# Cloudera Distributed Hadoop (CDH) Installation and Configuration on Virtual Box

By

Kavya Mugadur

W1014808

# Table of contents
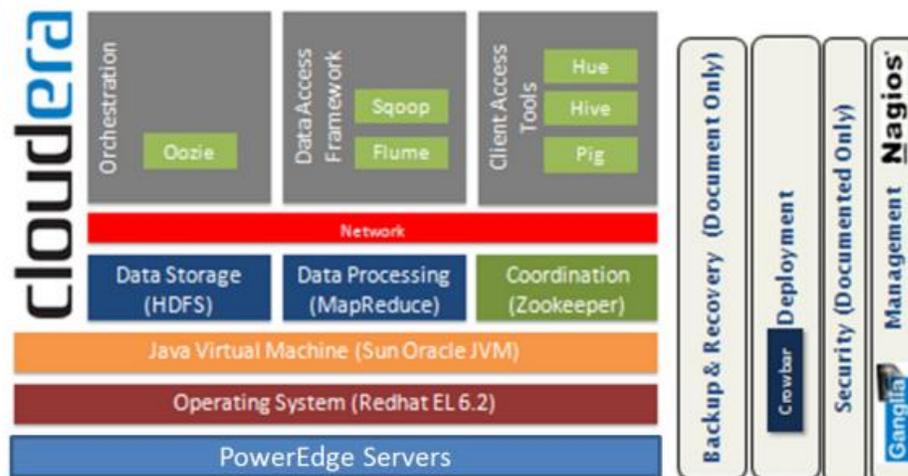
# 1.What is CDH ?

CDH (Cloudera Distribution Hadoop) is open-source Apache Hadoop distribution provided by Cloudera Inc which is a Palo Alto-based American enterprise software company.

CDH (Cloudera's Distribution Including Apache Hadoop) is the most complete, tested, and widely deployed distribution of Apache Hadoop. CDH is 100% open source and is the only Hadoop solution to offer batch processing, interactive SQL and interactive search as well as enterprise-grade continuous availability. More enterprises have downloaded CDH than all other distributions combined.

<u>CLOUDERA TAXONOMY</u>



The PowerEdge servers, the operating system, and the Java Virtual Machine make up the foundation on which the Hadoop software stack runs. The dark blue layer, depicting the core Hadoop components, comprises two frameworks:

• The Data Storage Framework is the file system that Hadoop uses to store data on the cluster nodes. Hadoop Distributed File System (HDFS) is a distributed, scalable, and portable file system.

• The Data Processing Framework (MapReduce) is a massively-parallel compute framework inspired by Google's MapReduce papers.

The next layer of the stack is the network layer. This is a dedicated cluster network, implemented from a blueprint using tested and qualified components. This implementation provides predictable high performance without interference from other applications.

The next three frameworks—the Orchestration, the Data Access Framework, and the Client Access Tools—are utilities that are part of the Hadoop ecosystem and provided by the CDH distribution.

3

## 2. Hadoop Basics

The Hadoop platform was designed to solve problems where you have a big data. It's for situations where you want to run analytics that are deep and computationally extensive, like clustering and targeting. The majority of this data will be "unstructured" – complex data poorly suited to management by structured storage systems like relational database.

Unstructured data comes from many sources and takes many forms  web logs, text files, sensor readings, user-generated content  like  product  reviews  or text messages, audio, video and still imagery and more

Dealing with big data requires two things:

> • Inexpensive, reliable storage; and

> • New tools for analyzing unstructured and structured data.

Apache Hadoop is a powerful open source software platform that addresses both of these problems. Hado op is an Apache Software Foundation project. **Cloudera** offers commercial support and services to Hadoop users.

### 2.1 Reliable Storage: HDFS

Hadoop includes a fault-tolerant storage system called the **Hadoop Distributed File System, or HDFS**. HDFS is able to store huge amounts of information, scale up incrementally and survive the failure of signi ficant parts of the storage infrastructure without losing data.

Hadoop creates clusters of machines and coordinates work among them. Clusters can be builtwith inexpe nsive computers. If one fails, Hadoop continues to operate the cluster without losing data or interrupting work, by shifting work to the remaining machines in the cluster.

HDFS manages storage on the cluster by breaking incoming files into pieces, called "blocks," and storing each of the blocks redundantly across the pool of servers.  In the common case, HDFS stores three comple te copies of each file by copying each piece to three different servers



Figure 1: HDFS distributes file blocks among servers

## 2.2 Hadoop for Big Data Analysis

Many popular tools for enterprise data management -relational database systems, for example –
are designed to make simple queries run quickly. They use techniques like indexing to examine just a small portion of all the available data in order to answer a question.

Hadoop is a different sort of tool. Hadoop is aimed at problems that require examination of all the available data. For example, text analysis and image processing generally require that every single record be read, and often interpreted in the context of similar records. Hadoop uses a technique called **MapReduce** to carry out this exhaustive analysis quickly.

In the previous section, we saw that HDFS distributes blocks from a single file among a large number of servers for reliability. Hadoop takes advantage of this data distribution by pushing the work involved in an analysis out to many different servers. Each of the servers runs the analysis on its own block from the file. Results are collated and digested into a single result after each piece has been analyzed.

Running the analysis on the nodes that actually store the data delivers much much better performance than reading data over the network from a single centralized server. Hadoop monitors jobs during execution, and will restart work lost due to node failure if necessary. In fact, if a particular node is running very slowly, Hadoop will restart its work on another server with a copy of the data.



**Figure 2: Hadoop pushes work out to the data**

Hadoop's MapReduce and HDFS use simple, robust techniques on inexpensive computer systems to deliver every high data availability and to analyze enormous amounts of information quickly. Hadoop offers enterprises a powerful new tool for managing big data.

5

## 3.Ways to Install CDH4

You can install CDH4 in any of the following ways:

    a.   Installing using Cloudera quickstart vm.
    b.   Automated method using Cloudera Manager. Cloudera Manager Free Edition automates the installation and configuration of CDH4 on an entire cluster if you have root or password-less sudo SSH access to your cluster's machines.
    c.   Manual methods described below:
- Download and install the CDH4 "1-click Install" package
- Add the CDH4 repository
- Build your own CDH4 repository
- Install from a CDH4 tarball

In this document we will explain the installation of CDH using first method.

## 4. Installation and Configuration of CDH on Virtual machine using Cloudera quickstart vm

Cloudera quickstart VM contains a sample of Cloudera's platform for "Big Data". The VM from Cloudera is available in VMware, VirtualBox and KVM flavors, and all require a 64 bit host OS. This VM runs CentOS 6.2 and includes CDH4.3, Cloudera Manager 4.6, Cloudera Impala 1.0.1 and Cloudera Search .9 Beta

In this document we have installed CDH on VirtualBox. Below are the steps to install CDH using Cloudera quickstart vm on Virtual box

Step 1:  Download the virtual box- executable file from https://www.virtualbox.org/wiki/Downloads
Download VirtualBox 4.2.16 for Windows hosts

Step 2: Install VirtualBox by double clicking on the downloaded file.

Click install to install VirtualBox with default settings. Installation is shown as below:

If the installation is successful you will see a Virtual Manager window to manage VMs.



Step 3: Download the Cloudera quickstart vm for VirtualBox

Go to the link - https://ccp.cloudera.com/display/SUPPORT/Cloudera+QuickStart+VM

Select quickVM for VirtualBox and click on download



Step 4: Unzip the downloaded file.

When you unzip the file cloudera-quickstart-vm-4.3.0-virtualbox.tar you will find these two files in the directory.



Step 5: Open VirualBox and click on "New" to create new virtual box



Give name for new virtual machine and select type as Linux and versions as Linux 2.6

11

Step 6 : Select Memory Size as 4GB and click Next.



Step 7: In the next page, VirtualBox asks to select Hard Drive for new VirualBox as shown in the screenshot. Create a virtual hard drive now is selected by default. But you have to select "Use an existing virtual hard drive file" option.

Select "Use an existing virtual hard drive file".



Step 8:  Click on the small yellow icon beside the dropdown to browse and select the cloudera-quickstart-vm-4.3.0-virtualbox-disk1.vm file (which is download in step 4).

Click on create to create Cloudera quickstart vm.



Step 9 : Your virtual box should look like following screen shots. We can see the new virtual machine named Cloudera Hadoop on the left side.



14

Step 10: Select Cloudera vm and click on "Start"    ⟹

Virtual Machine starts to boot

```
Entering non-interactive startup
Calling the system activity data collector (sadc):
Bringing up loopback interface:                              [  OK  ]
Starting auditd:                                             [  OK  ]
Starting portreserve:                                        [  OK  ]
Starting system logger:                                      [  OK  ]
Starting irqbalance:                                         [  OK  ]
Starting kdump:                                              [FAILED]
Starting system message bus:                                 [  OK  ]
Setting network parameters...                                [  OK  ]
Starting NetworkManager daemon:                              [  OK  ]
Starting Avahi daemon...                                     [  OK  ]
Starting cups: e1000: eth1 NIC Link is Up 1000 Mbps Full Duplex, Flow Control: R
X
                                                             [  OK  ]
Mounting other filesystems:  mount: sysfs already mounted or /sys busy
mount: according to mtab, /sys is already mounted on /sys
                                                             [FAILED]
Starting acpi daemon:                                        [  OK  ]
Starting HAL daemon:                                         [  OK  ]
Retrigger failed udev events                                 [  OK  ]
Adding udev persistent rules                                 [  OK  ]
Starting sshd:                                               [  OK  ]
hrtimer: interrupt took 3245317 ns
_
```

```
Starting system message bus:                                 [  OK  ]
Setting network parameters...                                [  OK  ]
Starting NetworkManager daemon:                              [  OK  ]
Starting Avahi daemon...                                     [  OK  ]
Starting cups: e1000: eth1 NIC Link is Up 1000 Mbps Full Duplex, Flow Control: R
X
                                                             [  OK  ]
Mounting other filesystems:  mount: sysfs already mounted or /sys busy
mount: according to mtab, /sys is already mounted on /sys
                                                             [FAILED]
Starting acpi daemon:                                        [  OK  ]
Starting HAL daemon:                                         [  OK  ]
Retrigger failed udev events                                 [  OK  ]
Adding udev persistent rules                                 [  OK  ]
Starting sshd:                                               [  OK  ]
hrtimer: interrupt took 3245317 ns
Starting mysqld:                                             [  OK  ]
DB initialization done.
pg_ctl: another server might be running; trying to start server anyway
waiting for server to start........ done
server started
Starting postfix:                                            [  OK  ]
Starting abrt daemon:                                        [  OK  ]
Starting cloudera-scm-agent:                                 [  OK  ]
Starting cloudera-scm-server: _
```

15

Step 11: System is loaded and CDH is installed on virtual machine.



Step 12: System redirects you to the index page of Cloudera.

Step 13: Select Cloudera Manager and Agree to the information assurance policy.



Step 14: Login to Cloudera Manager as admin. Password is admin.



17

Step 15:  We can see all the services running on our single node cluster.



Step 16:  Click on the Hosts tab and we can see that one host is running , version of CDH installed on it is 4 , health  of the host is good and last heart beat was listened 5.22s ago.

Step 17 : Click on the localhost.localdomain to se the detail information about the host



Step 18 : We can also chang the password for admin by selecting the adminstration tab and clicking on "Change Password" button.

## 5. Running MapReduce Program

Step 1: Write a MapReduce program. We have used the word count program for testing the CDH installation.

Write a java program using eclipse.



Step 2: Add external jars to compile the code. Right click on the project and select properties. Select the Java build path and then click on "Add External Jars". Select all the jars present in folder /usr/bin/hadoop/client

Step 3: Click "OK" o add those jars to your program.



Step 4 : Create a jar file of your program.

Right click on project select "Export" and then click on "Jar File" under Java folder.

Give the location path where you want to store your .jar file.



Now we can see the jar file in the given location.

Step 5: Input file

Open the terminal and create a input file which is a huge text file.

$vim input.txt

Step 6 : Check the input.txt. Execute command $ ls –ltr. The highlighted file is our input file.



Step 7 : Make a new file directory on HDFS (Hadoop Distributed File System)

$ sudo su hdfs

hadoop fs –mkdir /WordCountFiles

hadoop fs –ls /

Step 8 :Copy this file on the NameNode i.e., on HDFS

$ hdfs dfs –copyFromLocal input.txt /WordCountFiles



Step 9: Run the program using the hadoop command

$ hadoop jar ./WordCount.jar WordCount /WordCountFiles/input.txt /WordCountFiles/output

Where ./WordCount.jar is the path and name of the jar file we created in Step 4 and WordCount is the name of the program.

MapReduce program starts to run.We can see the percentage of mapping and reducing the program is doing on the command line.

We can see that the map and reduce percentage gradually increasing, this shows that the program is successfully running on CDH and using its MapReduce technique to count the frequency of each words.

When the program runs we can see on the command line the number of input bytes the program has read and number of tasks launched and other useful information on the command line.

Step 10 : Check the output

When the program runs successfully output directory is created. In our case the output directory name is "output" (as mentioned in the command )

$ hadoop fs –ls /WordCountFiles

$ hadoop fs –ls /WordCountFiles/ouput

We can see that there are three files/directory in the output directory. Our output is present in part-0000 file



Step 11 :  Display the output

$hadoop fs –cat /WordCountFiles/ouput/part-0000

Output file shows the word and the number of times it has occurred in the file. For example , word "Mirror" has occurred 22 times in the given input file.

```
Applications  Places  System                              Fri Jul 19, 9:29 PM   cloudera

                          cloudera@localhost:~                          _  □  X

File  Edit  View  Search  Terminal  Help

Mirror  22
Next    1
Not     1
Now     2
Prince  2
Queen   107
Queen,  36
Queen.  12
Seven   22
She     24
Since   11
Snow    176
Soon    11
The     117
Then    11
They    22
This    12
Time    1
When    11
White   124
White's 5
White,  33
White.  14
a       181
about   22
accept  2
after,  22
again   1
again,  1
again:  11
against 11
age     11
alive,  11
alive.  11
all     11
alone   11
also    11
always  22
an      22

[localhost.local...] [cloudera@loca...] [Snow White - ...] [cloudera@loc...] [Resource - Wo...] [eclipse]
```

# 6. Configuring hadoop in multi-tasking mode (Multi -Thread)

Step1 : Go to directory where task tracker's map reduce configuration file is found.

$ cd /var/run/Cloudera-scm-agent/process/11-mapreduce-TASKTRACKER



Step 2: Open the file mapred-site.xml for editing.

Default file looks as follows:

Step 3 : Change the Configuration.

Change mapreduce.tasktracker.map.tasks.maximum and mapreduce.tasktracker.reduce.tasks.maximum to 4. These parameters determine the maximum number of map/reduce tasks that will be run by task tracker in pseudo distributed mode.



Step 4 : The above params are just a hint for the tracker. If you want to enforce the number then add new properties mapred.map.tasks and mapred.reduce.tasks and set them to the desired value.

Step 5 : Restart the hadoop daemons

$ bin/hadoop-daemon.sh stop cloudera-scm-server

$ bin/hadoop-daemon.sh start cloudera-scm-server

$ bin/hadoop-daemon.sh stop cloudera-scm-agent

$ bin/hadoop-daemon.sh start cloudera-scm-agent

Now the new configurations will load and when the huge file data is given, the number of task to map and reduce will be increased.

# 7. Configuring Flume (Monitoring Configuration)

Apache Flume is a distributed, reliable, and available system for efficiently collecting, aggregating and moving large amounts of log data from many different sources to a centralized data store.

The use of Apache Flume is not only restricted to log data aggregation. Since data sources are customizable, Flume can be used to transport massive quantities of event data including but not limited to network traffic data, social-media-generated data, email messages and pretty much any data source possible.

Apache Flume is a top level project at the Apache Software Foundation.

Step 1 : Check Flume installation. Flume is installed as a part of quickstart VM.

$ rpm –qa | grep –i flume-ng

This displays the flume installation files

Step 3 :Check the flume template files under /etc/flume-ng/conf



```
[root@localhost conf]# pwd
/etc/flume-ng/conf
[root@localhost conf]# ls -ltr
total 12
-rw-r--r--. 1 root root 1197 Oct 29  2012 flume-env.sh.template
-rw-r--r--. 1 root root 1661 Oct 29  2012 flume-conf.properties.template
-rw-r--r--. 1 root root 2156 May 28 00:24 log4j.properties
-rw-r--r--. 1 root root    0 May 28 00:24 flume.conf
[root@localhost conf]# 
```

Step 4: Flume default configuration file.

$ vim flume.conf



```
[root@localhost conf]# cat flume-conf.properties.template
# Licensed to the Apache Software Foundation (ASF) under one
# or more contributor license agreements.  See the NOTICE file
# distributed with this work for additional information
# regarding copyright ownership.  The ASF licenses this file
# to you under the Apache License, Version 2.0 (the
# "License"); you may not use this file except in compliance
# with the License.  You may obtain a copy of the License at
#
#  http://www.apache.org/licenses/LICENSE-2.0
#
# Unless required by applicable law or agreed to in writing,
# software distributed under the License is distributed on an
# "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY
# KIND, either express or implied.  See the License for the
# specific language governing permissions and limitations
# under the License.


# The configuration file needs to define the sources,
# the channels and the sinks.
# Sources, channels and sinks are defined per agent,
# in this case called 'agent'

agent.sources = seqGenSrc
agent.channels = memoryChannel
agent.sinks = loggerSink

# For each one of the sources, the type is defined
agent.sources.seqGenSrc.type = seq

# The channel can be defined as follows.
agent.sources.seqGenSrc.channels = memoryChannel

# Each sink's type must be defined
agent.sinks.loggerSink.type = logger
```

34

Step 5 :Flume command and its options. We will be using the agent command.



Step 6:Changing flume.conf file

This configuration defines a single agent named a1. a1 has a source that listens for data on port 44444, a channel that buffers event data in memory, and a sink that logs event data to the console. The configuration file names the various components, then describes their types and configuration parameters. A given configuration file might define several named agents; when a given Flume process is launched a flag is passed telling it which named agent to manifest.

Step 7:Starting of a flume agent

An agent is started using a shell script called flume-ng which is located in the bin directory of the Flume distribution.

$ bin/flume-ng agent –conf conf –conf-file ./flume.conf –name a1 –Dflume.root.logger=INFO,console

Step 8 : Testing the new Configuration

Open the new terminal and then we can then telnet port 44444 and send Flume an event. The original
Flume terminal will output the event in a log message.

## 8. Running Pig project

**Apache Pig** is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs. The salient property of Pig programs is that their structure is amenable to substantial parallelization, which in turns enables them to handle very large data sets.

At the present time, Pig's infrastructure layer consists of a compiler that produces sequences of Map-Reduce programs, for which large-scale parallel implementations already exist . Pig's language layer currently consists of a textual language called Pig Latin which can be used to write queries.

Step 1: Checking the pig installation. Pig is installed as a part of quickstart vm

$rpm –qa |grep –i pig



Step 2: Open Pig Shell which display the grunt prompt. This is used to run user queries in high level query language called Pig Latin

$ pig –x local

Step 3:Running the sample pig latin query: This greps for a particular pattern in a input file.

Step 4 : Displays the success message on successful completion of query.

# 9. MAHOUT installation

The goal of Mahout is to build a vibrant, responsive, diverse community to facilitate discussions not only on the project itself but also on potential use cases.  Mahout comes installed with quickstart vm. There is no additional configuration required.
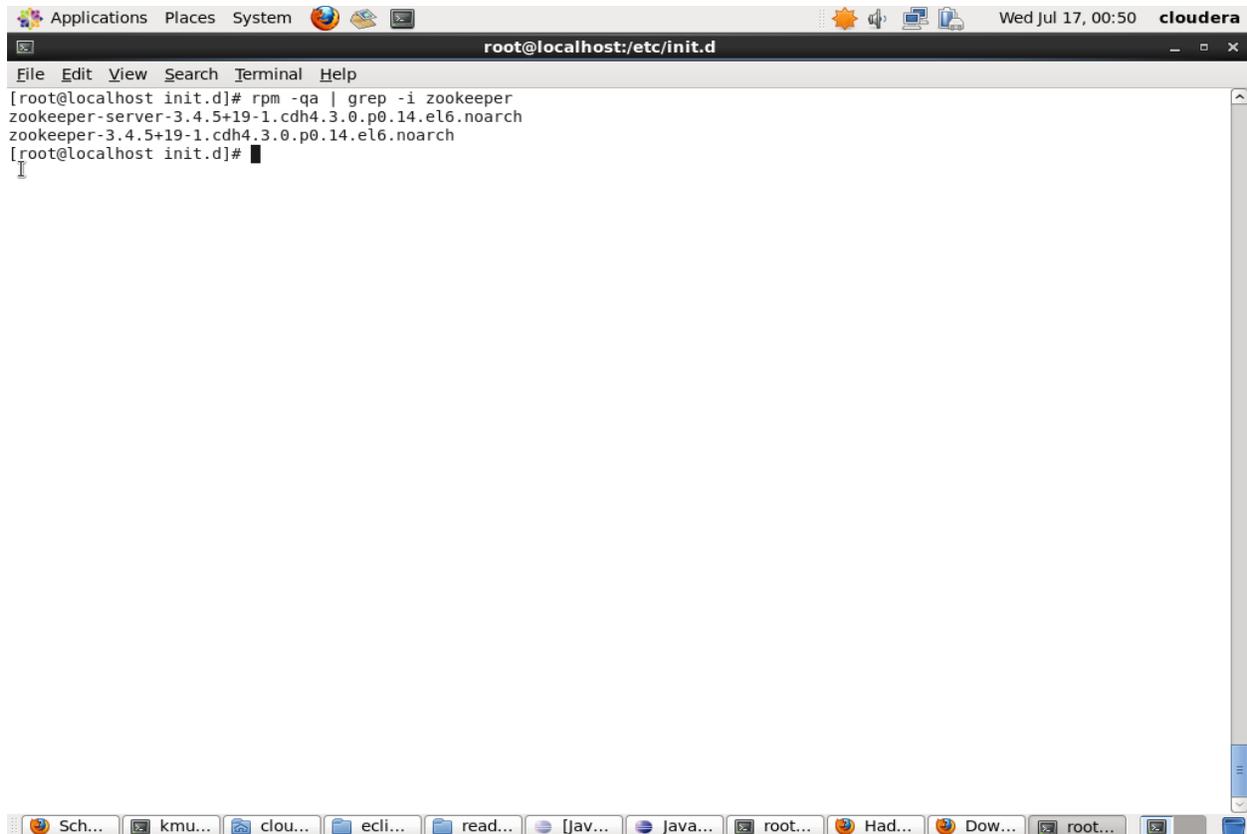
# 10. Zookeeper

ZooKeeper is a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services. All of these kinds of services are used in some form or another by distributed applications. ZooKeeper allows distributed processes to coordinate with each other through a shared hierarchal namespace which is organized similarly to a standard file system. All objects in this namespace are represented by znode which has data and attributes associated with it. We can perform all basic file operations like create, delete, modify etc on this node.

Step 1: Check for zookeeper installation. Zookeeper is installed as part of Cloudera quickstart vm



Step 2 : Start zookeeper

$service zookeeper.server start

Step 3: Creating a z-node and setting some data in that node

**References:**

www.cloudera.com

http://hadoop.apache.org/

http://developer.yahoo.com/hadoop/tutorial/

44