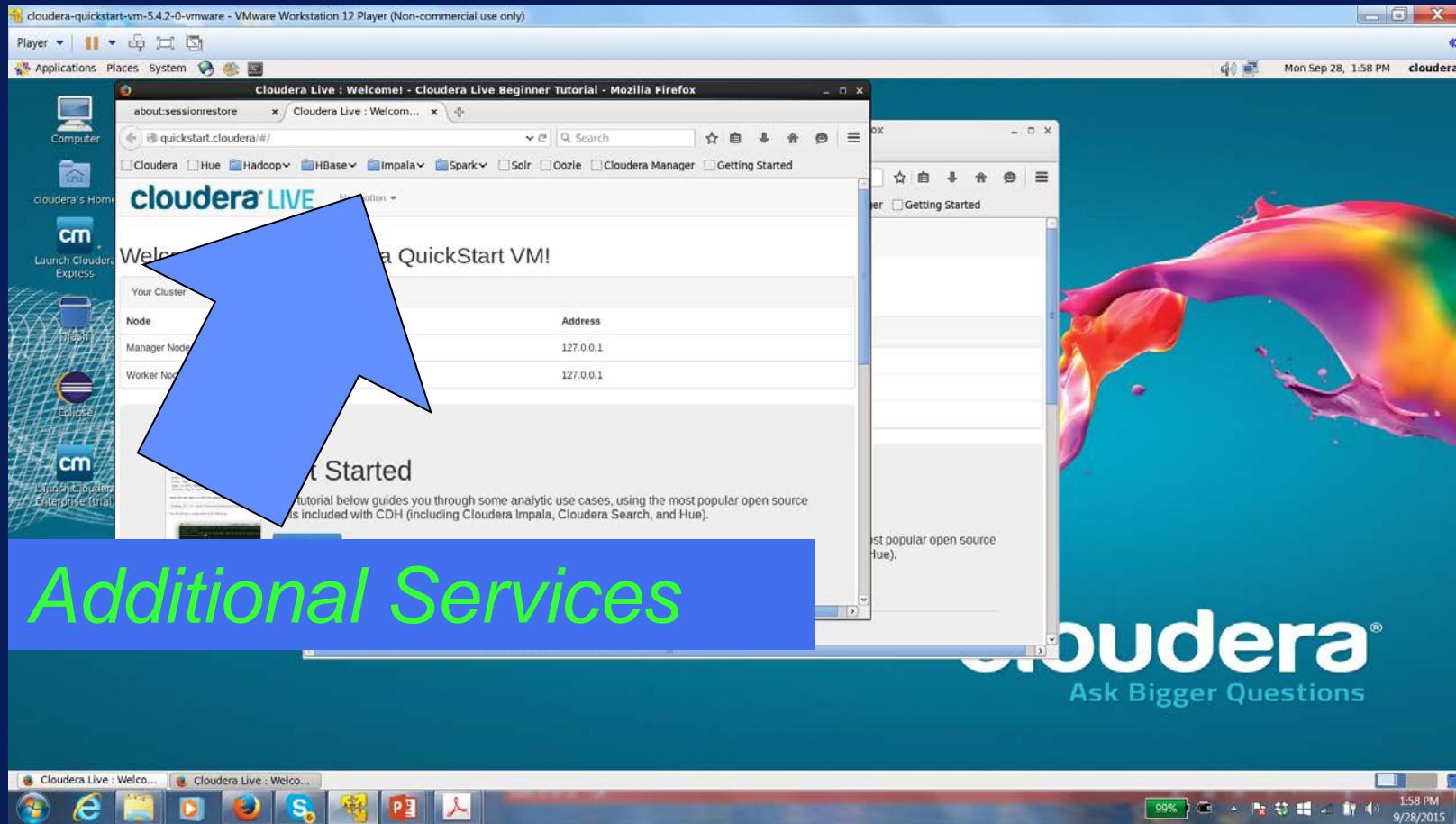
Launch the  
Cloudera  
quick start  
VM

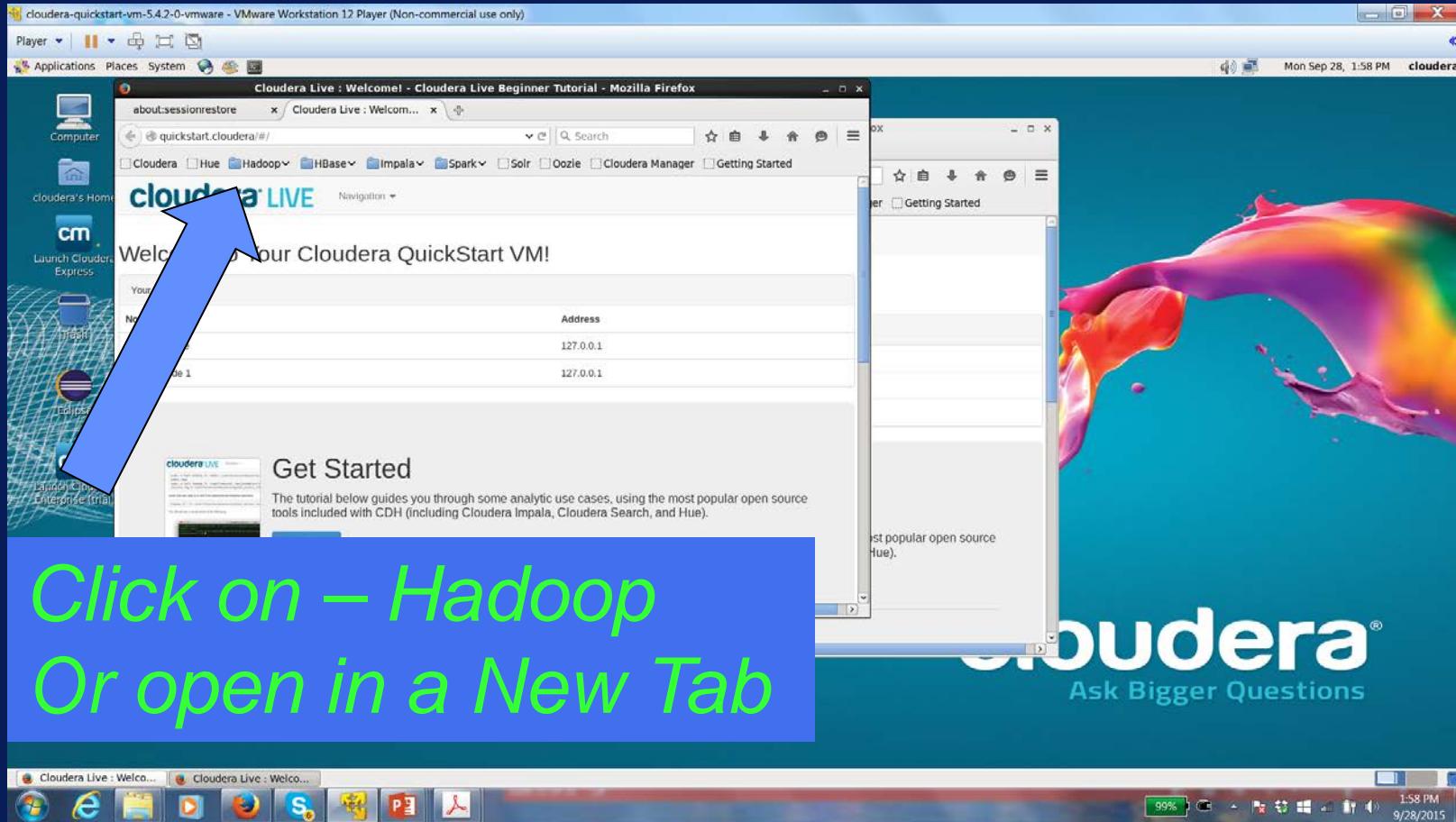
cloudera®  
Ask Bigger Questions

# Inside the VM









cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

Namenode Information - Mozilla Firefox

Cloudera Live : Welcome... | Namenode Information

quickstart.cloudera:5070/dfshealth.html#tab-overview

Cloudera | Hue | Hadoop | HBase | Impala | Spark | Solr | Oozie | Cloudera Manager | Getting Started

Hadoop Overview Datanodes Snapshot Startup Progress Utilities

Overview 'quickstart.cloudera:8020' (active)

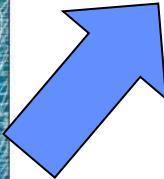
Started: Sun Sep 27 21:35:10 PDT 2015  
Version: 2.6.0-cdh5.4.2, r15b703c8725733b7b2813d2325659eb7d57e7a3f  
Compiled: 2015-05-20T00:03Z by jenkins from Unknown  
Cluster ID: CID-9b81ad2f-4ba7-43cb-95d7-e2390594e63c  
Block Pool ID: BP-286282631-127.0.0.1-1433865208026

Summary

Security is off.  
Safemode is off.  
458 files and directories, 384 blocks = 842 total filesystem object(s).  
Heap Memory used 238.7 MB of 480 MB Heap Memory. Max Heap Memory is 889 MB.  
Non Heap Memory used 39.91 MB of 40.5 MB Committed Non Heap Memory. Max Non Heap Memory is 130 MB.

Configured Capacity:	54.64 GB
DFS Used:	376.38 MB
Non DFS Used:	9.44 GB

Cloudera Live : Welcome... | Namenode Information... | 100% | 2:13 PM | 9/28/2015



*Overview of the Hadoop installation within the VM*

cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

Namenode Information - Mozilla Firefox

Cloudera Live : Welcome... | Namenode Information

quickstart.cloudera:50070/dfshealth.html#tab-overview

Cloudera Hue Hadoop Impala Spark Solr Oozie Cloudera Manager Getting Started

Hadoop Overview Datanodes Snapshot Startup Progress Utilities

Overall Status: 'start.cloudera:8020' (active)

Start Date: Sun Sep 27 21:35:10 PDT 2015

Version: 2.6.0-cdh5.4.2, r15b703c8725733b7b2813d2325659eb7d57e7a3f

Compiled: 2015-05-20T00:03Z by jenkins from Unknown

Cluster ID: CID-9b81ad2f-4ba7-43cb-95d7-e2390594e63c

Block Pool ID: BP-286282631-127.0.0.1-1433865208026

## Summary

Security is off.

Safemode is off.

458 files and directories, 384 blocks = 842 total filesystem object(s).

Heap Memory used 238.7 MB of 480 MB Heap Memory. Max Heap Memory is 889 MB.

Non Heap Memory used 39.91 MB of 40.5 MB Committed Non Heap Memory. Max Non Heap Memory is 130 MB.

Configured Capacity:	54.64 GB
DFS Used:	376.38 MB
Non DFS Used:	9.44 GB

Cloudera Live : Welcome... | Namenode Information

2:13 PM 9/28/2015



cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

Namenode information - Mozilla Firefox

Cloudera Live : Welcome | Namenode information

quickstart.cloudera:50070/dfshealth.html#tab-datanode

Cloudera | Hue | Hadoop | HBase | Impala | Spark | Solr | Oozie | Cloudera Manager | Getting Started

Hadoop Overview Datanodes Snapshot Startup Progress Utilities

## Datanode Information

In operation

Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
quickstart.cloudera (127.0.0.1:50010)	1	In Service	54.64 GB	376.30 MB	9.44 GB	44.82 GB	382	376.30 MB (0.67%)	0	2.6.0- cah5.4.2

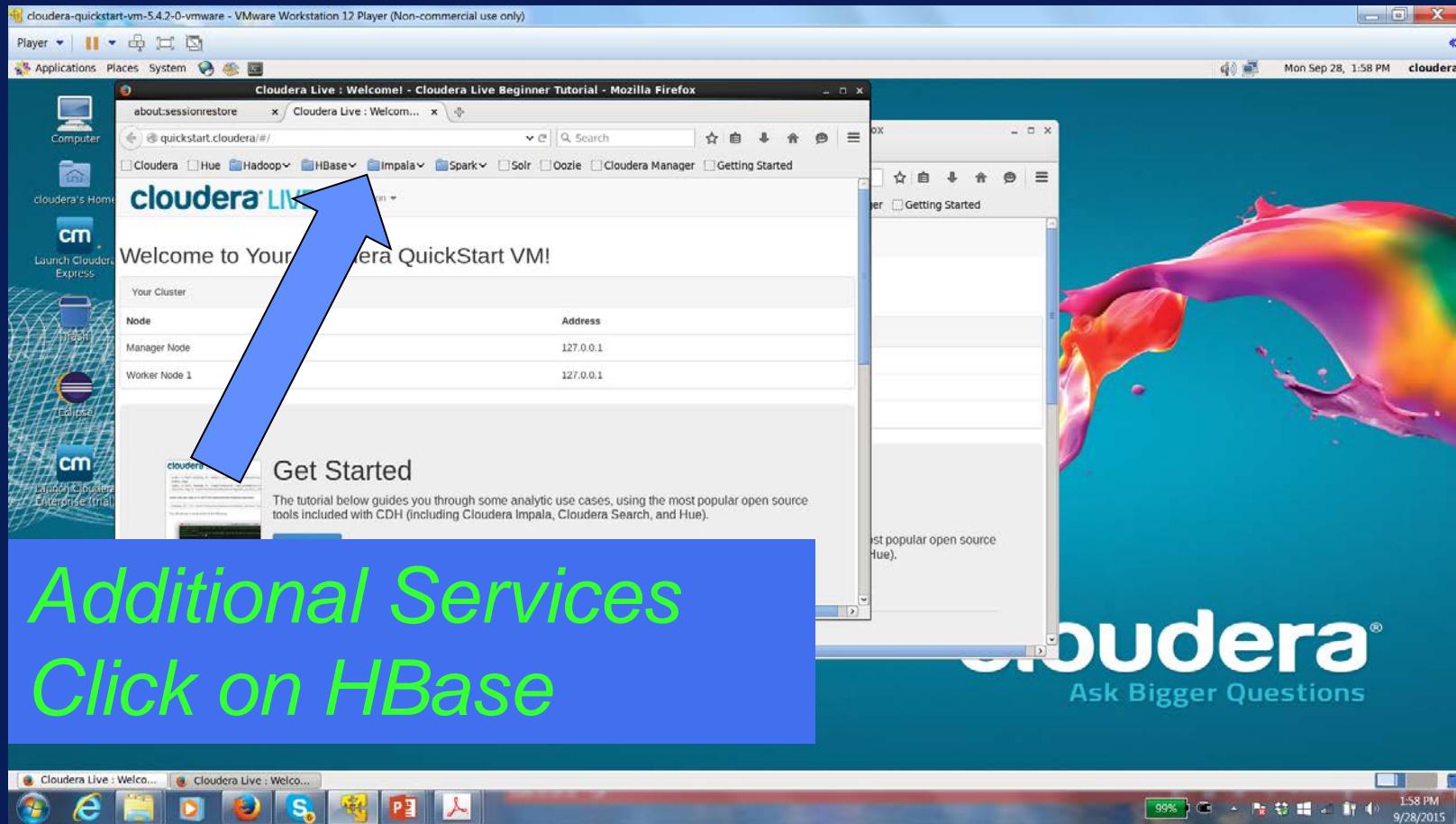
Decommissioning

Node	Last contact	Under replicated blocks	Blocks with no live replicas	Under Replicated Blocks In files under construction

Hadoop, 2014.

Legacy UI

Cloudera Live : Welcome | Namenode information



cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

HBase Region Server: quickstart.cloudera - Mozilla Firefox

Mon Sep 28, 2:29 PM cloudera

Cloudera Live : Welcome | HBase Region Server: quickstart.cloudera | Apache HBASE

quickstart.cloudera:60030/rs-status?filter=operation

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

APACHE HBASE Home Local Logs Log Level Debug Dump Metrics Dump HBase Configuration

## RegionServer quickstart.cloudera,60020,1443414987853

### Server Metrics

Base Stats Memory Requests WALS Storefiles Queues

Requests Per Second	Num. Regions	Block locality	Slow WAL Append Count
0	2	100	0

### Tasks

Show All Monitored Tasks Show non-RPC Tasks Show All RPC Handler Tasks Show Active RPC Calls Show Client Operations View as JSON

No tasks currently running on this node.

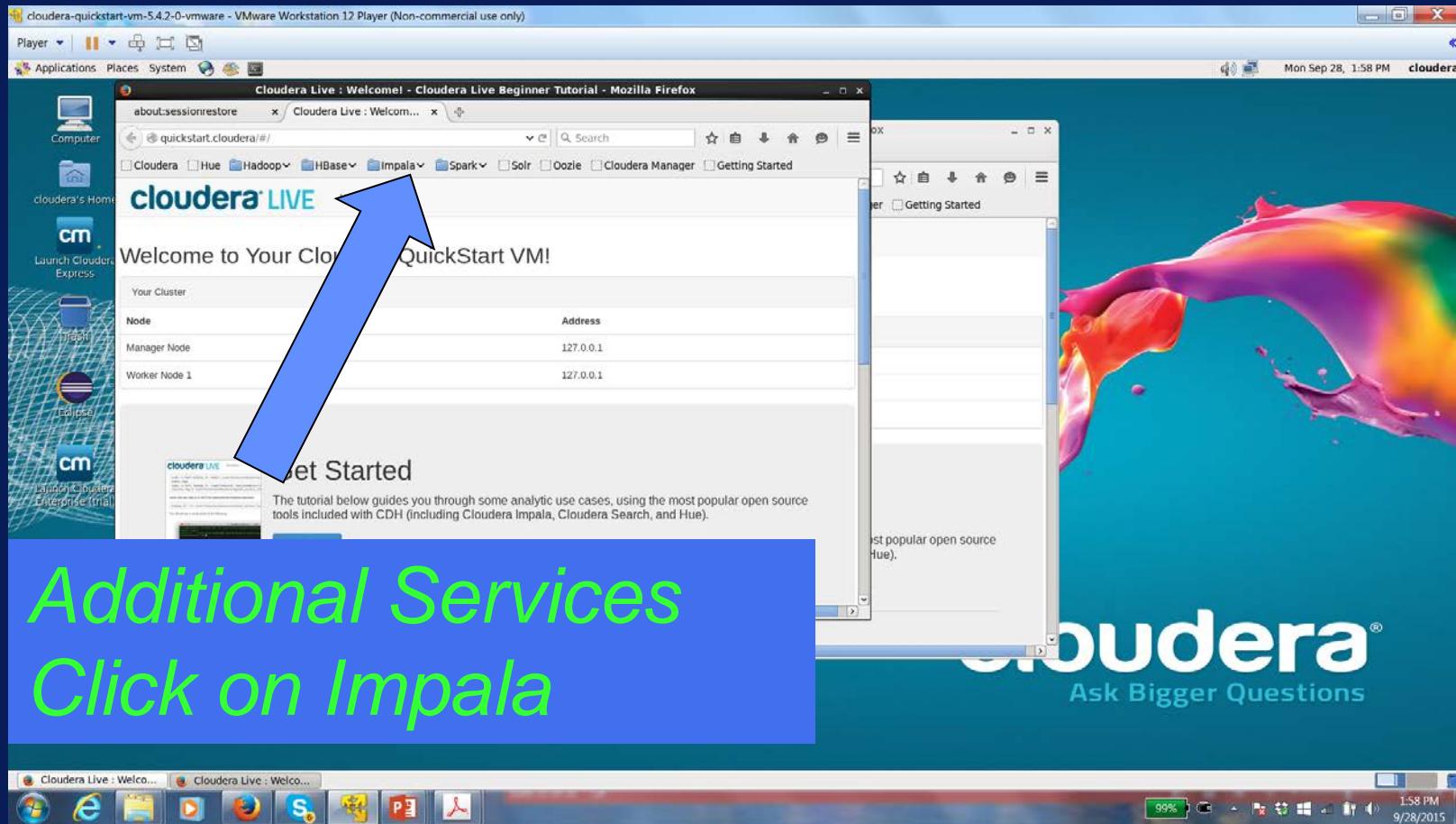
### Block Cache

Base Info Config Stats L1 L2

Attribute	Value	Description
Implementation	LruBlockCache	Block cache implementing class

Cloudera Live : Welcome | HBase Region Server: ...

100% 2:29 PM 9/28/2015



cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

Cloudera Impala - Mozilla Firefox

Cloudera Live : Welcome Cloudera Impala quickstart.cloudera:25000/queries

Cloudera Manager Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

impalad / /backends /catalog /hadoop-varz /logs /memz /metrics /queries /rpcz /sessions /threadz /varz

Queries

This page lists all running queries, plus any completed queries that are archived in memory. The size of that archive is controlled with the --query\_log\_size command line parameter.

0 queries in flight

User	Default Db	Statement	Query Type	Start Time	Scan Progress	State	Last Event	# rows fetched	Details	Action
------	------------	-----------	------------	------------	---------------	-------	------------	----------------	---------	--------

Last 25 Completed Queries

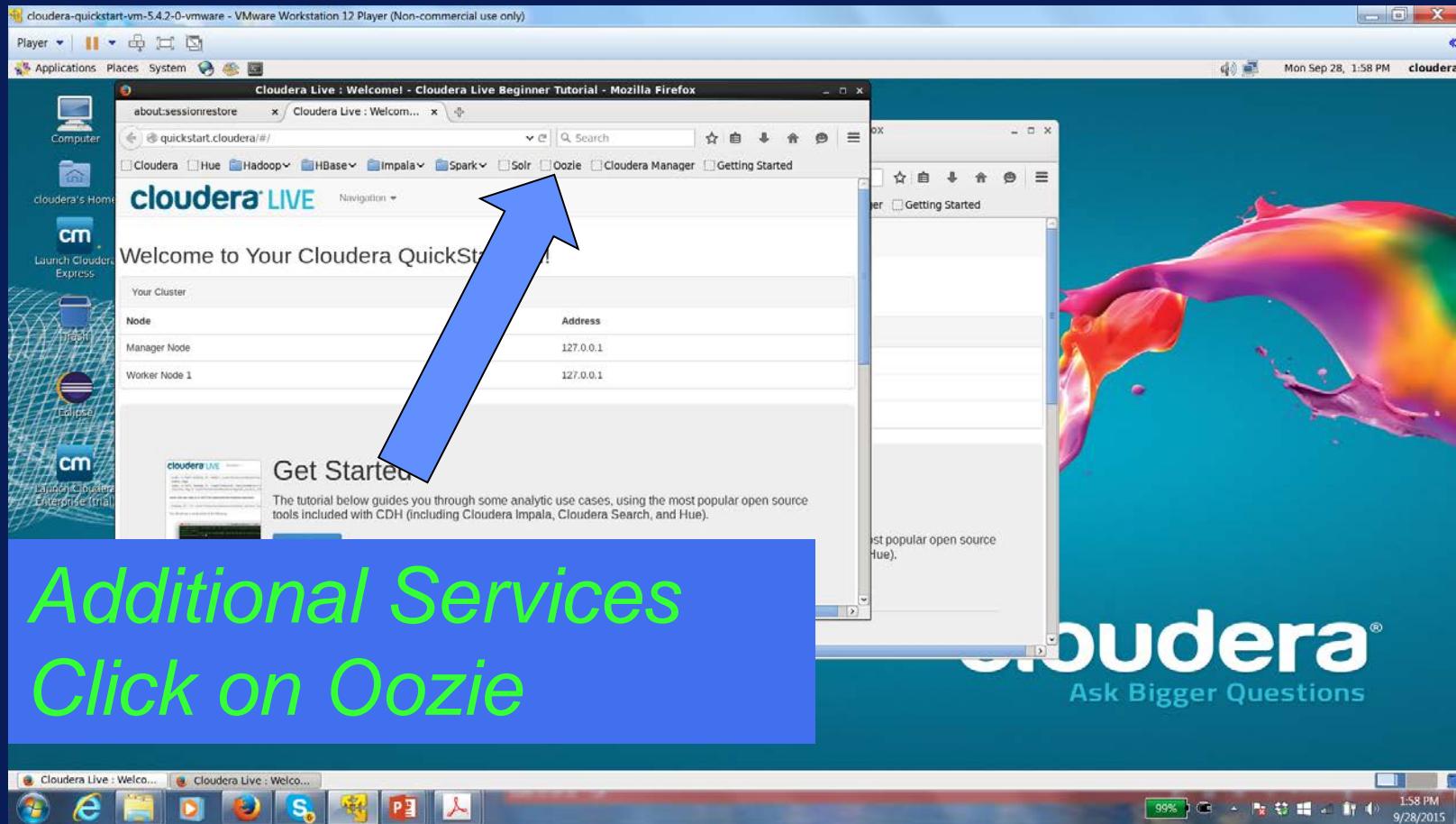
User	Default Db	Statement	Query Type	Start Time	End Time	Scan Progress	State	# rows fetched	Details
------	------------	-----------	------------	------------	----------	---------------	-------	----------------	---------

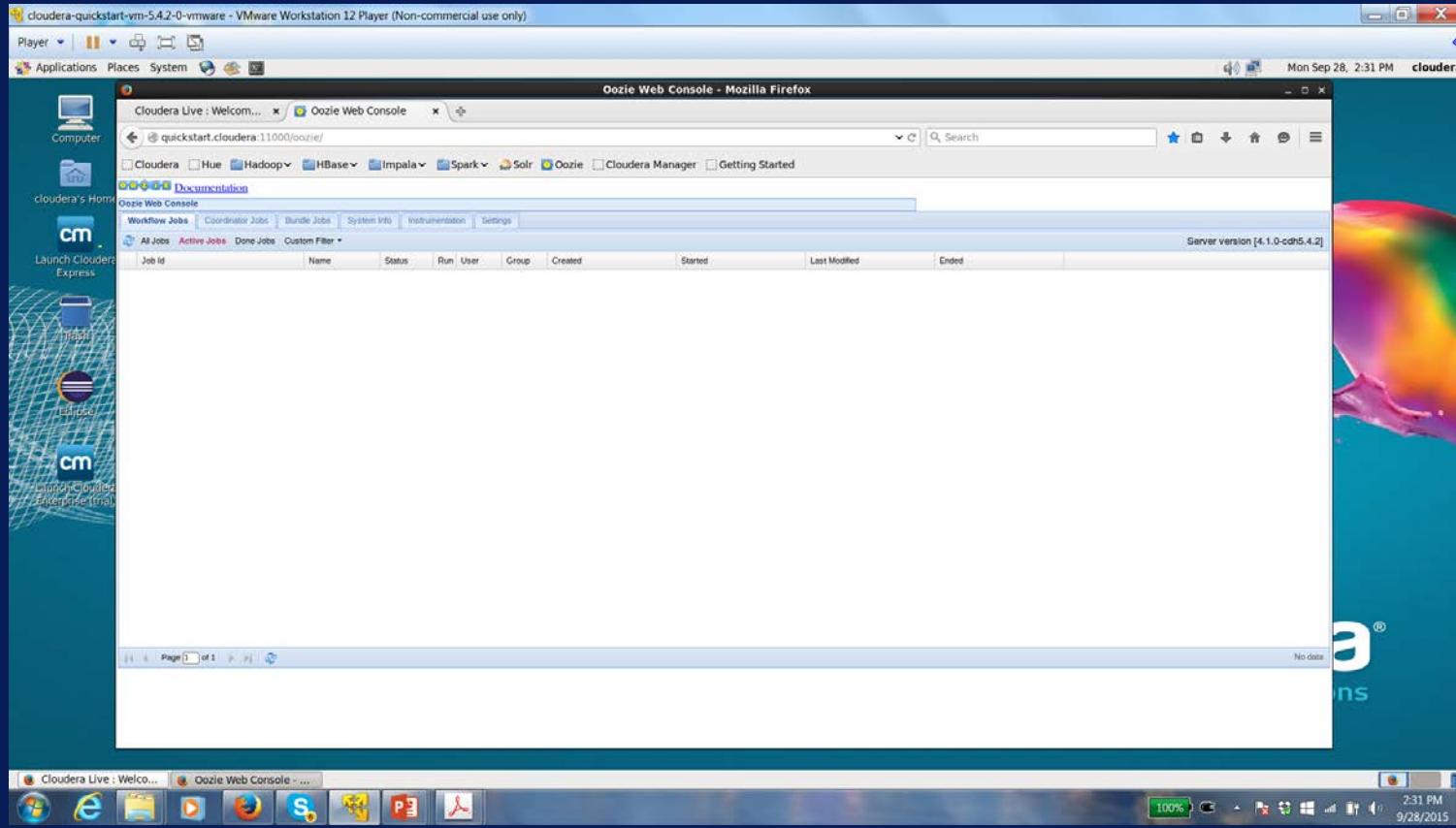
Query Locations

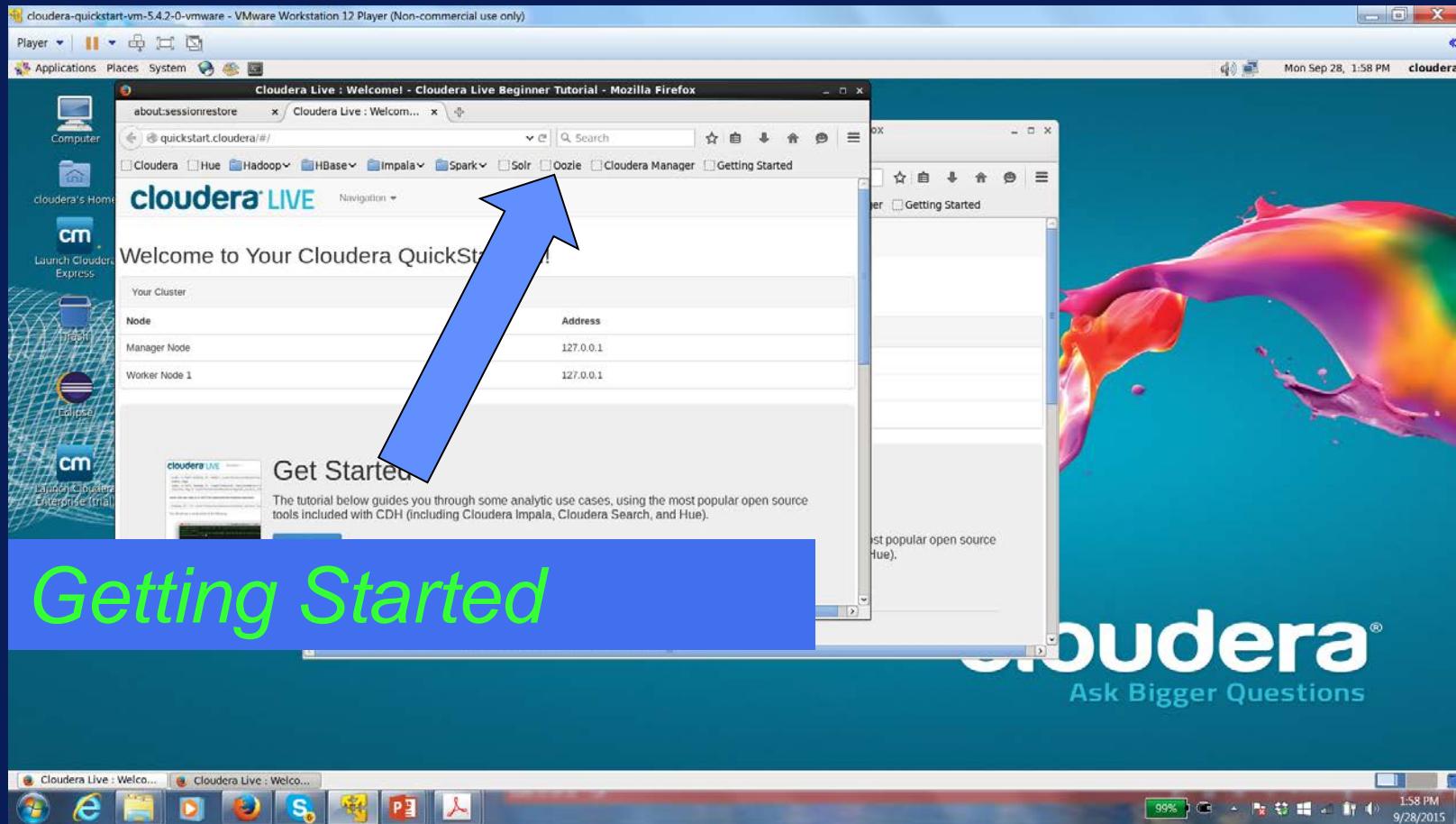
Location	Number of Fragments
----------	---------------------

Cloudera Live : Welcome Cloudera Impala - Mozilla Firefox

100% 2:31 PM 9/28/2015







cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System cloudera

Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial - Mozilla Firefox

Mon Sep 28, 2:33 PM

Computer cloudera's Home Launch Cloudera Express Manager Services cm Launch Cloudera Enterprise Trial

Cloudera Live : Welcome... Cloudera Live : Welcome...

quickstart.cloudera/#

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

cloudera LIVE Navigation

## Welcome to Your Cloudera QuickStart VM!

Your Cluster	
Node	Address
Manager Node	127.0.0.1
Worker Node 1	127.0.0.1

### Get Started

The tutorial below guides you through some analytic use cases, using the most popular open source tools included with CDH (including Cloudera Impala, Cloudera Search, and Hue).

[Start Tutorial](#)

### Analyze Your Data

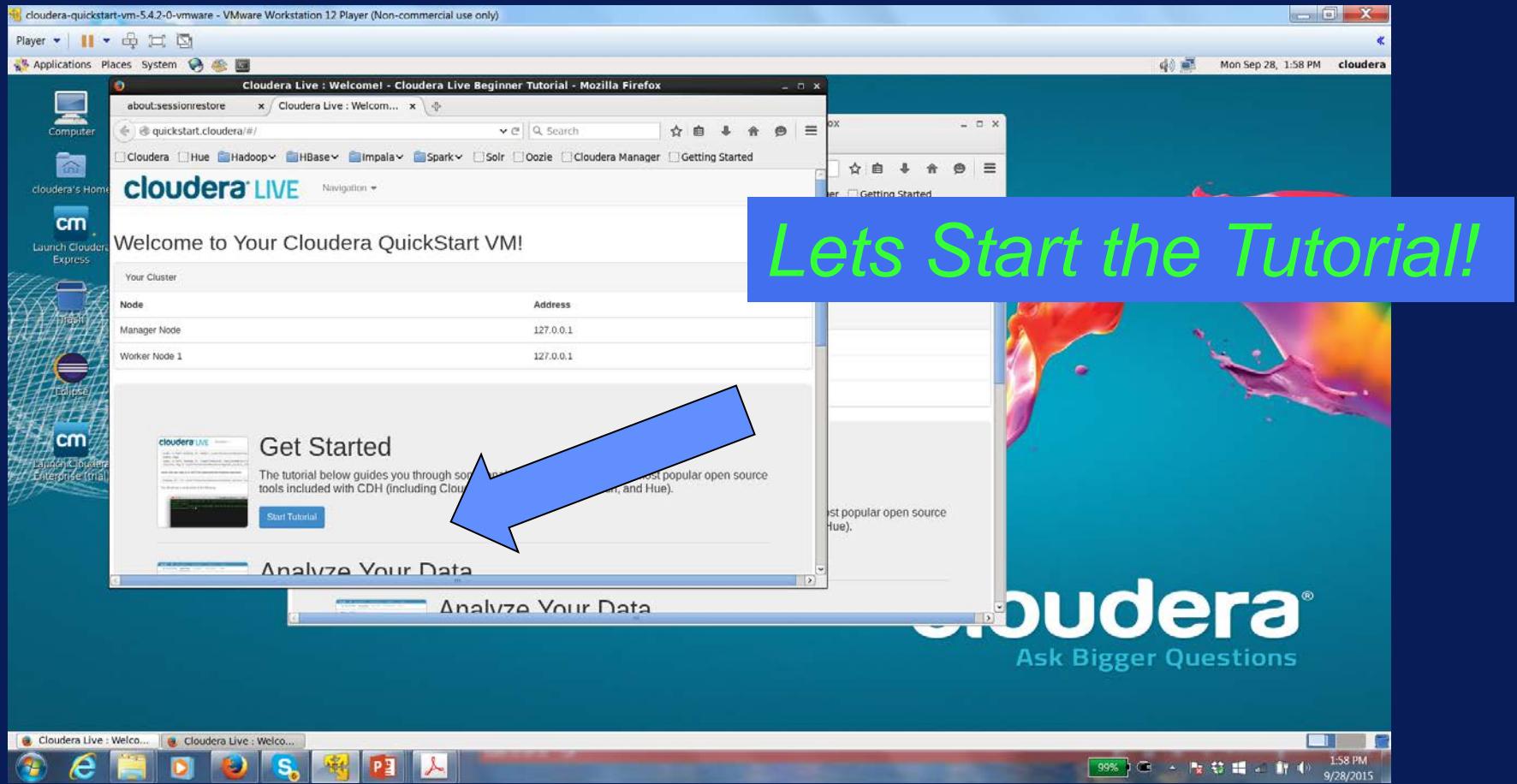
Hue is the open source web interface for Hadoop that lets you analyze your data. Simply load in your data and then easily begin to analyze, search, and visualize it. In the QuickStart VM, the administrative username for Hue is 'cloudera' and the password is 'cloudera'.

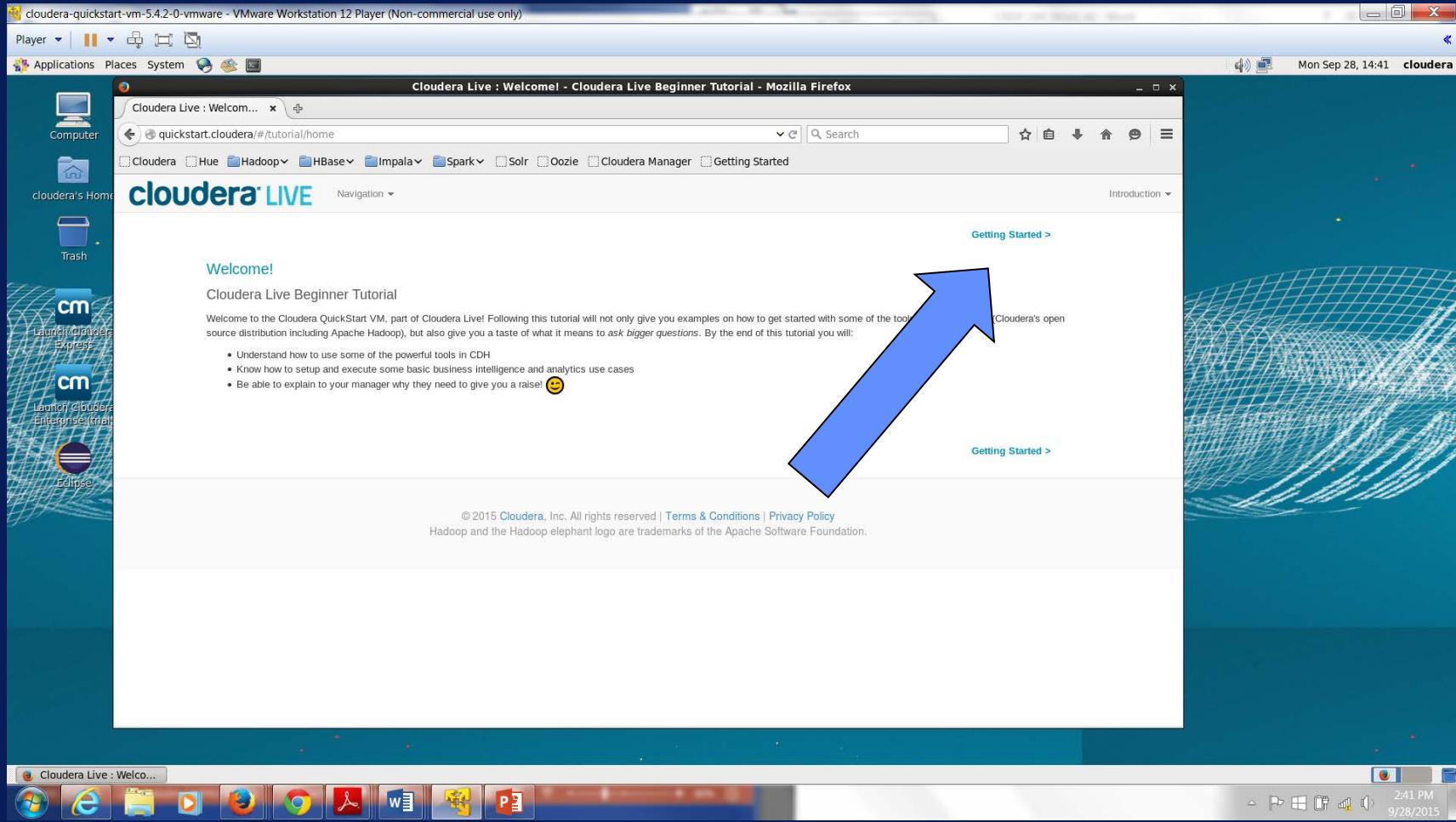
[Launch Hue UI!](#)

### Manage Your Cluster

Cloudera Live : Welcome... Cloudera Live : Welcome...

2:33 PM 9/28/2015





Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial - Mozilla Firefox

Cloudera Live : Welcome... quickstart.cloudera/#/tutorial/getting\_started

cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Getting Started

Tutorial Exercise 1 >

## Getting Started

### Define a Business Question

For the remainder of this tutorial, we will present examples in the context of a made up corporation called DataCo, and our mission is to help the organization get better insight by asking bigger questions.

### Scenario:

Your Management: is talking euphorically about *Big Data*...

You: are carefully skeptical, as it will most likely all land on your desk anyway. Alternatively it has already landed on you, with the nice project description of: Go figure this Hadoop thing out...

< Introduction Tutorial Exercise 1 >

### Good to Know

Any successful PoC needs to address something your organization cares about. Hence, the first thing you need to do is to: **define a business question**.

It won't just impress your manager that you *think big* and have perspective on the business needs of your organisation (which in English means you just helped your manager to look good in front of his management). It will also help you to go through a well scoped PoC and get the investments you need to be successful.

Without a well defined question, you won't know how to properly model your data, i.e. what structure to apply at query time, or what data sets and tools to use to best serve the use case.

cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player | Applications Places System cloudera

Mon Sep 28, 14:44 cloudera

Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial - Mozilla Firefox

Cloudera Live : Welcome... quickstart.cloudera/#/tutorial/ingest\_structured\_data

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Navigation Tutorial Exercise 1 Tutorial Exercise 2 >

## Tutorial Exercise 1

### Ingest Structured Data

In this scenario, DataCo's business question is: *What products do our customers like to buy?* To answer this question, the first thought might be to look at the transaction data, which should indicate what customers actually do buy and like to buy, right?

This is probably something you can do in your regular RDBMS environment, but a benefit with Cloudera's platform is that you can do it at greater scale at lower cost, on the same system that you may also use for many other types of analysis.

What this exercise demonstrates is how to do exactly the same thing you already know how to do, but in CDH. Seamless integration is important when evaluating any new infrastructure. Hence, it's important to be able to do what you normally do, and not break any regular BI reports or workloads over the dataset you plan to migrate.

```
graph TD; departments[departments] --- categories[categories]; departments --- products[products]; categories --- products; products --- order_items[order_items]; products --- customers[customers]; order_items --- customers
```

**About Sqoop:**

Apache Sqoop is a tool that uses MapReduce to transfer data between Hadoop clusters and relational databases very efficiently. It works by spawning tasks on multiple data nodes to download various portions of the data in parallel. When you're finished, each piece of data is replicated to ensure reliability, and spread out across the cluster to ensure you can process it in parallel on your cluster.

There are 2 versions of Sqoop included in Cloudera's platform. Sqoop 1 is a "thick client" and is what you use in this tutorial. The command you run will directly submit the MapReduce jobs to transfer the data. Sqoop 2 consists of a central server that submits the MapReduce jobs on behalf of clients, and a much lighter weight client that you use to connect to the server. The "Sqoop" you see in Cloudera Manager is the Sqoop 2 server, although Cloudera

Cloudera Live : Welcome...

2:44 PM 9/28/2015

cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player | Applications Places System cloudera

Mon Sep 28, 14:44 cloudera

Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial - Mozilla Firefox

Cloudera Live : Welcom... quickstart.cloudera/#/tutorial/ingest\_structured\_data

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

Navigation Tutorial Exercise 1 > Tutorial Exercise 2 >

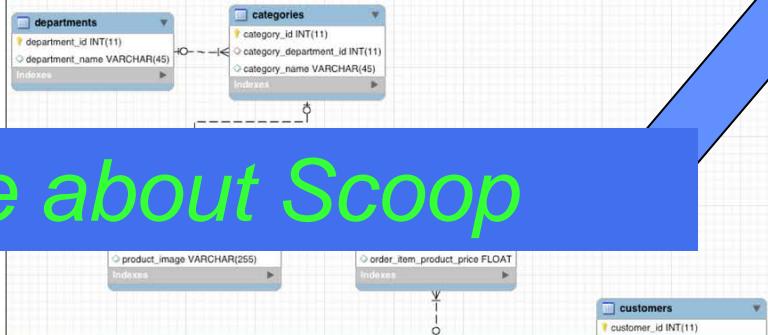
**Tutorial Exercise 1**

Ingest Structured Data

In this scenario, DataCo's business question is: *What products do our customers like to buy?* To answer this question, the first thought might be to look at the transaction data, which should indicate what customers actually do buy and like to buy, right?

This is probably something you can do in your regular RDBMS environment, but a benefit with Cloudera's platform is that you can do it at greater scale at lower cost, on the same system that you may also use for many other types of analysis.

What this exercise demonstrates is how to do exactly the same thing you already know how to do, but in CDH. Seamless integration is important when evaluating any new infrastructure. Hence, it's important to be able to do what you normally do, and not break any regular BI reports or workloads over dataset you plan to migrate.



**About Sqoop:**

Apache Sqoop is a tool that uses MapReduce to transfer data between Hadoop clusters and relational databases very efficiently. It works by spawning tasks on multiple data nodes to download various portions of the data in parallel. When you're finished, each piece of data is replicated to ensure reliability, and spread out across the cluster to ensure you can process it in parallel on your cluster.

There are 2 versions of Sqoop included in Cloudera's platform. Sqoop 1 is a "thick client" and is what you use in this tutorial. The command you run will directly submit the MapReduce jobs to transfer the data. Sqoop 2 consists of a central server that submits the MapReduce jobs on behalf of clients, and a much lighter weight client that you use to connect to the server. The "Sqoop" you see in Cloudera Manager is the Sqoop 2 server, although Cloudera

More about Scoop

2:44 PM 9/28/2015

# Scroll down for more

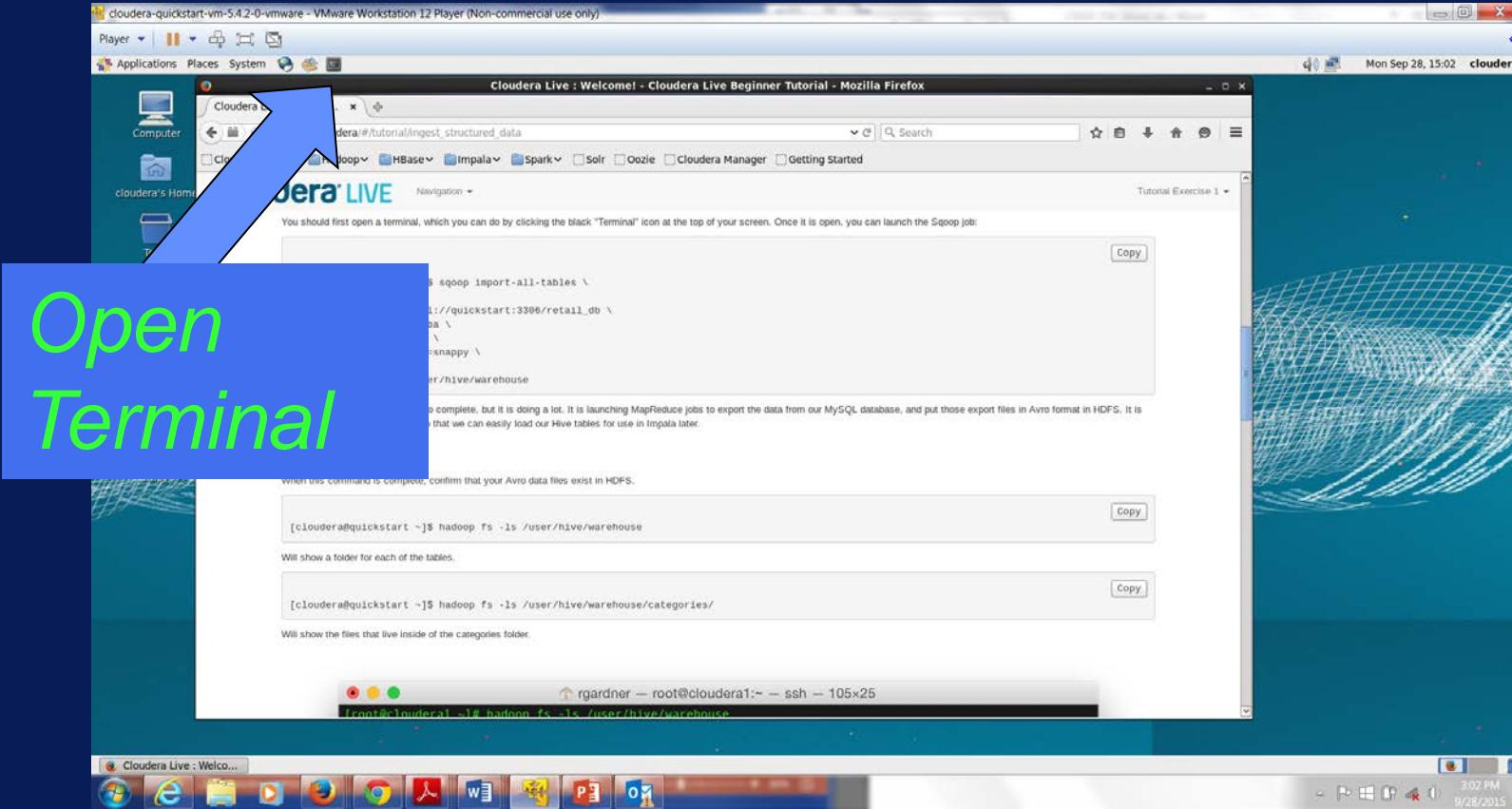
The screenshot shows a desktop environment with a VMware Workstation window titled "cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)". Inside the window, a Firefox browser displays the "Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial". The main content is a "Navigation" section with a "Tutorial Index" and a "Data Flow Diagram". The diagram illustrates a relational database schema with tables: departments, categories, products, order\_items, and customers. Each table has its primary key highlighted in yellow. The "order\_items" table is shown with a detailed view of its columns: order\_item\_id, order\_item\_order\_id, order\_item\_product\_id, order\_item\_quantity, order\_item\_subtotal, and order\_item\_product\_price. To the right of the diagram, a large blue arrow points downwards, with the text "Scroll down for more" overlaid in green. Below the diagram, there is explanatory text about using Sqoop to ingest MySQL data into HDFS.

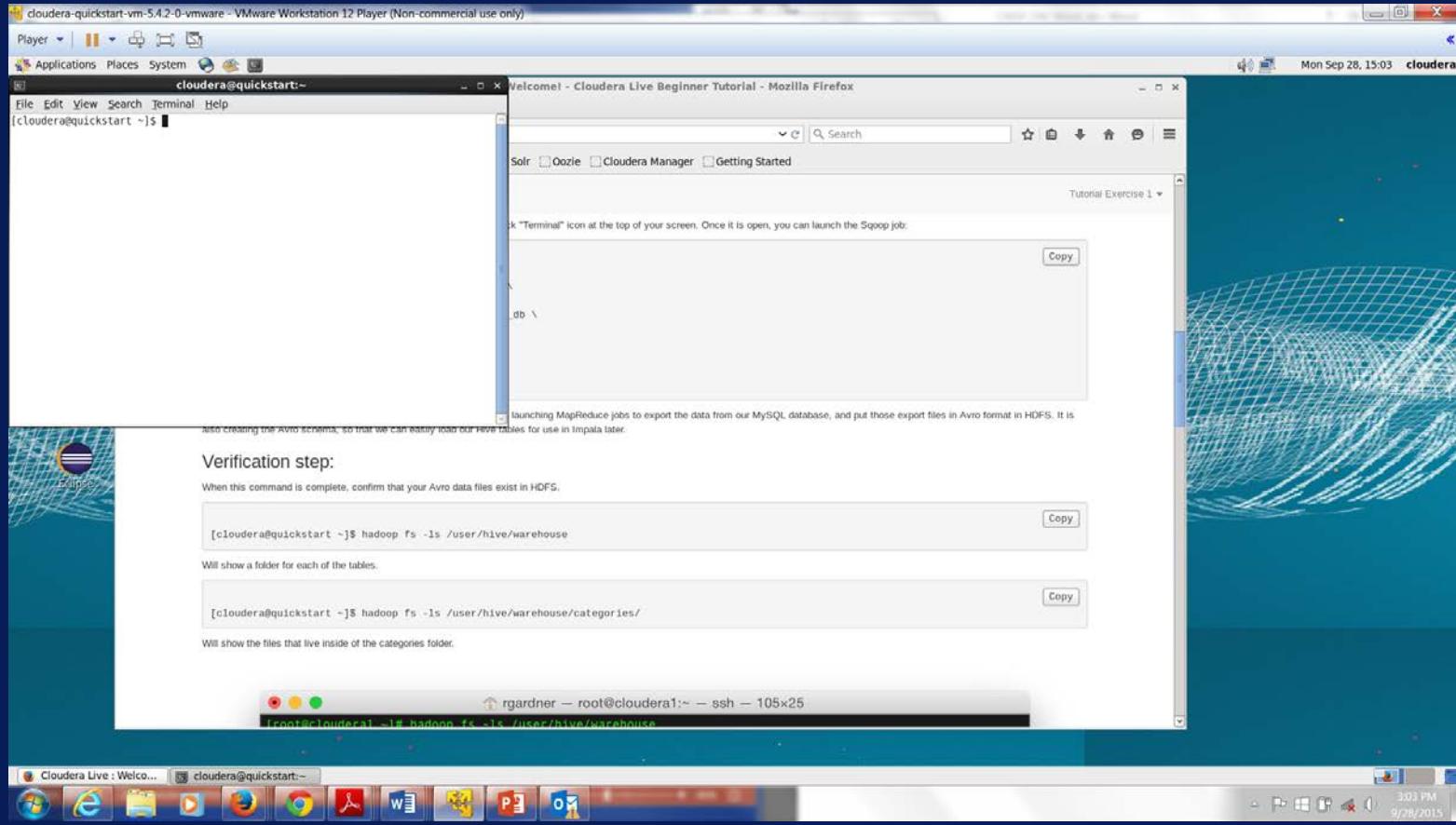
To analyze the transaction data in the new platform, we need to ingest it into the Hadoop Distributed File System (HDFS). We need to find a tool that easily transfers structured data from a RDBMS to HDFS, while preserving structure. That enables us to query the data, but not interfere with or break any regular workload on it.

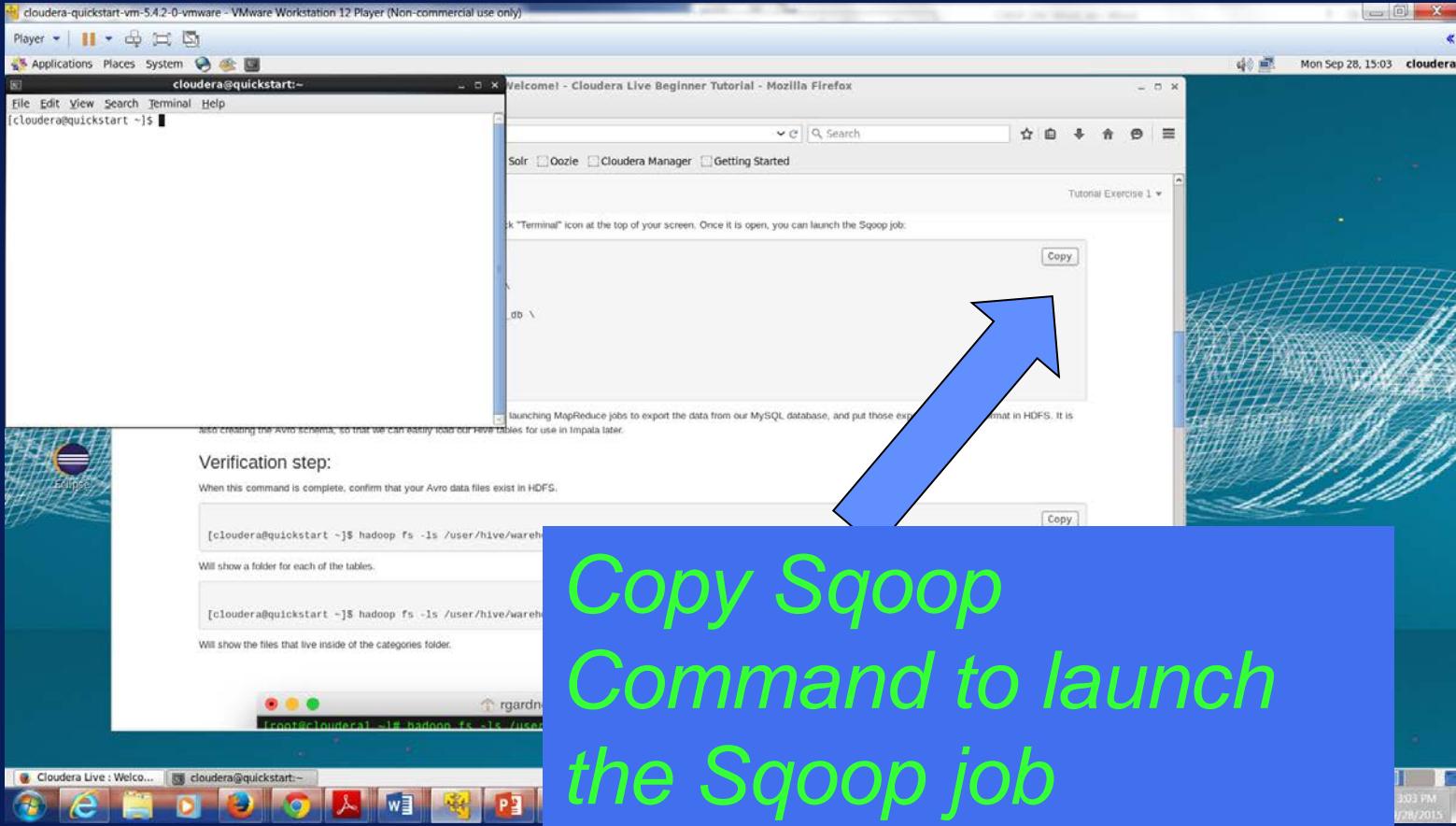
Apache Sqoop, which is part of CDH, is that tool. The nice thing about Sqoop is that we can automatically load our relational data from MySQL into HDFS, while preserving the structure.

finished, each piece of data is replicated to ensure reliability, and spread out across the cluster to ensure you can process it in parallel on your cluster.

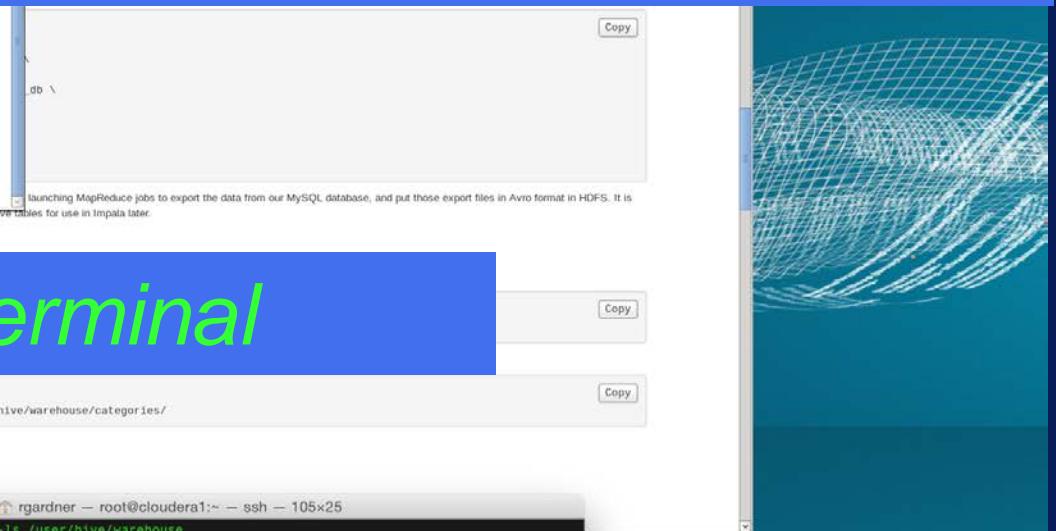
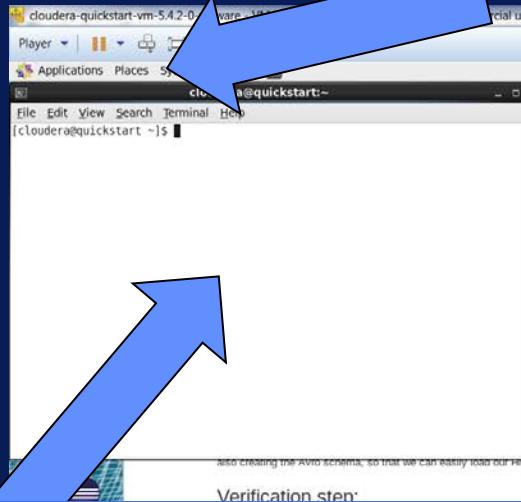
There are 2 versions of Sqoop included in Cloudera's platform. Sqoop 1 is a "thick client" and is what you use in this tutorial. The command you run will directly submit the MapReduce jobs to transfer the data. Sqoop 2 consists of a central server that submits the MapReduce jobs on behalf of clients, and a much lighter weight client that you use to connect to the server. The "Sqoop" you see in Cloudera Manager is the Sqoop 2 server, although Cloudera Manager will make sure that both the "sqoop" and "sqoop2" commands are correctly configured on all your machines.







*Click Edit -> Paste  
Or  
Shift-Ctrl-V*



*Paste into the Terminal*



```
Total committed heap usage (bytes)=223870976
File Input Format Counters
  Bytes Read=8
File Output Format Counters
  Bytes Written=53535
15/09/28 15:07:28 INFO mapreduce.ImportJobBase: Transferred 52.2003 KB in 15.771
7 seconds (3.3148 KB/sec)
15/09/28 15:07:28 INFO mapreduce.ImportJobBase: Retrieved 1345 records.
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ hadoop fs -ls /user/hive/warehouse
Found 6 items
drwxr-xr-x  - cloudera hive      0 2015-09-28 15:05 /user/hive/warehouse/ca
tmporiles
drwxr-xr-x  - cloudera hive      0 2015-09-28 15:06 /user/hive/warehouse/cu
stomers
drwxr-xr-x  - cloudera hive      0 2015-09-28 15:06 /user/hive/warehouse/de
partments
drwxr-xr-x  - cloudera hive      0 2015-09-28 15:06 /user/hive/warehouse/or
der_items
drwxr-xr-x  - cloudera hive      0 2015-09-28 15:06 /user/hive/warehouse/or
ders
drwxr-xr-x  - cloudera hive      0 2015-09-28 15:07 /user/hive/warehouse/pr
ducts
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$
```

Verification step:  
When this command is complete, confirm that your Avro data files exist in HDFS.

# Confirm that data files exist in HDFS



cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System Mozilla Firefox Mon Sep 28, 15:24 cloudera

File Edit View Search Terminal Help

```
cloudera@quickstart:~
```

```
drwxr-xr-x - cloudera hive 0 2015-09-28 15:05 /user/hive/warehouse/categories
drwxr-xr-x - cloudera hive 0 2015-09-28 15:06 /user/hive/warehouse/customers
drwxr-xr-x - cloudera hive 0 2015-09-28 15:06 /user/hive/warehouse/departments
drwxr-xr-x - cloudera hive 0 2015-09-28 15:06 /user/hive/warehouse/oder_items
drwxr-xr-x - cloudera hive 0 2015-09-28 15:07 /user/hive/warehouse/orders
drwxr-xr-x - cloudera hive 0 2015-09-28 15:07 /user/hive/warehouse/products
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ hadoop fs -ls /user/hive/warehouse/categories/
Found 2 items
-rw-r--r-- 1 cloudera hive 0 2015-09-28 15:05 /user/hive/warehouse/categories/_SUCCESS
-rw-r--r-- 1 cloudera hive 1344 2015-09-28 15:05 /user/hive/warehouse/categories/part-m-00000.avro
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -mkdir /user/examples
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -chmod +rw /user/examples
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ hadoop fs -copyFromLocal ~/*.avsc /user/examples/
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$
```

[cloudera@quickstart ~]\$ sudo -u hdfs hadoop fs -mkdir /user/examples
[cloudera@quickstart ~]\$ sudo -u hdfs hadoop fs -chmod +rw /user/examples
[cloudera@quickstart ~]\$ hadoop fs -copyFromLocal ~/\*.avsc /user/examples/

Now that we have the data, we can prepare it to be queried. We're going to do this in the next section using Impala, but you may notice we imported this data into Hive's directories. Hive and Impala both read their data from HDFS, and they even share metadata about the tables. The difference is that Hive executes queries by compiling them to MapReduce jobs. As you will see later, this means it can be more flexible, but is much slower. Impala is an MPP query engine that reads the data directly from the file system itself. This allows it to execute queries fast enough for interactive analysis and exploration.

If one of these steps fails, please reach out to our [Cloudera Live Forum](#) and get help.

## CONCLUSION

Now you have gone through the first basic steps to Sqoop structured data into HDFS, transform it into Avro file format (you can read about the benefits of Avro as a common format in Hadoop [here](#)), and import the schema files for use when we query this data.

< Getting Started Tutorial Exercise 2 >

Cloudera Live : Welcome cloudera@quickstart:~

3:24 PM 9/28/2015

cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial - Mozilla Firefox

quickstart.cloudera#tutorial-ingest\_structured\_data

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

cloudera LIVE

Navigation

departments

- department\_id INT(11)
- department\_name VARCHAR(40)

categories

- category\_id INT(11)
- category\_department\_id INT(11)
- category\_name VARCHAR(45)

products

- product\_id INT(11)
- product\_name VARCHAR(45)
- product\_description VARCHAR(255)
- product\_price FLOAT
- product\_image VARCHAR(255)

order\_items

- order\_item\_id INT(11)
- order\_item\_order\_id INT(11)
- order\_item\_product\_id INT(11)
- order\_item\_quantity TINYINT(4)
- order\_item\_subtotal FLOAT
- order\_item\_product\_price FLOAT

customers

- customer\_id INT(11)
- customer\_name VARCHAR(45)
- customer\_email VARCHAR(45)
- customer\_password VARCHAR(45)
- customer\_street VARCHAR(255)
- customer\_city VARCHAR(45)
- customer\_state VARCHAR(45)
- customer\_zipcode VARCHAR(45)

orders

- order\_id INT(11)
- order\_date DATETIME
- order\_customer\_id INT(11)
- order\_status VARCHAR(45)

Tutorial Exercise 1

finished, each piece of data is replicated to ensure redundancy and spread out across the cluster to ensure it can process it in parallel.

There are two main approaches to moving data from a relational database to a Hadoop platform: "MapReduce" and "Sqoop". We will focus on Sqoop in this tutorial.

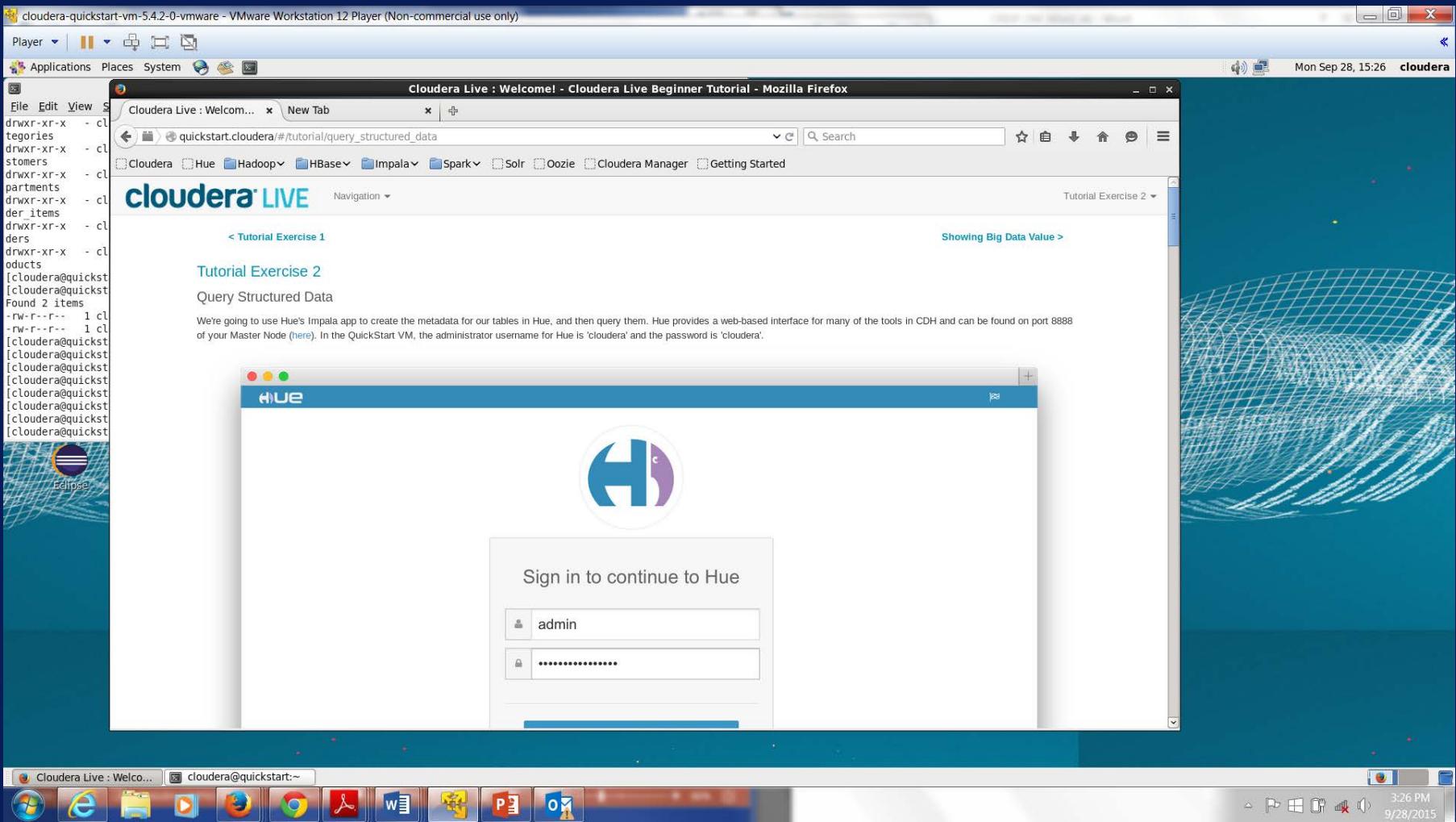
When you run Sqoop, it will submit the MapReduce job to Hadoop. Sqoop has two parts: a client and a server. The client connects to the database and submits the MapReduce job to the server. The server then runs the job on behalf of clients, and a

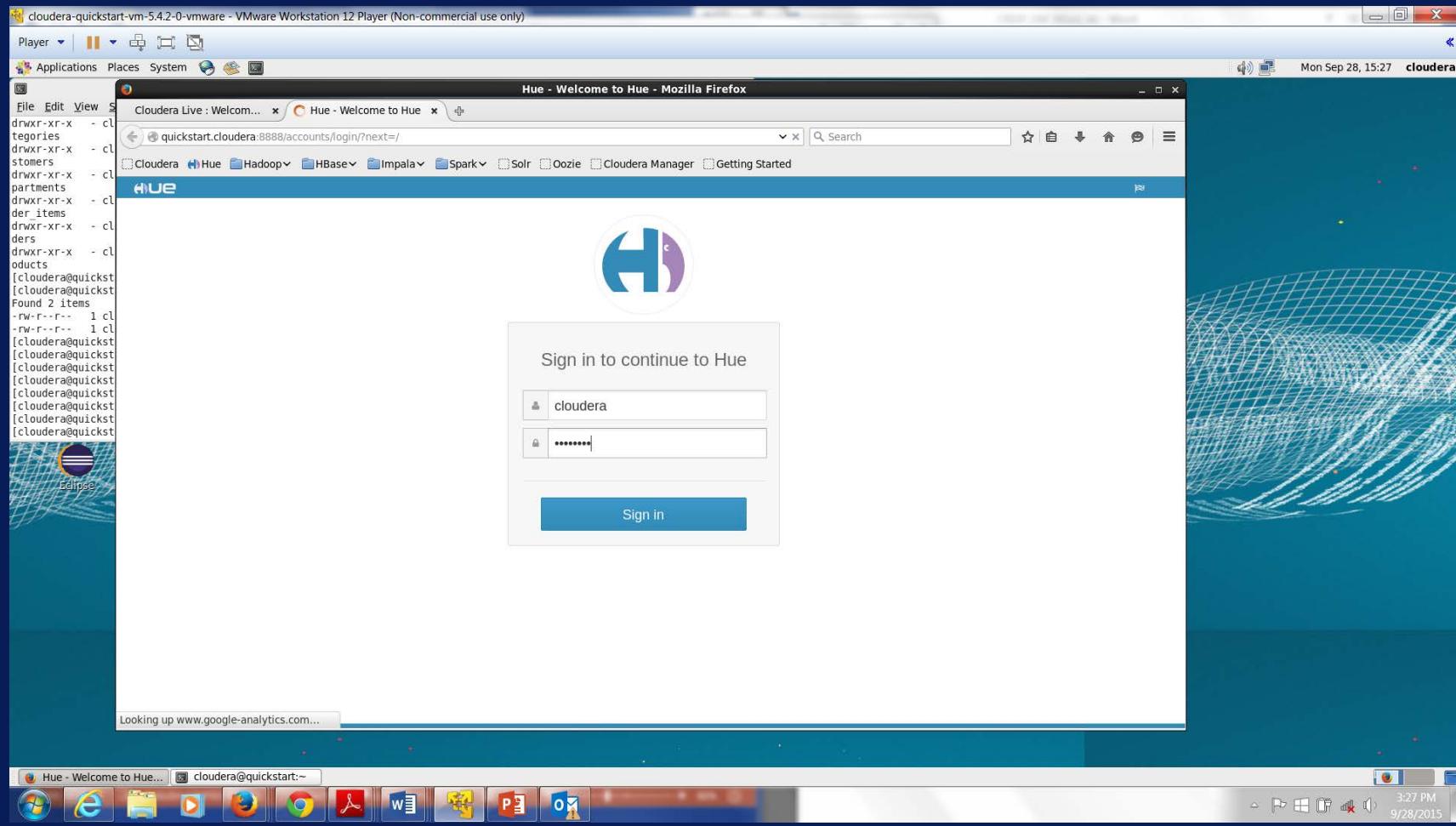
To analyze the transaction data in the new platform, we need to ingest it into the Hadoop Distributed File System (HDFS). We need to find a tool that easily transfers structured data from a RDBMS to HDFS, while preserving structure. That enables us to query the data, but not interfere with or break any regular workload on it.

Apache Sqoop, which is part of CDH, is that tool. The nice thing about Sqoop is that we can automatically load our relational data from MySQL into HDFS, while preserving the structure.

Next step

The screenshot shows a desktop environment with a VMware player window open. Inside the window, a Firefox browser displays a Cloudera Live tutorial page. The page features a database schema diagram with tables like departments, categories, products, order\_items, customers, and orders. A large blue arrow points from the text box on the right towards the schema diagram. The text box contains instructions about using Sqoop to move data from a relational database to HDFS. A large green 'Next step' button is overlaid on the bottom left of the text box. The desktop taskbar at the bottom shows various application icons, and the system tray indicates the date and time as Mon Sep 28, 14:57 2015.





Player

Applications Places System

Mon Sep 28, 15:31 cloudera

Cloudera Live : Welcome to Hue - Quick Start

127.0.0.1:8888/about/admin\_wizard

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

HUE Quick Start Configuration Server Logs

About Quick Start Configuration Examples Users Go!

Checking current configuration

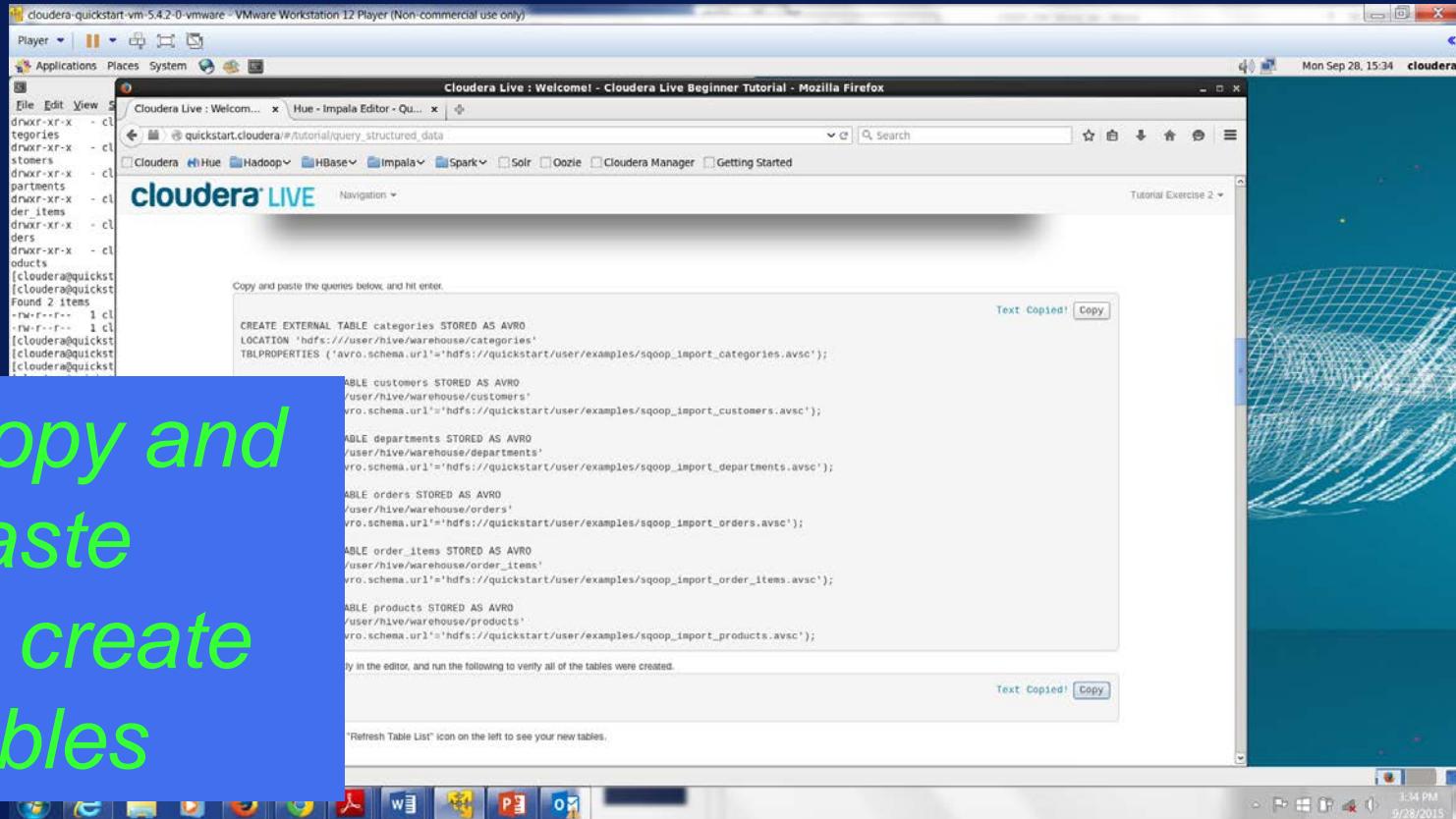
Back Next

Looking up www.google-analytics.com...

Hue and the Hue logo are trademarks of Cloudera, Inc.



Click on Impala



Copy and  
paste  
to create  
tables

cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

Hue - Impala Editor - cloudera@quickstart:~

File Edit View S  
Categories  
Customers  
Departments  
Order Items  
Orders  
Products

Cloudera Live : Welcome | Hue - Impala Editor - ... | 127.0.0.1:8888/impala/execute/query/3#query/results

HUE Query Editors Data Browsers Workflows Search Security File Browser Job Browser cloudera

Impala Query Editor My Queries Saved Queries History

DATABASE default

Table name...  
categories  
customers  
departments  
order\_items  
orders  
products

1 show tables;  
2  
3

Execute Save as... Explain or create a New query

Recent queries Query Log Columns Results Chart

	name
0	categories
1	customers
2	departments
3	order_items
4	orders
5	products

Looking up www.google-analytics.com...

Hue - Impala Editor - cloudera@quickstart:~

3:33 PM 9/28/2015

First delete the queries currently in the editor

Copy the show tables command into the window

The screenshot shows the Hue Impala Editor interface. In the top-left corner, there's a terminal window displaying file permissions and directory contents. The main area is the Query Editor, which has a 'show tables;' command entered. Below the editor is a results table showing six rows of table names. A large blue arrow points from the text 'Copy the show tables command into the window' to the 'show tables;' command in the editor. Another blue arrow points from the text 'First delete the queries currently in the editor' to the results table.

Hue - Impala Editor - Query - Mozilla Firefox

File Edit View S... Cloudera Live : Welcome... / Hue - Impala Editor - Qu... 127.0.0.1:8888/impala/execute/query/4#query/results Mon Sep 28, 15:40 cloudera

Player Applications Places System

Cloudera Live : Welcome... / Hue - Impala Editor - Qu... 127.0.0.1:8888/impala/execute/query/4#query/results

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

HUE Query Editors Data Browsers Workflows Search Security File Browser Job Browser

Impala Query Editor My Queries Saved Queries History

Assist Settings

DATABASE default

Table name: categories

categories customers departments order\_items orders products

```
1 -- Most popular product categories
2 select category_name, count(order_item.quantity) as count
3 from order_items
4 inner join products p on oi.order_item.product_id = p.product_id
5 inner join categories c on c.category_id = p.product_category_id
6 group by c.category_name
7 order by count desc
8 limit 10;
9
10
11;
```

Execute Save as... Explain or create a New query

Recent queries Query Log Columns Results Chart

	category_name	count
0	Cleats	24551
1	Men's Footwear	22246
2	Women's Apparel	21035
3	Indoor/Outdoor Games	19298
4	Fishing	17325
5	Water Sports	15540

Hue - Impala Editor - ... cloudera@quickstart:~

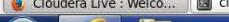
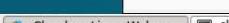
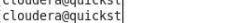
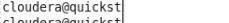
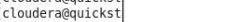
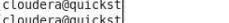
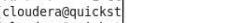
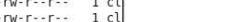
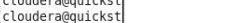
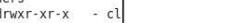
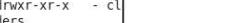
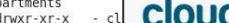
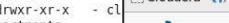
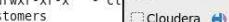
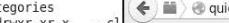
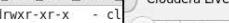
Revenue per product

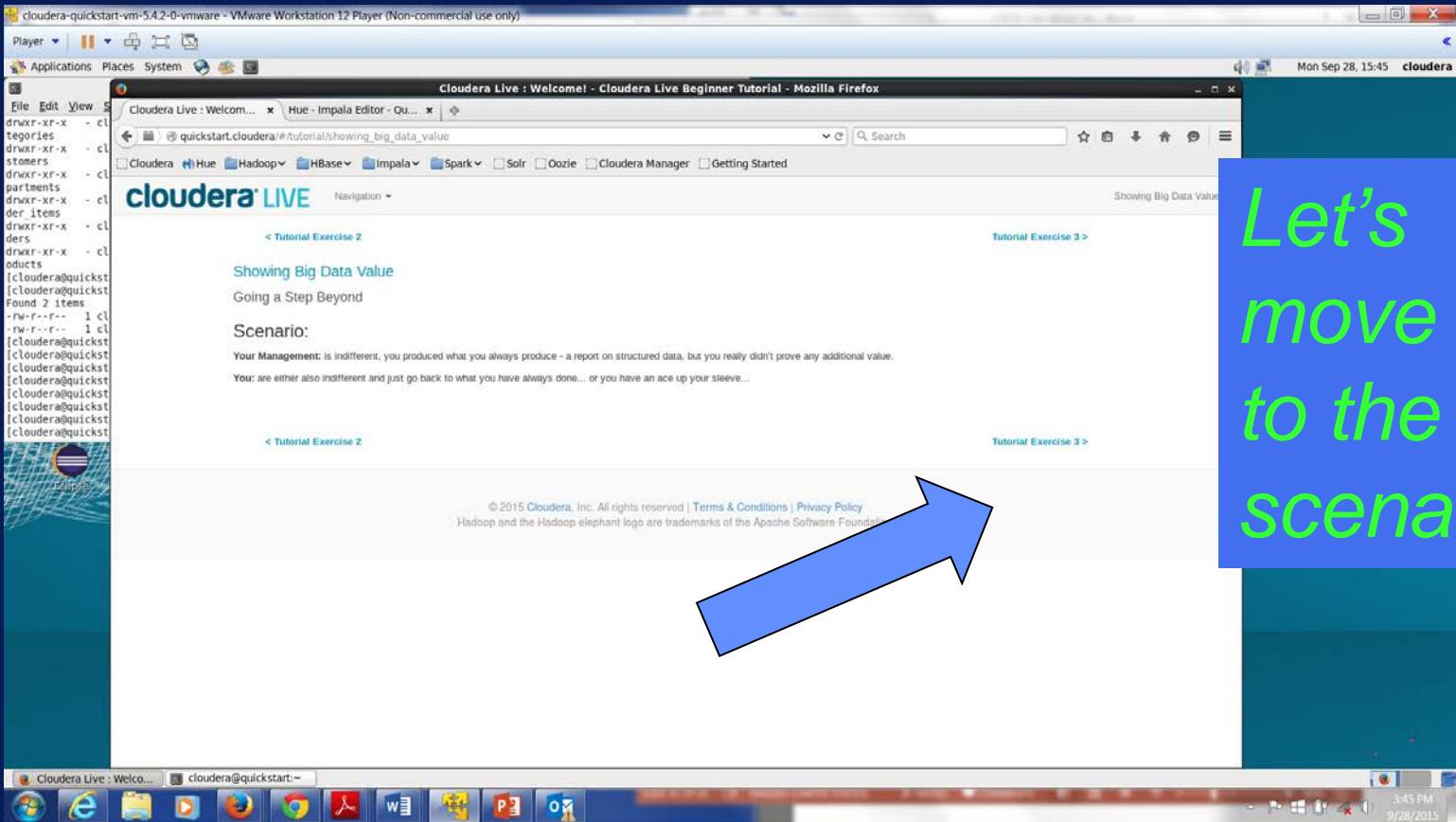


Player

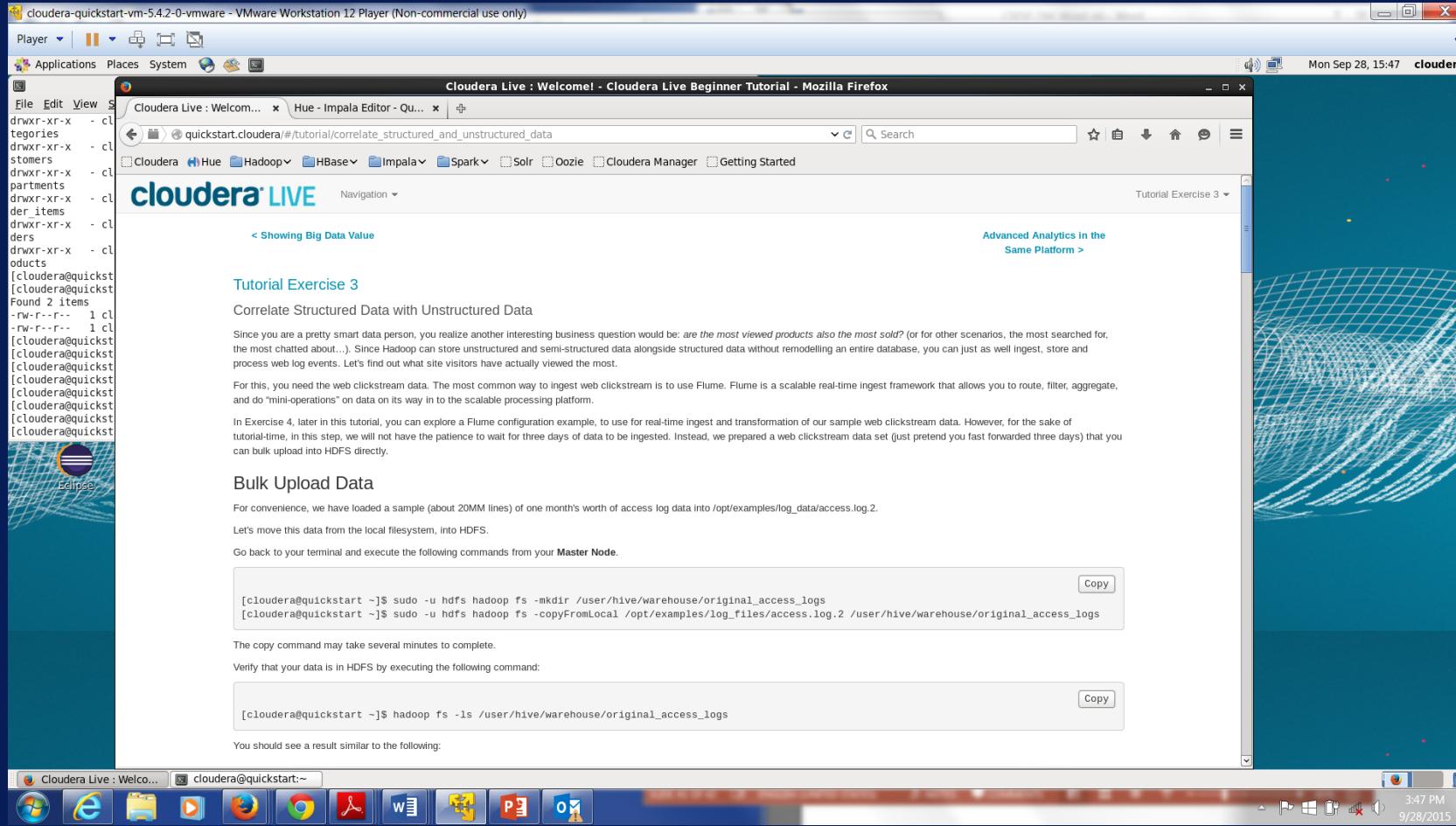
Applications Places System

Mon Sep 28, 15:43 cloudera





*Let's move on to the next scenario!!*



Player

Applications Places System

cloudera@quickstart:~

```
File Edit View Search Terminal Help
products
[cloudera@quickstart ~]$ hadoop fs -ls /user/hive/warehouse/categories/
Found 2 items
-rw-r--r-- 1 cloudera hive          0 2015-09-28 15:05 /user/hive/warehouse/categories/_SUCCESS
-rw-r--r-- 1 cloudera hive    1344 2015-09-28 15:05 /user/hive/warehouse/categories/part-m-00000.avro
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -mkdir /user/examples
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -chmod +rw /user/examples
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -chown root.root /user/examples
[cloudera@quickstart ~]$ hadoop fs -copyFromLocal ~/*.avsc /user/examples/
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -mkdir /user/hive/warehouse/original_access_logs
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -copyFromLocal /opt/examples/log_files/access.log.2 /user/hive/warehouse/
original_access_logs
[cloudera@quickstart ~]$ 
[cloudera@quickstart ~]$ 
```

Go back to your terminal and execute the following commands from your **Master Node**.

```
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -mkdir /user/hive/warehouse/original_access_logs
[cloudera@quickstart ~]$ sudo -u hdfs hadoop fs -copyFromLocal /opt/examples/log_files/access.log.2 /user/hive/warehouse/original_access_logs
```

Text Copied!

Copy

The copy command may take several minutes to complete.

Verify that your data is in HDFS by executing the following command:

```
[cloudera@quickstart ~]$ hadoop fs -ls /user/hive/warehouse/original_access_logs
```

Copy

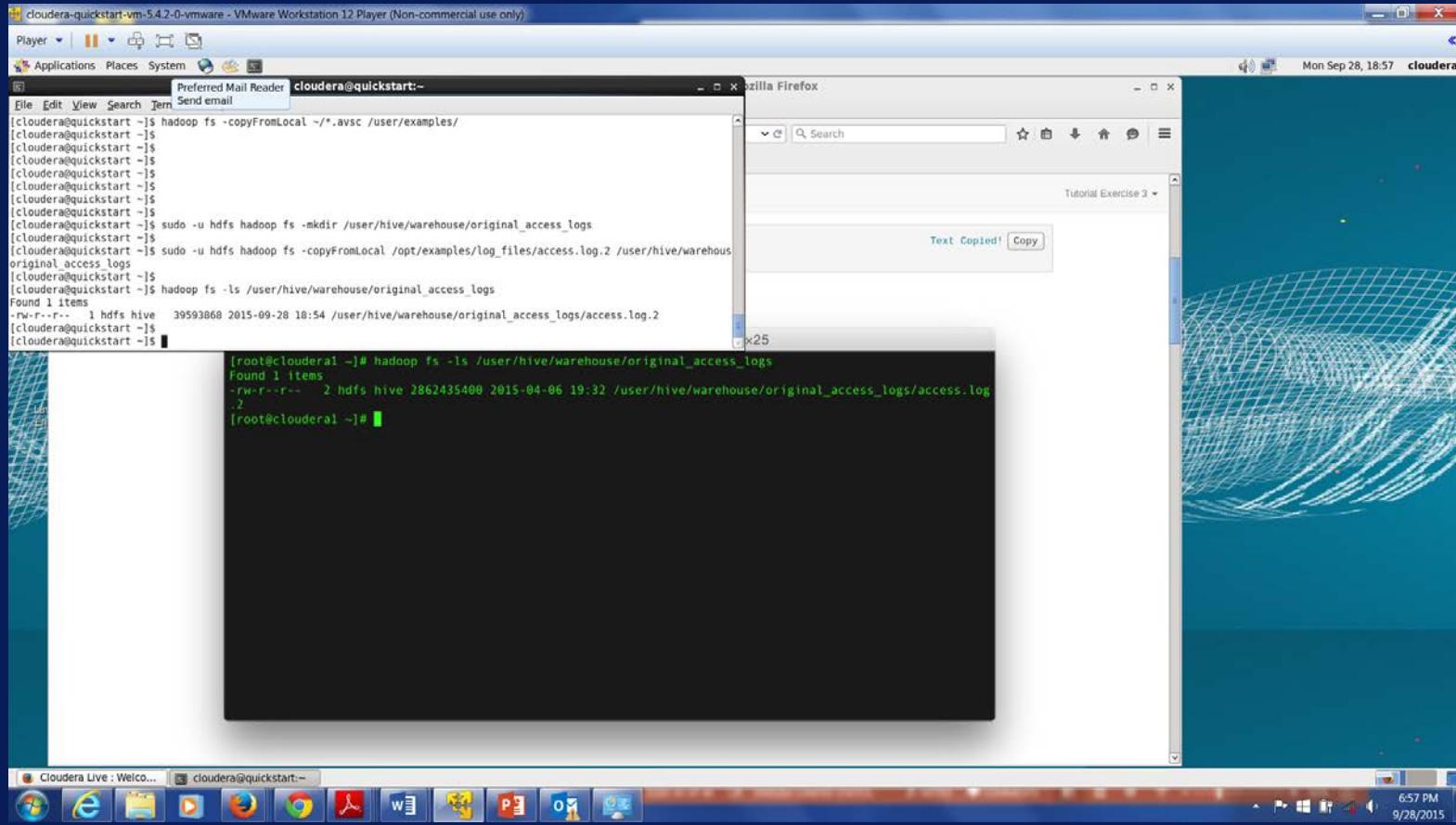
You should see a result similar to the following:

```
[root@cloudera1 ~]# hadoop fs -ls /user/hive/warehouse/original_access_logs
Found 1 items
-rw-r--r-- 2 hdfs hive 2862435400 2015-04-06 19:32 /user/hive/warehouse/original_access_logs/access.log
```

Cloudera Live : Welcome cloudera@quickstart:~



Copy data into HDFS



Player

Applications Places System



Mon Sep 28, 19:04

cloudera

## Cloudera Live : Welcome! - Cloudera Live Beginner Tutorial - Mozilla Firefox

File E Cloudera Live : Welcom... x Hue - Impala Editor - Qu... x +  
0: jdbc  
0: jdbc ↺ ↻ quickstart.cloudera/#/tutorial/correlate\_structured\_and\_unstructured\_data  
0: jdbc  
0: jdbc  
0: jdbc Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started



Navigation ▾

Tutorial Exercise 3 ▾

logs into individual fields using a regular expression. Second, you'll transfer the data from this intermediate table to one that does not require any special SerDe. Once the data is in this table, you can query it much faster and more interactively using Impala.

We'll query Hive using a command-line JDBC client for Hive called Beeline. You can invoke it from the terminal with the following:

```
[cloudera@quickstart ~]$ beeline -u jdbc:hive2://quickstart:10000/default -n admin -d org.apache.hive.jdbc.HiveDriver
```

Text Copied! Copy

Once the Beeline shell is connected, run the following queries:

```
0: jdbc:hive2://quickstart:10000/default> CREATE EXTERNAL TABLE intermediate_access_logs (
    ip STRING,
    date STRING,
    method STRING,
    url STRING,
    http_version STRING,
    code1 STRING,
    code2 STRING,
    dash STRING,
    user_agent STRING)
ROW FORMAT SERDE 'org.apache.hadoop.hive.contrib.serde2.RegexSerDe'
WITH SERDEPROPERTIES (
    'input.regex' = '([^\"]* - - \[\[(\[\\" \])*\\]\] ([^\"]*) ([^\"]*)" ([^\"]*)" (\\"d*) (\\"d*) "([^\"]*)" "([^\"]*)"',
    'output.format.string' = "%1$s %2$s %3$s %4$s %5$s %6$s %7$s %8$s %9$s"
)
LOCATION '/user/hive/warehouse/original_access_logs';
```

Text Copied! Copy

```
0: jdbc:hive2://quickstart:10000/default> CREATE EXTERNAL TABLE tokenized_access_logs (
    ip STRING,
    date STRING,
    method STRING,
    url STRING,
    http_version STRING,
    code1 STRING,
    code2 STRING,
    dash STRING,
    user_agent STRING)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
```

Cloudera Live : Welco...

cloudera@quickstart:~



cloudera-quickstart-vm-5.4.2-0-vmware - VMware Workstation 12 Player (Non-commercial use only)

Player Applications Places System

File Edit View Search Terminal Help

```
bash: url: command not found
[cloudera@quickstart ~]$ bash: http_version STRING,
bash: http_version: command not found
[cloudera@quickstart ~]$ bash: code1 STRING,
bash: code1: command not found
[cloudera@quickstart ~]$ bash: code2 STRING,
bash: code2: command not found
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ dash STRING,
dash: Can't open STRING,
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ user agent STRING)
bash: syntax error near unexpected token `)'
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
bash: ROW: command not found
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ LOCATION '/user/hive/warehouse/tokenized_access_logs';
bash: LOCATION: command not found
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ [cloudera@quickstart ~]$ ADD JAR /usr/lib/hive/lib/hive-contrib.jar;
bash: ADD: command not found
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ [cloudera@quickstart ~]$ [cloudera@quickstart ~]$ INSERT OVERWRITE TABLE tokenized_access_logs SELECT * FROM intermediate_access_logs;
bash: INSERT: command not found
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$ [cloudera@quickstart ~]$ [cloudera@quickstart ~]$ !quit
bash: !quit: event not found
[cloudera@quickstart ~]$ [cloudera@quickstart ~]$
```

where url like '%\\product\\%'  
group by url order by count(\*) desc;

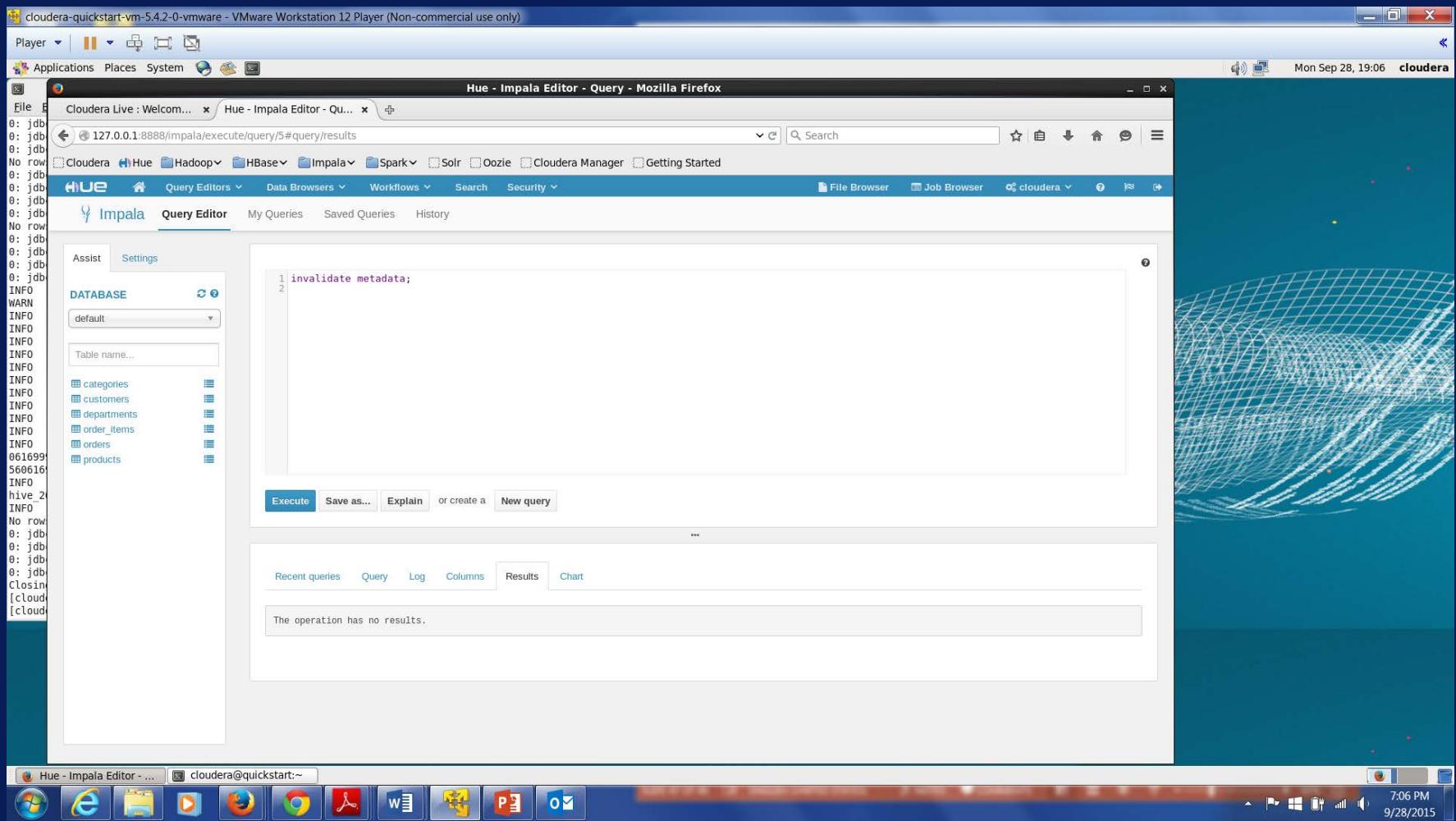
You should see a result similar to the following:

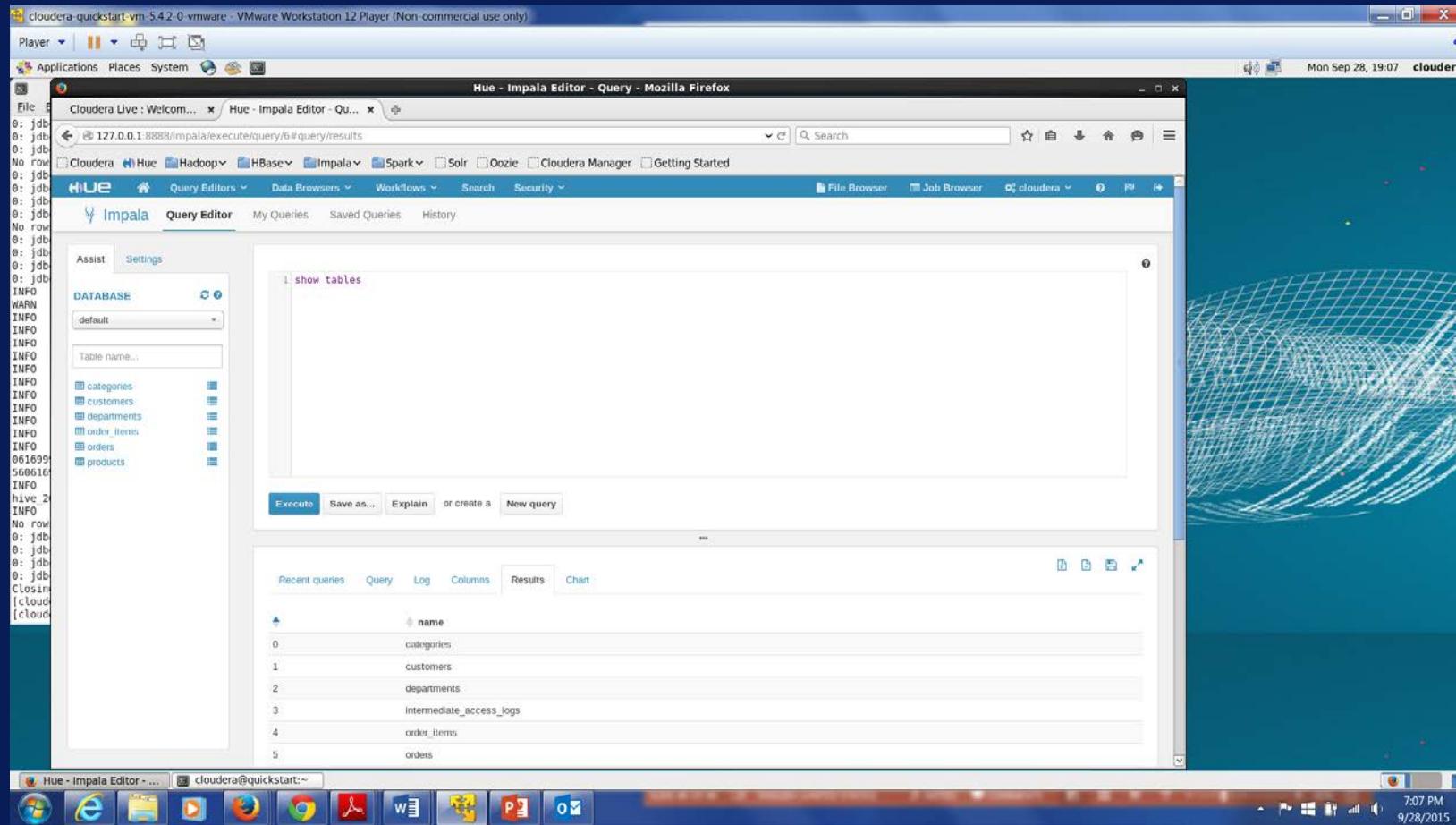
HUE Query Editors Data Browsers Workflows Search Security

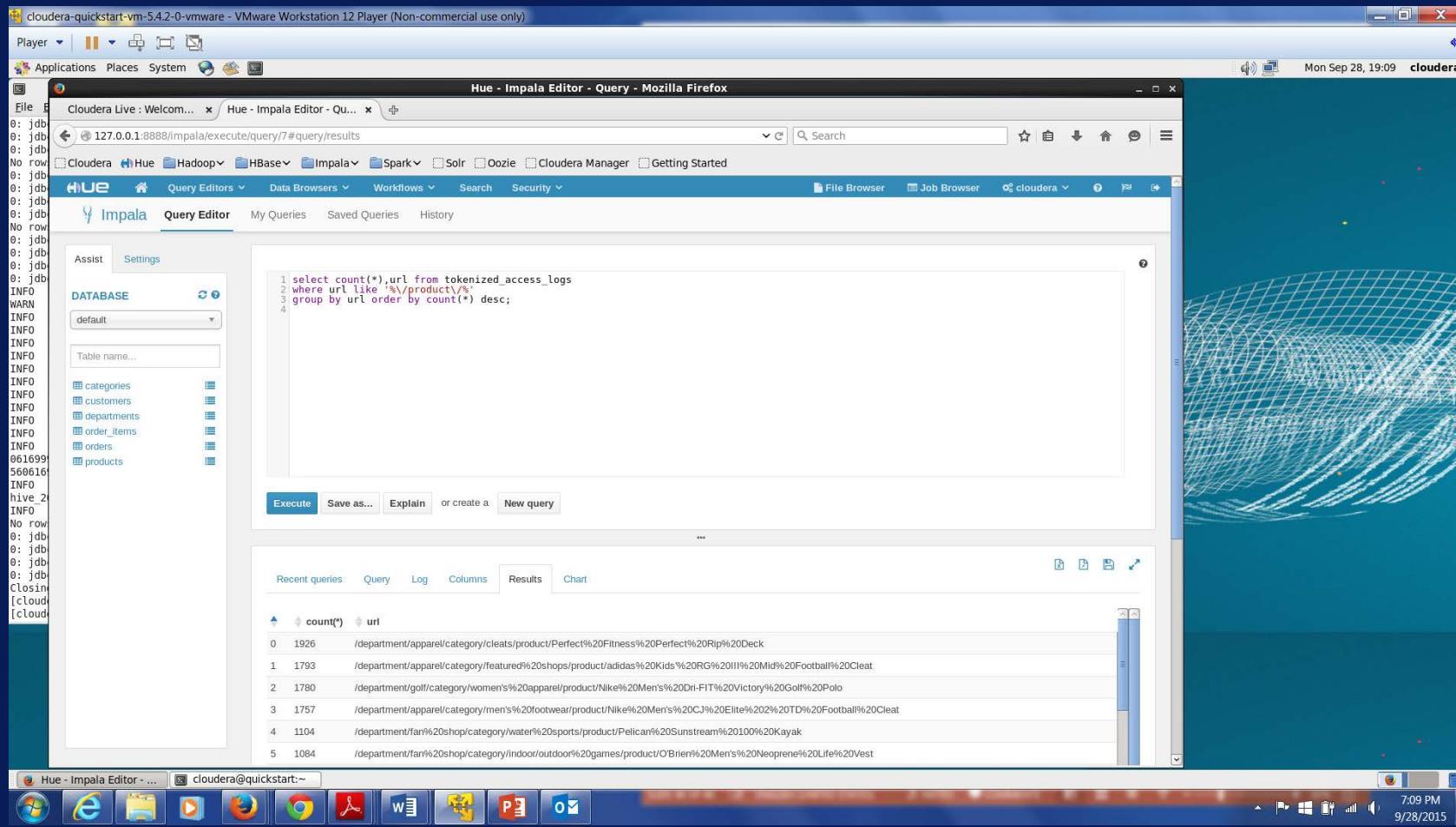
Impala Query Editor My Queries Saved Queries History

Cloudera Live : Welcome cloudera@quickstart:~

6:59 PM 9/28/2015







Well, in our example with DataCo, once these odd findings are presented to your manager, it is immediately escalated. Eventually, someone figures out that on that view page, where most visitors

# Assignment