



A BLENDED TECHNOLOGY INITIATIVE OF INCORPORATING BIG DATA AND CLOUD – HD INSIGHT CLUSTER CREATION AND SETUP USING MICROSOFT AZURE PLATFORM WITH BIG DATA HADOOP FRAMEWORK

Mr.D.R.Sudharsan– III CSE / Knowledge Institute of Technology, Salem
Mr.E.Arun – IIICSE / Knowledge Institute of Technology, Salem
Mr.R.J.Vigneswaran – IIICSE / Knowledge Institute of Technology, Salem
Mr.K.Sabari – IIICSE / Knowledge Institute of Technology, Salem
Mr.N.Vasanth- IIICSE / Knowledge Institute of Technology, Salem
Mr.S.S.Aravinth – ASP/ CSE / Knowledge Institute of Technology, Salem

ABSTRACT:

Big data is one of the emerging technology in this current scenario. In order to handle the data in big data a well-defined lab environment is mandatory. That lab environment should consists of a huge amount of commodity hard wares to form Name node, Secondary name node and set of data nodes. Even after this the next tedious part of managing the big data is configuration of hadoop cluster and map reduce strategy. Whenever it's needed to monitor the hadoop cluster the user have to access that particular lab environment. These all issues can be overcome by the solution called HD INSIGHTS (Hadoop on azure)

a blended technological platform provided by the Microsoft .By using this the whole lab environment and hadoop cluster can be setup in cloud without need of any configuration in linux and the user can access the particular hadoop cluster via RDP (Remote Desktop Protocol) with highly secured and authenticated fashion from word wide.

I) INTRODUCTION:

1. Cloud Computing

Cloud computing is the term which defines that any of the service that can be provided via internet, that service may be a storage or a platform or a



infrastructure or anything else. Therefore these services are broadly classified into three types namely

Infrastructure as a service

Platform as a service

Software as a service

2. Big Data

The term big data is used to represent the collection of complex and large sets of structured, unstructured and semi structured data. Such that it's very difficult to capture, process, store, search and analyze these kinds of data by using the conventional data bases management tools and traditional database management systems.

3. Intention of Blended Technology

Blended technology is a influence of different domains in which consists of a combination of one or more technology to give the solution for a real world issue. Lets think about IBM Bluemix it is a blended technology with combination of web and cloud where it provides a sufficient platform to develop and deploy the user's web page. It is not needed to install any local web servers and text editors in a machine. The only need is an internet connectivity

4. Microsoft Azure

Azure is the cloud computing platform provided by the Microsoft. Especially azure is focused on platform as a service. Because of that it provides sufficient platform for emerging technologies like big data, web

of things, machine learning etc., In domain of analytics it offers lot of services like

HD insight, Data factory etc.

5. Hadoop

Hadoop is the open source java framework managed by apache software foundation has proven to be very helpful in storing and managing vast amount of complex data like structured, unstructured and semi structured data cheaply and efficiently. Which allows you to store large amount of data of a cluster in a low cost commodity hardware. It is very important to understand that hadoop is not a database it is a distributed file system (HDFS-hadoop distributed file system) which distributes the data over number of data nodes and also provides you capability to process that distributed data using a simple programming model.

II) PHYSICAL HADOOP CLUSTER CREATION AND SETUP:

In order to create a physical hadoop cluster hardware is one of the main requirement and it's required to get a clear idea on linux machine working and there is lot of configuring procedures as shown below.

Step-1

Download the JDK setup file for linux

Make a directory in following location /usr/lib/jvm by using command

mkdir /usr/lib/jvm

step2

Move the downloaded JDK file into the JVM directory created in /usr/lib/jvm



Step-3

Install java by using the following command

Update-alternatives /usr/bin/java java

/usr/lib/jvm/java

Update-alternatives /usr/bin/javac java

/usr/lib/jvm/javac

Step-4

Open .bashrc file by using

nano .bashrc

Setup the java home path in .bashrc file

Export JAVA_HOME=/usr/lib/jvm/jdk-1.8.0

Export PATH=\$PATH:\$JAVA_HOME/bin

```
suz@ubuntu: ~  
File Edit View Search Terminal Help  
GNU nano 2.2.4 File: .bashrc Modified  
  
if [ -f /etc/bash_completion ] && ! shopt -oq posix; then  
  . /etc/bash_completion  
fi  
  
export JAVA_HOME=/usr/lib/jvm/jdk1.8.0_45  
export PATH=$PATH:$JAVA_HOME/bin  
[  
  
^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^X Cut Text ^C Cur Pos  
^X Exit ^J Justify ^W Where Is ^V Next Page ^U UnCut Text ^T To Spell
```

Step-5

Verify whether java had installed or not by the command

Java -version

```
suz@ubuntu: ~  
File Edit View Search Terminal Help  
root@ubuntu:/home# nano .bashrc  
No command 'nano' found, did you mean:  
Command 'nano' from package 'nano' (main)  
Command 'nona' from package 'hugin-tools' (universe)  
Command 'mono' from package 'mono-runtime' (main)  
nano: command not found  
root@ubuntu:/home# nano .bashrc  
root@ubuntu:/home# su - suz  
suz@ubuntu:~$ nano .bashrc  
suz@ubuntu:~$ java -version  
java version "1.8.0_45"  
Java(TM) SE Runtime Environment (build 1.8.0_45-b14)  
Java HotSpot(TM) Client VM (build 25.45-b02, mixed mode)  
suz@ubuntu:~$
```

Step-6

Install SSH client and server by using the command

sudo apt-get install openssh-client openssh-server

step-7

Download the hadoop setup and the place it to any of the directory you want

Make the following editing which shown in a picture given below in files

Conf/core-site.xml



```
suz@ubuntu: ~/Desktop/hadoop-1.0.3/conf
File Edit View Search Terminal Help
GNU nano 2.2.4 File: core-site.xml

<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>fs.default.name</name>
<value>hdfs://localhost:50000</value>
</property>
</configuration>

[ Read 11 lines ]
^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^K Cut Text ^C Cur Pos
^X Exit ^J Justify ^W Where Is ^V Next Page ^U UnCut Text ^T To Spell
```

```
suz@ubuntu: ~/Desktop/hadoop-1.0.3/conf
File Edit View Search Terminal Help
GNU nano 2.2.4 File: hdfs-site.xml

<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>dfs.name.dir</name>
<value>/home/ip/hadoop/hadoop-1.1.2-name-dir</value>
</property>
<property>
<name>dfs.data.dir</name>
<value>/home/suz/home/Desktop/hadoop-1.0.3</value>
</property>
<property>
<name>dfs.permissions</name>
<value>>false</value>
</property>
</configuration>

[ Read 24 lines ]
^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^K Cut Text ^C Cur Pos
^X Exit ^J Justify ^W Where Is ^V Next Page ^U UnCut Text ^T To Spell
```

Conf/mapred-site.xml

```
suz@ubuntu: ~/Desktop/hadoop-1.0.3/conf
File Edit View Search Terminal Help
GNU nano 2.2.4 File: mapred-site.xml

<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>mapred.job.tracker</name>
<value>hdfs://localhost:50001</value>
</property>
</configuration>

^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^K Cut Text ^C Cur Pos
^X Exit ^J Justify ^W Where Is ^V Next Page ^U UnCut Text ^T To Spell
```

Conf.hdfs-site.xml

Step-8

Setup the hadoop environment path in a file
conf/hadoop-env.sh

```
suz@ubuntu: ~/Desktop/hadoop-1.0.3/conf
File Edit View Search Terminal Help
GNU nano 2.2.4 File: hadoop-env.sh Modified

# Set Hadoop-specific environment variables here.

# The only required environment variable is JAVA_HOME. All others are
# optional. When running a distributed configuration it is best to
# set JAVA_HOME in this file, so that it is correctly defined on
# remote nodes.

# The java implementation to use. Required.
export JAVA_HOME=/usr/lib/jvm/jdk1.8.0_45

# Extra Java CLASSPATH elements. Optional.
# export HADOOP_CLASSPATH="<extra_entries:$HADOOP_CLASSPATH>"

# The maximum amount of heap to use, in MB. Default is 1000.
# export HADOOP_HEAPSIZE=2000

# Extra Java runtime options. Empty by default.
# export HADOOP_OPTS="server"

[ Read 54 lines ]
^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^K Cut Text ^C Cur Pos
^X Exit ^J Justify ^W Where Is ^V Next Page ^U UnCut Text ^T To Spell
```

Step-9

Setup hadoop home directory path in .bashrc
As shown below

While this process it asks you for the ssh password key.

Step-11

Check the java process status by the command

Step-10

Start formatting the namenode by the following command

Hadoop namenode -format (inside the bin directory)

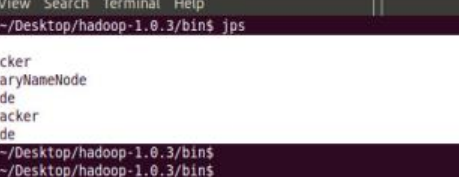
```

suz@ubuntu: ~/Desktop/hadoop-1.0.3/bin
File Edit View Search Terminal Help

suz@ubuntu:~/Desktop/hadoop-1.0.3/bin$ ./hadoop namenode -format
Warning: SHADOOP_HOME is deprecated.

15/07/27 20:54:38 INFO namenode.NameNode: STARTUP MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = ubuntu/127.0.1.1
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 1.0.3
STARTUP_MSG: build = https://svn.apache.org/repos/asf/hadoop/common/branches/
ranch-1.0 -r 1335192; compiled by 'hortonfo' on Tue May 8 20:31:25 UTC 2012
*****/
15/07/27 20:54:40 INFO util.GSet: VM type = 32-bit
15/07/27 20:54:40 INFO util.GSet: 2% max memory = 19.33375 MB
15/07/27 20:54:40 INFO util.GSet: capacity = 2^22 = 4194304 entries
15/07/27 20:54:40 INFO util.GSet: recommended=4194304, actual=4194304
15/07/27 20:55:02 INFO namenode.FSNamesystem: fsOwner=suz
15/07/27 20:55:02 INFO namenode.FSNamesystem: supergroup=supergroup
15/07/27 20:55:02 INFO namenode.FSNamesystem: isPermissionEnabled=false
15/07/27 20:55:02 INFO namenode.FSNamesystem: dfs.block.invalidate.limit=100
15/07/27 20:55:02 INFO namenode.FSNamesystem: isAccessTokenEnabled=false access
eyUpdateInterval=0 min(s), accessTokenLifetime=0 min(s)
15/07/27 20:55:02 INFO namenode.NameNode: Caching file names occurring more than
10 times

```



The screenshot shows a terminal window with the title bar "suz@ubuntu: ~/Desktop/hadoop-1.0.3/bin". The menu bar includes "File", "Edit", "View", "Search", "Terminal", and "Help". The command prompt is "suz@ubuntu:~/Desktop/hadoop-1.0.3/bin\$". The user has entered the command "jps", and the output is displayed as follows:

```
5829 Jps
4519 JobTracker
4137 SecondaryNameNode
5546 NameNode
4686 TaskTracker
3967 DataNode
```

Below the output, the user has entered the command "jps" multiple times, but no new output is shown, suggesting the command is being executed repeatedly without waiting for the previous output to scroll off the screen.

III) AZURE HD INSIGHT CLUSTER

CREATION:

These kinds of hadoop installation by using linux commands will take much time to configure and also the process of importing and exporting data is also a problematic one

Instead of working with these complicated configuration for hadoop cluster, the HD INSIGHT a service provided Microsoft azure is used to create our hadoop cluster with in a few minutes.

Step-1

Login into the www.manage.azure.com web portal by using the Microsoft account

Step-2

Step-10

Start all the daemons by the command **./start-all.sh** (inside the bin directory)



Then in the next page all the service provided by the Microsoft will be displayed like cloud, web apps, virtual machines etc.,
Among that select the STORAGE
After that create your storage account.

URL
newstorage13 ✓
*.core.windows.net

LOCATION/AFFINITY GROUP
East US

REPLICATION
Geo-Redundant

Step-3

Select option hd insight and setup the hadoop cluster in the created storage account with number of data nodes and name node.
it is shows in the next page.

CLUSTER NAME
newcluster ✓
*.azurehdinsight.net

CLUSTER SIZE ?
4 data nodes

HTTP USER NAME: admin
CONFIRM PASSWORD
.....

STORAGE ACCOUNT ?
newstorage13

OPERATING SYSTEM: Windows Server 2012 R2 Datacenter

CREATE HDINSIGHT CLUSTER ✓

Step-4

In next page select configuration option at top of the page and switch on the hadoop services



newcluster

[DASHBOARD](#) [MONITOR](#) [CONFIGURATION](#) [SCALE](#)

cluster connectivity

HADOOP SERVICES

ON

OFF

then in bottom select enable remote option the Remote Desktop Protocol will be downloaded.

Step-5

The it will ask you to setup then user name and password for RDP.

CONFIGURE HDINSIGHT

Configure Remote Desktop

RDP USER NAME

hadoop

PASSWORD

.....

CONFIRM PASSWORD

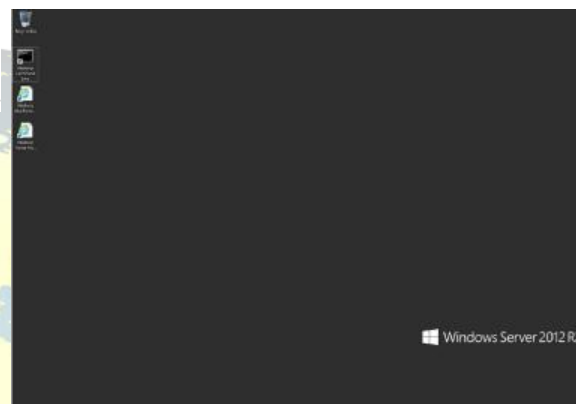
.....

EXPIRES ON ?

2015-07-31

Step-6

When the user opens the downloaded RDP it will ask you for the password created by you. After the verification of password then it opens up the virtual machine with the operating system of windows 2012. Here you are allotted with the memory space of 1000 GB only for your hadoop cluster in a cloud.



NAME NODE, MAP REDUCE, JOB TRACKER STATUS

Name node status.

The Hadoop data storage layer is the Hadoop Distributed File System (HDFS), and HDFS can be accessed from the management portal, the Hadoop command line, or from the Hadoop web service at port 30070 by the URL

<http://headnodehost:30070/dfshealth.jsp>



http://headnodehost:30070/dfshealth.html#tab-overv

Namenode information

Hadoop

Overview

Datanodes

Snapshot

Startup Progress

Utilities

jobtrackerhost Hadoop Map/Reduce Administration

State: RUNNING
Started: Wed May 07 17:16:01 GMT 2014
Version: 1.2.0.1.3.7.1-01293, 180cb428e0d99e348c9f5493440aa205bd9694972
Compiled: Thu Apr 03 18:39:04 Pacific Daylight Time 2014 by jenkins
Identifier: 201405071715
SafeMode: OFF

Overview 'headnodehost:9000' (active)

Cluster Summary (Heap Size is 239.19 MB/3.56 GB)

Started:	Tue Jul 28 13:00:26 GMT 2015	Running Map Tasks	Running Reduce Tasks	Total Submissions	Nodes	Occupied Map Slots	Occupied Reduce Slots	Reserved Map Slots	Reserved Reduce Slots	Map Task Capacity	Reduce Task Capacity
Version:	2.4.0.2.1.15.1-0001, rd6f928e863d5d5fb8510e2553a7d80f29e	14	1	6	4	14	1	0	0	16	8
Compiled:	2015-06-09T20:27Z by jenkins from (detached from d6f928e)										
Cluster ID:	CID-47d52a1b-6e44-47f3-b1d1-0f21245f41cd										
Block Pool ID:	BP-2058924255-100.112.68.138-1438088417327										

Job tracker status

it is a daemon which distributes the data towards the data node and also it stores the path that where a particular data is stored in data node.

Task Tracker

Task tracker is the daemon which perform the process in a part of data node and return the result to the job tracker. In this the scheduler information (queue status) can be known by the following link

Map reduce status

The Hadoop data processing layer has two components, JobTracker and TaskTracker. JobTracker usually runs on the same server as the name node. However, if needed, it can also run on a different machine. The role of JobTracker is to accept the jobs submitted to an HDInsight cluster from any source, break the job into tasks, and assign those tasks to the data nodes in the cluster. TaskTracker, which runs on each data node, sends status messages to JobTracker as it works on each job. To see the status of any job running on an HDInsight cluster, use the Hadoop web server on port 50030, following the http://your_hadoop_server:50030/ URL template http://your_hadoop_server:50030/, or the Hadoop MapReduce Status shortcut.

root

default

joblauncher

0.0% used

0.0% used

0.0% used

'joblauncher' Queue Status

Queue State: RUNNING
Used Capacity: 0.0%
Absolute Used Capacity: 0.0%
Absolute Capacity: 5.0%
Absolute Max Capacity: 50.0%
Used Resources: <memory:0, vCores:0>
Num Schedulable Applications: 0
Num Non-Schedulable Applications: 0
Num Containers: 0
Max Applications: 500
Max Applications Per User: 5000
Max Schedulable Applications: 6
Max Schedulable Applications Per User: 10
Configured Capacity: 5.0%
Configured Max Capacity: 50.0%
Configured Minimum User Limit Percent: 100%
Configured User Limit Exceeded: 10.0%



almost all of the hadoop operations can be possible in
HD INSIGHTS

Hadoop command line & windows power shell

When the user opens the command prompt in windows 2012 server a virtual machine which is opened via RDP, it will display the prompt called hadoop command line. Windows power shell is the shell which support all hadoop operations by this HD insight cluster all the hadoop operations can be possible in a very easy manner

IV) BENEFITS OF PHYSICAL TO CLOUD

When compared to physical machine the HD insight cluster is cost wise efficient because it doesn't need commodity hardware that required memory is allocated from cloud in HD insight hadoop cluster and one of the main advantage is that it avoids the usage of hardware and tedious linux commands for hadoop configuration. In a physical hadoop cluster whenever the user needs to monitor the data then each and every time he need to access the same physical machine but user can access HD insight hadoop cluster from anywhere all over the world by using RDP the only need is internet access.

V) CONCLUSION:

This Paper gives a research on blended technology in a combination of big data and cloud called HD INSIGHTS (HADOOP OF AZURE) the service provided by Microsoft and gives the brief and elaborative knowledge on it, and it is known fact that

References:

- <https://azure.microsoft.com/en-in/documentation/articles/hdinsight-hadoop-tutorial-get-started-windows/>
- <https://azure.microsoft.com/en-in/services/hdinsight/>
- <https://channel9.msdn.com/Series/Getting-started-with-Windows-Azure-HDInsight-Service>