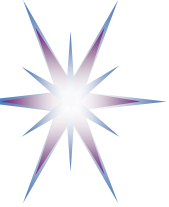# Big Data Infrastructure and Technologies for Data Analytics  (BDIT4DA)

## Practice Guidelines

## Hadoop cluster Installation:
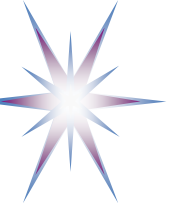## Cloudera Quickstart on VirtualBox

# Outline

- Hadoop cluster VMs
  - Cloudera Qyuickstart 5.13 – 5.3 GB (discontinued but version 5.13 available)
  - Hortonworks Sandbox – 10 GB
  - Oracle Big Data Lite 4.11 – 22 GB
  - Bitnami Hadoop Stack – 1.5 GB (only CLI and MapReduce)

- Oracle VirtualBox and pre-requisite configuration
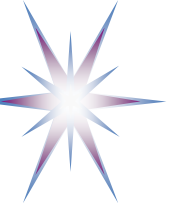- Installing Cloudera Quickstart Hadoop cluster on VirtualBox

Preparation

- Download Cloudera Quickstart 5.13 from below location (not available at Cloudera website since 2020)
  - https://surfdrive.surf.nl/files/index.php/s/2OvUuw5chV42Zwz
  - Directory "vm-software"
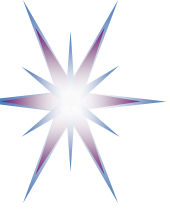- Configure VirtualBox Host-Only Network Adapter

# Cloudera Quickstart VM for VirtualBox

- QuickStart virtual machines (VMs) provide sufficient functionality to try CDH, Cloudera Manager, Impala, and Cloudera Search
  - Minimum required RAM on host machine is 4GB, recommended 8GB
- Cloudera Manager is installed in the VM but is turned off by default.
  - Cloudera strongly recommends that if you use Cloudera Manager, you configure the VM with a minimum of 8 GB RAM and two virtual CPU cores.
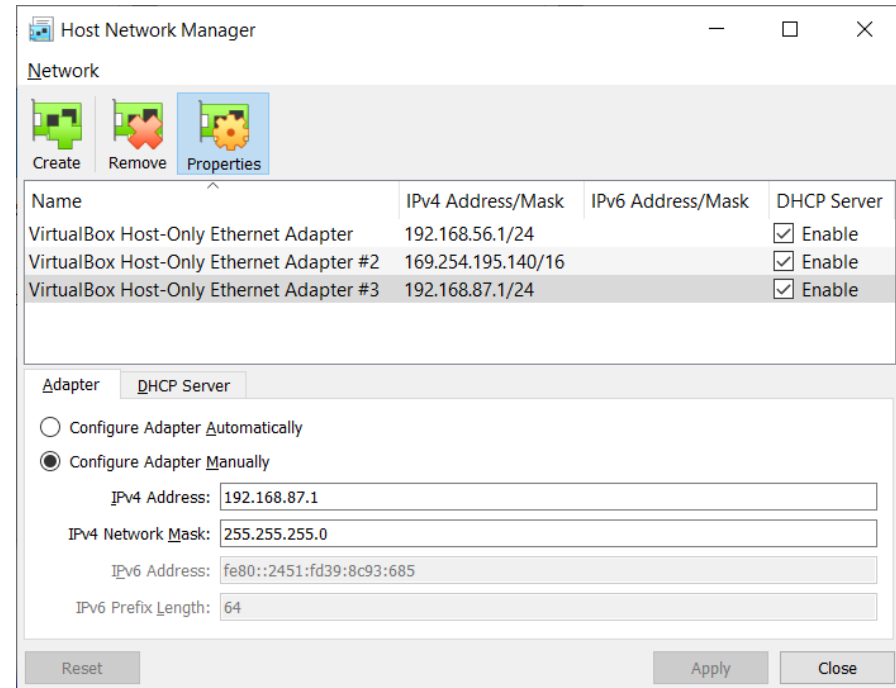
# VirtualBox Configuration – Host requirements

- Pre-requisite and requirement to host

- Hosting Cloudera Quickstart 5.13 VM require
  - Single node minimum RAM 4GB, recommended 8GB
  - Single node + Cloudera Manager minimum 8GB

- Host machine HardDrive space:
  - Quickstart 5.13 VM for VirtualBox – 5.3 GB
  - VM image installed on the host machine – 9-12 GB
  - Plus additionally for a snapshot 1.2-2.5 GB

- VirtualBox Installation
  - Note: For Windows, you can configure VirtualBox VMs location on drive different from C: (which may be have limited space)

# Accessing your VirtualBox Guest from your Host OS: Adding Host-Only network interface

- You need to do is to setup a VirtualBox Host-Only Network and get our guest connected to it.
  - If your guest machine is running, shut it down first.
- New VirtualBox versions
  - Click on File->Host Network Manager in the VirtualBox menu-bar.
  - Use Create/Properties Network Adapter
- Older VirtualBox versions
  - Click on File->Preferences
  - Select the Network option from the side menu and add network adapter

- Configure Adapter Automatically or Manually for better control
  - The default options for the newly-created Host-only network should be fine.
  - If not, you can add the data manually, by clicking on the Edit button in the DHCP Server tab.
  - Save all the settings in Preferences.
- Now open up the settings of your Guest machine and navigate to the Network option from the side menu and click on the Adapter 2 tab.
  - Don't forget to check the Enable Network Adapter option.
  - Save these settings and boot into your Guest machine.
  - After logging in, type ifconfig. Note the new IP, it should be under a new interface like eth0 or vboxnet0. Now you can use this IP to SSH, view the webpages on your machine's Apache Server, etc..

# Importing VM: Select+Click on .ova File

# VirtualBox Host-Only Network Configuration



- File > Host Network Manager
- > Create Host-Only Network

- Configure Host-Only Adapters with addresses of your local network addresses
- Beneficially do it e.g. with address range of home and work network

# Guest VM Network Configuration



- Settings > Network
- Adapter 1 > Host-Only Network
- Adapter 2 > NAT

- Configure Host-Only Adapters with addresses of your local network addresses
- Beneficially do it e.g. with address range of home and work network

BDIT4DA HOL                    Cloudera Q

# Add Shared Folder: e.g. Exercise_Files_all



- Advice: Place all your working files in one folder on your drive
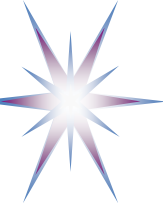
# Cloudera Quickstart VM for VirtualBox



Accounts

- Once you launch the VM, you are automatically logged in as the cloudera user. The account details are:
  - username: cloudera
  - password: cloudera
- The cloudera account has sudo privileges in the VM. The root account password is cloudera.
- The root MySQL password (and the password for other MySQL user accounts) is also cloudera.
- Hue and Cloudera Manager use the same credentials.

# Find IP address of your VM

- Debian
  - $ ip address show
- Ubuntu
  - $ ipconfig
- Example:
  - 192.168.102.154 – LAN
  - 192.168.163.3 – VirtualBox Host only
- Use for configuring your SSH or SCP client

# Add Permission for User "cloudera"
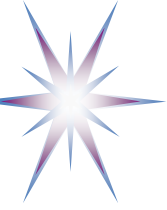


- Start Cloudera
- Open Terminal
- sudo gedit /etc/group
- Add line
  - vboxsf:x:474:cloudera
- Save
- Logout from the system & Login
  - System > Logout cloudera

# Working with Cloudera Quickstart

- Mount your local/host directory with exercises to Cloudera HDFS
- Use Hue visual interface
  - Similar to AWS EMR
- Use Query dropdown menu to select Query>Editor Pig, Hive, Java, Scala, etc
- Use Top-left dropdown menu to select:
  - Apps: Editor, Scheduler
  - Browser: Documents, Files, Tables, Indexes, Jobs
- Use Browsers > Files to Upload your datasets or script
  - Select right button Upload > Files
  - When connected with browser via SSH Tunnel, you will have access to local host directories

# Exercises local directory mounted to Cloudera FS (typically in the /media share)

# Hue: File Browser/Directory and Editor view



- When accessing local cluster via browser, you can upload file/data via File Browser in Hue from your local machine

# Use File Browser to Upload files

# Accessing CDH cluster from your browser



- Use SSH tunnel with convenient port mapping
- Benefit from accessing local file system and uploading files and data via browser

# Cloudera Quickstart with enabled Cloudera Manager



- To run Cloudera Manager Launch Cloudera Express from Desktop
- Be aware about requirement for configuring at least 2 virtual CPU and minimum memory 8 GB

# Configuring Tunnels for known Ports



- Using PuTTY SSH client
- First, configure dynamic tunnel to EMR cluster
  - Port 8157
  - Using Source port and Dynamic radio button
  - Example:
    - D5147
    - Source L8088 Destination Localhost:8088
    - Source L8888 Destination Localhost:8888

- Configure other local ports
  - See EMR ports table

- This will create a secure tunnel by forwarding a port (the "destination port") on the remote server to a port (the "source port") on the local host (127.0.0.1 or localhost).

# Configuring SSH tunnel on Windows

https://docs.bitnami.com/virtual-machine/faq/get-started/access-phpmyadmin/

- In order to access phpMyAdmin via SSH tunnel, you need an SSH client. In the instructions below we have selected PuTTY, a free SSH client for Windows and UNIX platforms.
  - The first step is to configure PuTTY.
- Once you have your SSH client correctly configured and you have confirmed that you can successfully access your instance using SSH, you need to create an SSH tunnel in order to access phpMyAdmin. Follow these steps:
  - In the "Connection -> SSH -> Tunnels" section, add a new forwarded port by introducing the following values:
    - Source port: 8888
    - Destination: localhost:8888
    - Remember that if you are redirecting HTTP requests to the HTTPS port, you must use destination port 443 instead of 80.
  - This will create a secure tunnel by forwarding a port (the "destination port") on the remote server to a port (the "source port") on the local host (127.0.0.1 or localhost).
  - Click the "Add" button to add the secure tunnel configuration to the session. You'll see the added port in the list of "Forwarded ports".
- PuTTY configuration
  - In the "Session" section, save your changes by clicking the "Save" button.
  - Click the "Open" button to open an SSH session to the server. The SSH session will now include a secure SSH tunnel between the two specified ports.
- Access the phpMyAdmin console through the secure SSH tunnel you created, by browsing to http://127.0.0.1:8888/phpmyadmin.
- Log in to phpMyAdmin by using the following credentials:
  - Username: root
  - Password: application password. (Refer to our FAQ to learn how to find your application credentials).

# Running MapReduce examples on CDH5.13

Follow step, adjust jar name and commands
https://docs.cloudera.com/documentation/other/tutorial/CDH5/topics/ht_usage.html#topic_5_2

- $ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar wordcount /user/cloudera/wordcount/output

- $ hadoop fs –rm –r  /user/cloudera/wordcount/output

- $ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 10 10000

- $ hadoop fs –rm –r  /user/cloudera/wordcount/output

- $hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar wordmean /user/cloudera/wordcount/output

# Commands/Programs in hadoop-mapreduce-examples.jar

Jobs in the examples JAR file hadoop-mapreduce-examples.jar

- **aggregatewordcount**: An Aggregate-based map/reduce program that counts the words in the input files.
- **aggregatewordhist:** An Aggregate-based map/reduce program that computes the histogram of the words in the input files.
- **bbp:** A map/reduce program that uses Bailey-Borwein-Plouffe to compute the exact digits of pi.
- **dbcount**: An example job that counts the pageview counts from a database.
- **distbbp**: A map/reduce program that uses a BBP-type formula to compute the exact bits of pi.
- **grep:** A map/reduce program that counts the matches to a regex in the input.
- **join:** A job that effects a join over sorted, equally partitioned data sets.
- **multifilewc**: A job that counts words from several files.
- **pentomino:** A map/reduce tile laying program to find solutions to pentomino problems.
- **pi:** A map/reduce program that estimates pi using a quasi-Monte Carlo method.
- **randomtextwriter:** A map/reduce program that writes 10 GB of random textual data per node.
- **randomwriter**: A map/reduce program that writes 10 GB of random data per node.
- **secondarysort**: An example defining a secondary sort to the reduce.
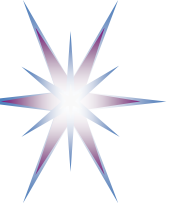- **sort:** A map/reduce program that sorts the data written by the random writer.
- **sudoku:** A Sudoku solver.
- **teragen:** Generate data for the terasort.
- **terasort:** Run the terasort.
- **teravalidate:** Check the results of the terasort.
- **wordcount:** A map/reduce program that counts the words in the input files.
- **wordmean:** A map/reduce program that counts the average length of the words in the input files.
- **wordmedian:** A map/reduce program that counts the median length of the words in the input files.
- **wordstandarddeviation:** A map/reduce program that counts the standard deviation of the length of the words in the input files.

# Do you need Reduce?

| ID | Username | Category | Amount |
|----|----------|----------|--------|
| 1 | Janani | Books | 200 |
| 2 | Swetha | Clothing | 450 |
| 3 | Shreya | Electronics | 300 |
| 4 | Jitu | Books | 700 |

## Which users spent >300?

## How many of them bought Books?

## Selecting a **specific** set of records from a dataset

# SQL Query can do it

| ID | Username | Category | Amount |
|----|----------|----------|--------|
| 1 | Janani | Books | 200 |
| 2 | Swetha | Clothing | 450 |
| 3 | Shreya | Electronics | 300 |
| 4 | Jitu | Books | 700 |

**Selecting a specific set of records**

**If this were a database table**

# An SQL query

```
select * from <table name>
where <condition>
```

| User | Followers |
|------|-----------|
| 1 | 30 |
| 2 | 30000 |
| 3 | 20 |
| 4 | 40 |
| 5 | 50 |
| 6 | 6000 |

## Get the top N records

**If this were a database table**

# An SQL query
### order by

```
select * from <table name>
where <condition>
order by <column name>
```