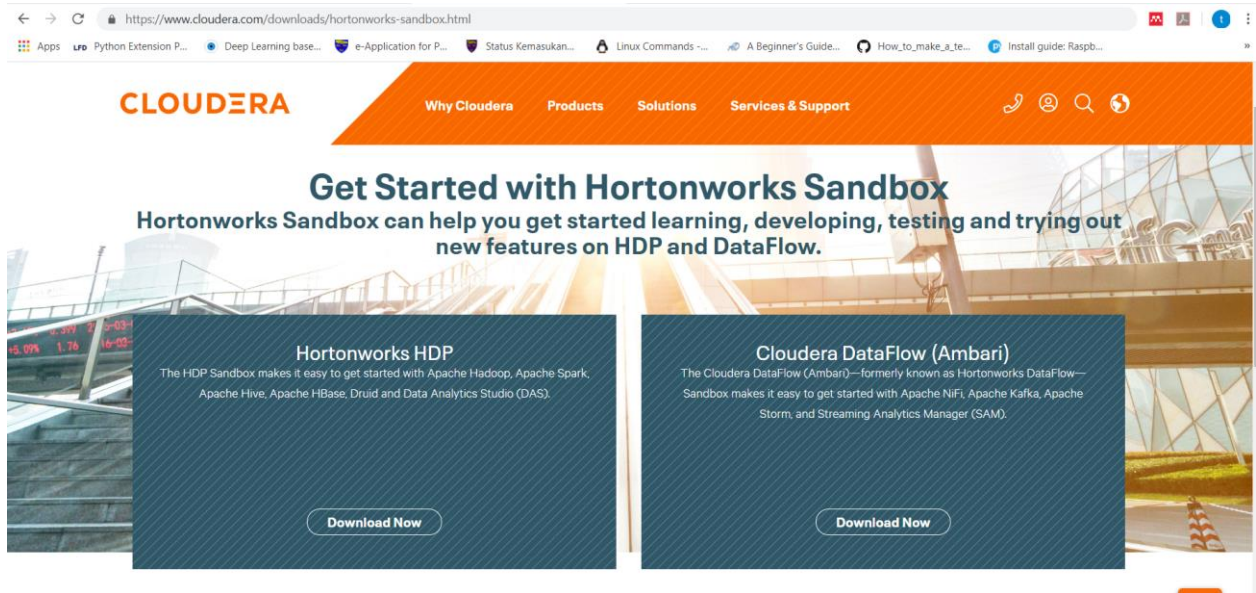Tan Sia Hong & Tan Chang Jung
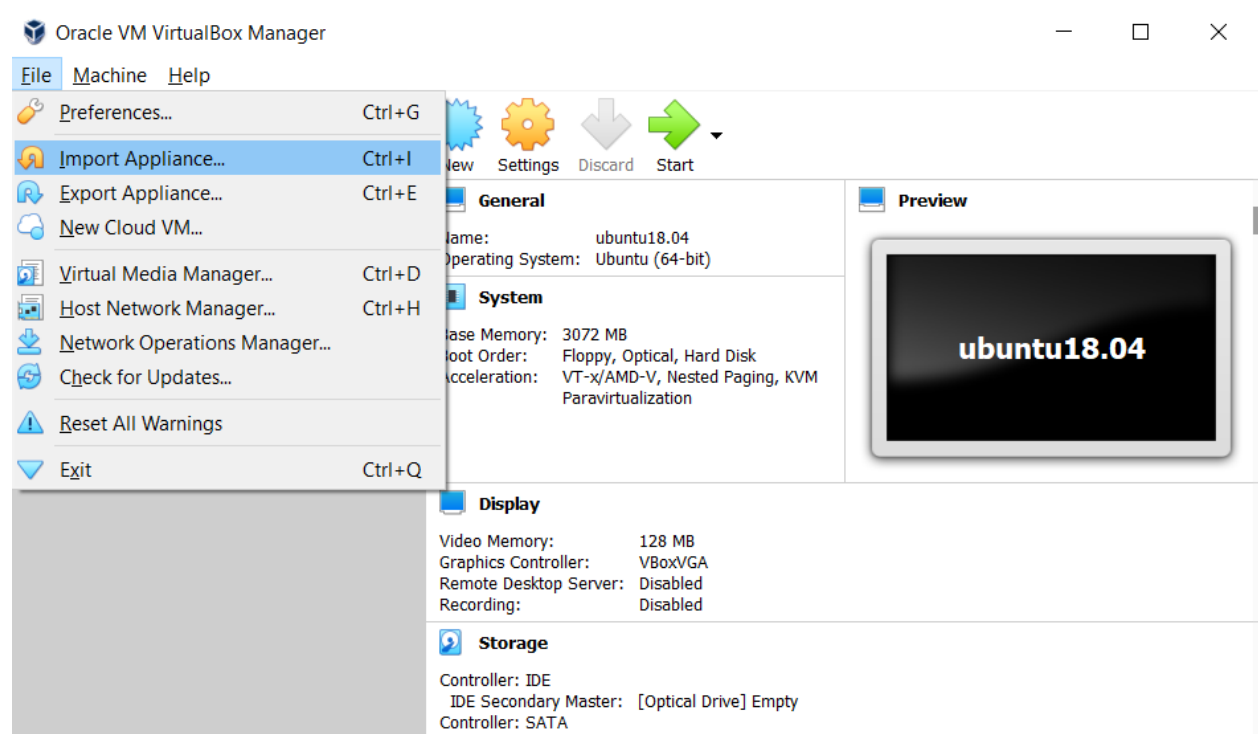
# Milestone 2: Stone data set into Hive data warehouse

## Part 1: Install Hortonworks HDP

1.  Download Hortonworks HDP from https://www.cloudera.com/downloads/hortonworks-sandbox.html (20GB)



2.  Import Hortonworks.ova into VM VirtualBox.
    VM VirtualBox can download from https://www.virtualbox.org/wiki/Downloads

Select the hortonworks.ova

? ✕

← Import Virtual Appliance

## Appliance to import

Please choose the source to import appliance from. This can be a local file system to import OVF archive or one of known cloud service providers to import cloud VM from.

Source: | Local File System | ▼

Please choose a file to import the virtual appliance from. VirtualBox currently supports importing appliances saved in the Open Virtualization Format (OVF). To continue, select the file to import below.

File: | D:\Program file\HDP_3.0.1_virtualbox_181205.ova |

Expert Mode | Next | Cancel

← Import Virtual Appliance

## Appliance settings

These are the virtual machines contained in the appliance and the suggested settings of the imported VirtualBox machines. You can change many of the properties shown by double-clicking on the items and disable others using the check boxes below.

| Virtual System 1 | |
|---|---|
| 🍀 Name | Hortonworks Sandbox HDP 3.0 1 |
| 🪟 Guest OS Type | Red Hat (64-bit) |
| 🖥 CPU | 4 |
| RAM | 8192 MB |
| 🖊 USB Controller | ☑ |
| 🔊 Sound Card | ☑ ICH AC97 |
| 🖧 Network Adapter | ☑ Intel PRO/1000 MT Desktop (82540EM) |
| ◈ Storage Controller (IDE) | PIIX4 |
| ⌄ ◈ Storage Controller (IDE) | PIIX4 |
| 💿 Virtual Disk Image | Hortonworks Sandbox HDP 3.0-disk001.vmdk |
| 📁 Base Folder | C:\Users\Tan Chang Jung\VirtualBox VMs |
| 🅰 Primary Group | / |

Machine Base Folder:  📁 C:\Users\Tan Chang Jung\VirtualBox VMs    ⌄

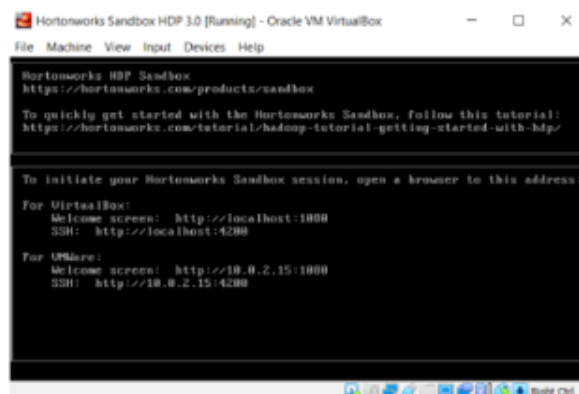MAC Address Policy: Include only NAT network adapter MAC addresses    ▾

Additional Options: ☑ Import hard drives as VDI

Appliance is not signed

[ Restore Defaults ]   [ Import ]   [ Cancel ]

*Optimize the number of CPU cores and RAM resources before import.

3.  Run the Hortonworks Sandbox HDP 3.0 for extraction and installation.  Notice that first time installation will take a few minutes.



## *Installation is completed

4. Use the Google Browser to access http://localhost:4200/ . The default root user login credentials will be:

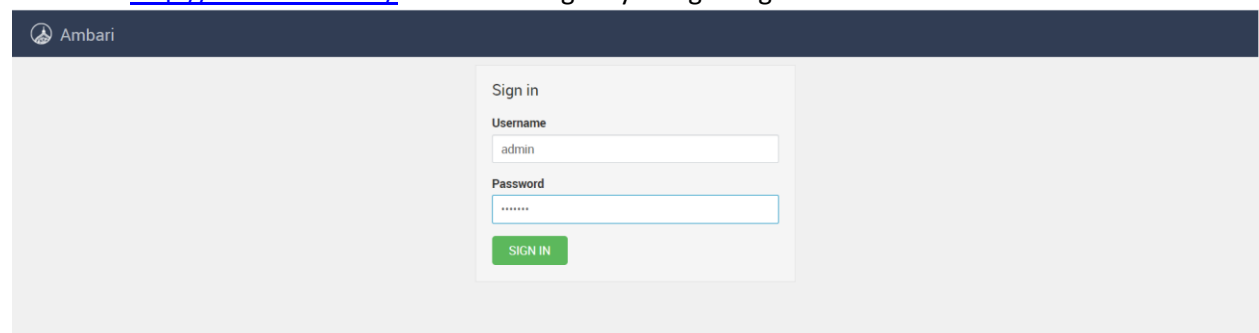User: root

Password: hadoop

After logging by default password, you will request to change password. Please change the password for root user.

```
sandbox-hdp login: root
root@sandbox-hdp.hortonworks.com's password:
You are required to change your password immediately (root enforced)
Last login: Tue Jun  2 14:43:46 2020 from 172.18.0.2
Changing password for root.
(current) UNIX password:
New password:
Retype new password:
[root@sandbox-hdp ~]#
```

5. Ambari enables system administrators to provision, manage and monitor a Hadoop cluster. Now, Type 'ambari-admin-password-reset' to reset the Ambari' administrator password.

```
sandbox-hdp login: root
root@sandbox-hdp.hortonworks.com's password:
Last login: Wed Jun  3 23:18:14 2020
[root@sandbox-hdp ~]# ambari-admin-password-reset
Please set the password for admin:
```

6. Access to http://localhost:8080/ for Ambari login by using Google Brower.

7. After logging in, the sandbox is taking some time to starting all the required services that wait time is depend on laptop performance.



Background Operations ×

0 Background Operations Running                                                                              ALL (10)  ▾

| Operations | Status | | User | Start Time | Duration |
|---|---|---|---|---|---|
| ✔ Start All Services | | 100% | raj_ops | Today 13:14 | 43m 47s |
| ✔ Start All Services | | 100% | raj_ops | Sat Jun 06 2020 23:38 | 49m 11s |
| ✔ Start All Services | | 100% | raj_ops | Tue Jun 02 2020 22:13 | 13m 13s |
| ✔ Start All Services | | 100% | raj_ops | Tue Jun 02 2020 21:31 | 15m 51s |
| ✔ Restart all components for HDFS | | 100% | admin | Sun May 31 2020 16:39 | 3m 30s |
| ✔ Start All Services | | 100% | raj_ops | Sun May 31 2020 16:01 | 11m 52s |

☐ Do not show this dialog again when starting a background operation                                          OK

8. The Hive and HDFS are completely set up and ready to be used.



## Summary                                                                                                   🔔 0

Components

⊘ Started            ⊘ Started
NAMENODE             SNAMENODE

37m 57s              12.3%
NAMENODE UPTIME      123.9 MB / 1011.3 MB
                     NAMENODE HEAP

1/1 Started          0/0 Live             0/0 Started
DATANODES            JOURNALNODES         NFSGATEWAYS

DATANODES STATUS

1                    0                    0
Live                 Dead                 Decommissioning

# PART 2: Store data into Hive database

1. Download and install WinSCP from https://winscp.net/eng/download.php



2. Then, logging in the WinSCP by Host name is 127.0.0.1, port number is 2222, user name is root. Note that File protocol must need SCP type.



3. Key in your password of sandbox hadoop.

4. Send local Windows' data into sandbox (VM) by using WinSCP.



5. Go to web shell client

```
hdfs fs -put /root/data/*.csv /user/root/datamining/data/
```

This is to copy the file from root directory into another directory HDFS.

6. Open hive in web shell client, use the database preferred, then create a table under the database.
   Eg. Table 'bitcoin'

```
CREATE TABLE bitcoin (MarketDate DATE, Open DOUBLE, High DOUBLE, Low DOUBLE, Close
DOUBLE, Volume DOUBLE, MarketCapacity DOUBLE)

ROW FORMAT DELIMITED

FIELDS TERMINATED BY ','

STORED AS TEXTFILE

TBLPROPERTIES("skip.header.line.count"="1");
```

Then load the csv data from HDFS into bitcoin table

```
LOAD DATA INPATH '/user/root/datamining/data/coin.csv' INTO TABLE bitcoin;
```

7. After load the data into Hive table, you can run query 'show tables;' to perform all table as below:

```
0: jdbc:hive2://sandbox-hdp.hortonworks.com:2> show tables;
INFO  : Compiling command(queryId=hive_20200610065014_3f96eded-047d-48b2-b6b3-4ff454c57525): show tables
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:tab_name, type:string, comment:from deserializer)], properties:null)
INFO  : Completed compiling command(queryId=hive_20200610065014_3f96eded-047d-48b2-b6b3-4ff454c57525); Time taken: 0.147 seconds
INFO  : Executing command(queryId=hive_20200610065014_3f96eded-047d-48b2-b6b3-4ff454c57525): show tables
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_20200610065014_3f96eded-047d-48b2-b6b3-4ff454c57525); Time taken: 0.04 seconds
INFO  : OK
+------------------+
|     tab_name     |
+------------------+
| binancecoin      |
| bitcoin          |
| bitcoincash      |
| bitcoinsv        |
| cardano          |
| chainlink        |
| cryptocomcoin    |
| dash             |
| eos              |
| ethereum         |
| ethereumclassic  |
| huobitoken       |
| litecoin         |
| monero           |
| neo              |
| stellar          |
| tether           |
| tron             |
| unussedleo       |
| xrp              |
+------------------+
```

8. The data show out from Hive table by SQL query 'Select * from bitcoin limit 10;'

```
INFO  : Executing command(queryId=hive_20200610065413_1e4351e1-b617-4bc9-b330-5c5f2a99d57f): select * from bitcoin limit 10
INFO  : Completed executing command(queryId=hive_20200610065413_1e4351e1-b617-4bc9-b330-5c5f2a99d57f); Time taken: 0.112 seconds
INFO  : OK
+---------------------+----------------+----------------+---------------+----------------+--------------------+------------------------+
| bitcoin.marketdate  | bitcoin.open   | bitcoin.high   | bitcoin.low   | bitcoin.close  | bitcoin.volume     | bitcoin.marketcapacity |
+---------------------+----------------+----------------+---------------+----------------+--------------------+------------------------+
| 2020-05-25          | 8786.11        | 8951.01        | 8719.67       | 8906.93        | 3.1288157264E10    | 1.63760453116E11       |
| 2020-05-24          | 9212.28        | 9288.4         | 8787.25       | 8790.37        | 3.25188033E10      | 1.61610414643E11       |
| 2020-05-23          | 9185.06        | 9302.5         | 9118.11       | 9209.29        | 2.7727866812E10    | 1.6930549244E11        |
| 2020-05-22          | 9080.33        | 9232.94        | 9008.64       | 9182.58        | 2.9810773699E10    | 1.68807619957E11       |
| 2020-05-21          | 9522.74        | 9555.24        | 8869.93       | 9081.76        | 3.9326160532E10    | 1.66947987864E11       |
| 2020-05-20          | 9725.33        | 9804.79        | 9447.2        | 9522.98        | 3.6546239703E10    | 1.75050963475E11       |
| 2020-05-19          | 9727.06        | 9836.05        | 9539.62       | 9729.04        | 3.9254288955E10    | 1.78831635026E11       |
| 2020-05-18          | 9675.69        | 9906.03        | 9570.36       | 9726.57        | 4.1827139896E10    | 1.78779483464E11       |
| 2020-05-17          | 9374.93        | 9823.0         | 9349.55       | 9670.74        | 4.0084250663E10    | 1.7774540447E11        |
| 2020-05-16          | 9333.24        | 9564.2         | 9260.69       | 9377.01        | 3.6164766408E10    | 1.72340956579E11       |
+---------------------+----------------+----------------+---------------+----------------+--------------------+------------------------+
10 rows selected (2.187 seconds)
```