

LAMBTON COLLEGE



A Project on  
[Hadoop Ecosystem]

121 Brunel Rd, Mississauga  
ON L4Z 3E9

A Group project for analyzing stack overflow data in Hive on Top of Hadoop  
Big Data Analytics DSMM

Under the supervision  
of  
Professor Teresa Zhu

**Submitted BY:Group E**

Aadarsha Chapagain (C0825975)  
Onyinye Mbanefo (C0831578)  
Roshan Acharya (C0831342)  
Anjana Kuriakose (C0829580)

**Submitted To:**

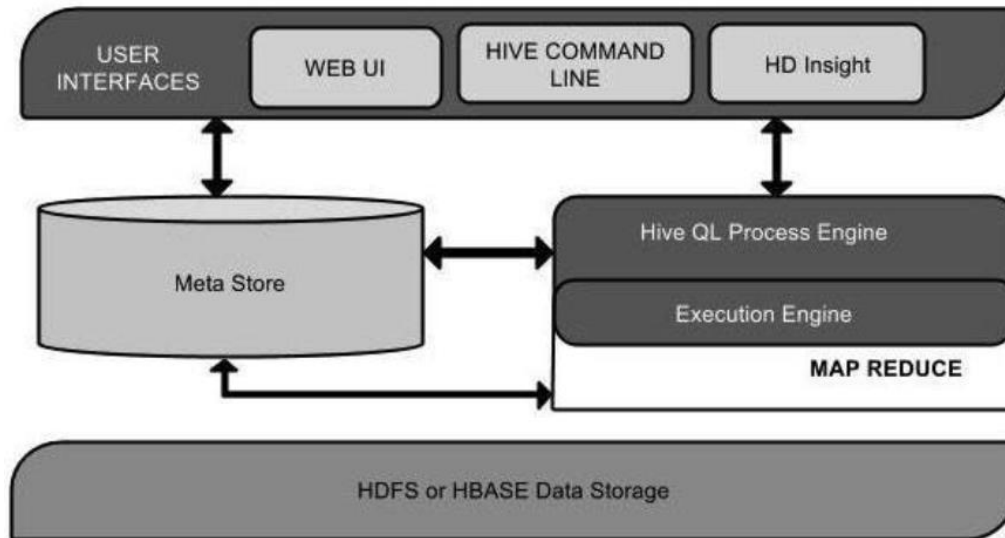
Lambton College  
Professor Teresa Zhu

**Submission Date:**  
27<sup>th</sup> September 2022

## Step1:

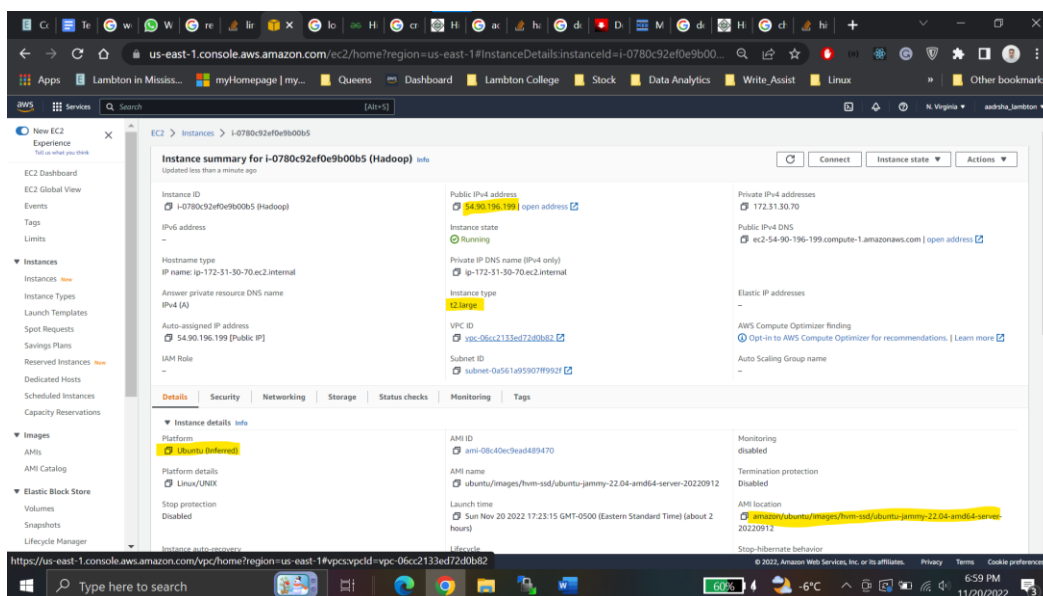
**1. Download the zip file from the following location & move it to your Hadoop environment. Place this data under /LDZ/data/ in Hadoop**

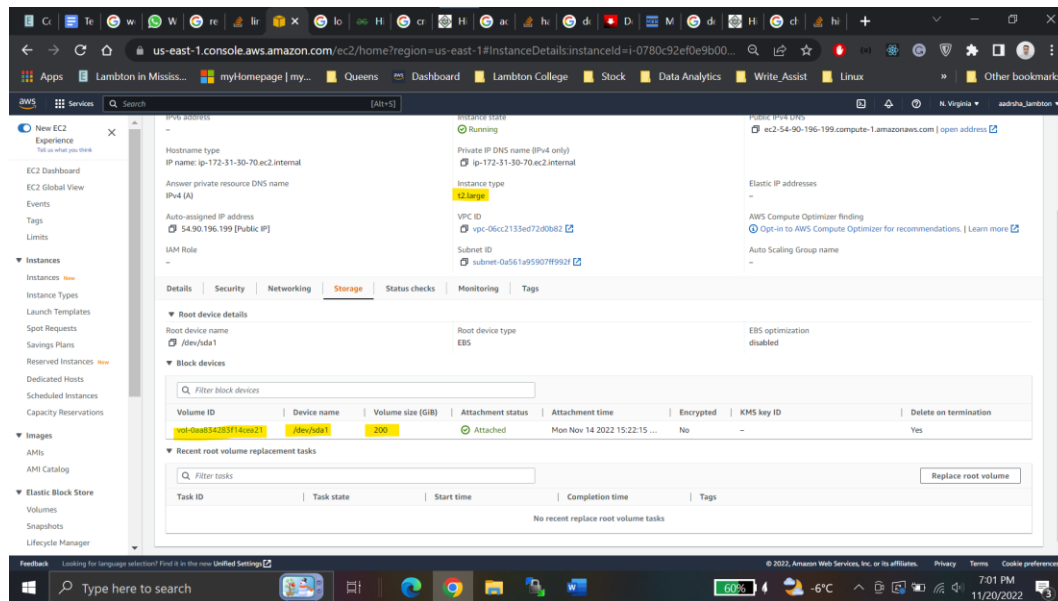
## Architecture for Hive



## Hadoop on Cloud Environment

Ubuntu 18.04 was installed on AWS ec2(t2.large) with 8gb of memory and 300gb of storage. Hadoop 2.6, Hive and mysql was installed on the machine.





## Single Node on Hadoop was set up using following commands on AWS ec2.

1. create a new 'hadoop' user in ubuntu.
2. create '/home/hadoop/work' and '/home/hadoop/work/hadoopdata' folders

```
mkdir /home/hadoop/work
```

```
mkdir /home/hadoop/work/hadoopdata
```

3. Download 'hadoop-2.6.0.tar.gz' file from Apache mirrors  
<https://archive.apache.org/dist/hadoop/core/hadoop-2.6.0/>, copy 'hadoop-2.6.0.tar.gz' file into this '/home/hadoop/work' directory and extract the tar file into same directory.

```
tar -xvzf hadoop-2.6.0.tar.gz
```

4. Open the '~/.bashrc' file on all the machines and add the following lines at the end and save:

command: gedit ~/.bashrc

```
export JAVA_HOME=/usr/lib/jvm/java-1.8.0-openjdk-amd64
```

```
export HADOOP_HOME=/home/hadoop/work/hadoop-2.6.0
```

```
export HADOOP_COMMON_HOME=$HADOOP_HOME
```

```
export HADOOP_HDFS_HOME=$HADOOP_HOME
```

```
export HADOOP_MAPRED_HOME=$HADOOP_HOME
```

```
export YARN_HOME=$HADOOP_HOME
```

```
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
export YARN_CONF_DIR=$HADOOP_HOME/etc/hadoop
export
PATH=$JAVA_HOME/bin:$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH
```

5. Enter the below commands on terminal:

```
ssh localhost

ssh-keygen -t rsa -P " " -f ~/.ssh/id_rsa

cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
```

6. Update the '/home/hadoop/work/hadoop-2.6.0/etc' folder files 'hadoop-env.sh','core-site.xml','hdfs-site.xml','mapred-site.xml', 'yarn-env.sh', 'yarn-site.xml','masters' and 'slaves' files as per the below configurations

hadoop-env.sh

```
=====

# The java implementation to use.

export JAVA_HOME=/usr/lib/jvm/java-1.8.0-openjdk-amd64

export JAVA_HOME=${JAVA_HOME}
```

core-site.xml

```
=====

<configuration>

  <property>

    <name>fs.defaultFS</name>

    <value>hdfs://localhost:8020</value>

  </property>

</configuration>
```

hdfs-site.xml

```
=====

<configuration>
```

```
<property>
    <name>dfs.replication</name>
    <value>1</value>
</property>
<property>
    <name>dfs.namenode.name.dir</name>
    <value>file:/home/hadoop/work/hadoopdata/dfs/name</value>
</property>
<property>
    <name>dfs.datanode.data.dir</name>
    <value>file:/home/hadoop/work/hadoopdata/dfs/data</value>
</property>
<property>
    <name>dfs.namenode.checkpoint.dir</name>
    <value>file:/home/hadoop/work/hadoopdata/dfs/namesecondary</value>
</property>
</configuration>
```

mapred-env.sh

=====

```
# export JAVA_HOME=/home/y/libexec/jdk1.8.0/
export JAVA_HOME=/usr/lib/jvm/java-1.8.0-openjdk-amd64
```

mapred-site.xml

=====

```
<configuration>
    <property>
        <name>mapreduce.framework.name</name>
```

```
        <value>yarn</value>
    </property>
</configuration>
```

yarn-env.sh

=====

```
# export JAVA_HOME=/home/y/libexec/jdk1.8.0/
export JAVA_HOME=/usr/lib/jvm/java-1.8.0-openjdk-amd64
```

yarn-site.xml

=====

```
<configuration>
  <property>
    <name>yarn.resourcemanager.hostname</name>
    <value>localhost</value>
  </property>
  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>
  <property>
    <name>yarn.nodemanager.aux-services.mapreduce_shuffle.class</name>
    <value>org.apache.hadoop.mapred.ShuffleHandler</value>
  </property>
</configuration>
```

slaves

=====

localhost

7. Format the 'namenode' from current machine using this command:

```
hadoop namenode -format or hdfs namenode -format
```

8. Start the hadoop by using this command on current machine:

```
start-dfs.sh
```

```
start-all.sh (depricated)
```

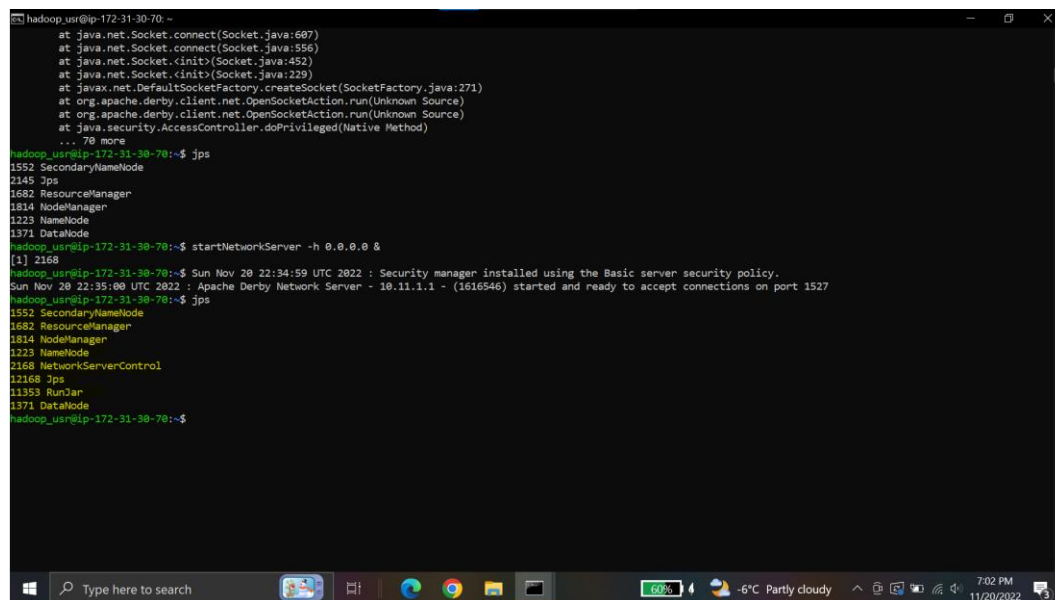
9. Stop the hadoop by using this command on current machine:

```
start-yarn.sh
```

```
stop-all.sh (depricated)
```

10. jps

With JPS it can be seen that namenode and resource manager are started.



```
hadoop_usr@ip-172-31-30-70:~$  
at java.net.Socket.connect(Socket.java:607)  
at java.net.Socket.connect(Socket.java:556)  
at java.net.Socket.<init>(Socket.java:452)  
at java.net.Socket.<init>(Socket.java:229)  
at javax.net.DefaultSocketFactory.createSocket(SocketFactory.java:271)  
at org.apache.derby.client.net.OpenSocketAction.run(Unknown Source)  
at org.apache.derby.client.net.OpenSocketAction.run(Unknown Source)  
at java.security.AccessController.doPrivileged(Native Method)  
... 78 more  
hadoop_usr@ip-172-31-30-70:~$ jps  
1552 SecondaryNameNode  
2145 Jps  
1682 ResourceManager  
1814 NodeManager  
1223 NameNode  
1371 DataNode  
hadoop_usr@ip-172-31-30-70:~$ startNetworkServer -h 0.0.0.0 &  
[1] 2168  
hadoop_usr@ip-172-31-30-70:~$ Sun Nov 20 22:34:59 UTC 2022 : Security manager installed using the Basic server security policy.  
Sun Nov 20 22:35:00 UTC 2022 : Apache Derby Network Server - 10.11.1.1 - (1616546) started and ready to accept connections on port 1527  
hadoop_usr@ip-172-31-30-70:~$ jps  
1552 SecondaryNameNode  
1682 ResourceManager  
1814 NodeManager  
1223 NameNode  
2168 NetworkServerControl  
2168 Jps  
11359 RunJar  
1371 DataNode  
hadoop_usr@ip-172-31-30-70:~$
```

To check hadoop storage hdfs dfs -df -h command can be used

```
hadoop_usr@ip-172-31-30-70: ~/work
hadoop_usr@ip-172-31-30-70:~/home/ubuntu$ cd
hadoop_usr@ip-172-31-30-70:~$ hdfs dfs -ls /LDZ/data/
hadoop_usr@ip-172-31-30-70:~$ hdfs dfs -ls /LDZ/
Found 2 items
drwxr-xr-x 1 hadoop_usr supergroup 0 2022-11-20 23:43 /LDZ/data
-rwxrwxrwx 1 hadoop_usr supergroup 46019 2022-11-20 23:31 /LDZ/hivexmlserde-1.0.2.0.jar
hadoop_usr@ip-172-31-30-70:~$ hdfs dfs -ls /LDZ/data/
hadoop_usr@ip-172-31-30-70:~$ ls
derby.log jar_files.zip myderby work
hadoop_usr@ip-172-31-30-70:~$ hdfs dfs -ls /LDZ/data/
hadoop_usr@ip-172-31-30-70:~$ ls
derby.log jar_files.zip myderby work
hadoop_usr@ip-172-31-30-70:~$ cd work/
hadoop_usr@ip-172-31-30-70:~/work$ ls
Posts.xml
hadoop_usr@ip-172-31-30-70:~/work$ hdfs dfs -put miniposts.xml /LDZ/data/
hadoop_usr@ip-172-31-30-70:~/work$ hdfs dfs -ls /LDZ/data/
Found 1 items
-rwxr-xr-x 1 hadoop_usr supergroup 8192 2022-11-20 23:45 /LDZ/data/miniposts.xml
hadoop_usr@ip-172-31-30-70:~/work$ hdfs dfs -du
du: '.': No such file or directory
hadoop_usr@ip-172-31-30-70:~/work$ hdfs dfs -df
Filesystem          Size      Used Available Use%
hdfs://localhost:8020 20722917440 167936 104889376768 0%
hadoop_usr@ip-172-31-30-70:~/work$ hdfs dfs -df -h
Filesystem          Size      Used Available Use%
hdfs://localhost:8020 193.6 G 164 K 97.7 G 0%
hadoop_usr@ip-172-31-30-70:~/work$
```

Once hadoop is installed create directory named LDZ and /LDZ/data in Hadoop and get the data in local system, unzip it and transfer it to hadoop .

**Download the zip file from the following location & move it to your Hadoop environment. Place this data under /LDZ/data/ in Hadoop**

```
hadoop_usr@ip-172-31-30-70: ~/work
hadoop_usr@ip-172-31-30-70:~$ ls
derby.log metastore_db work
hadoop_usr@ip-172-31-30-70:~$ pwd
/home/hadoop_usr
hadoop_usr@ip-172-31-30-70:~$ ls
derby.log metastore_db work
hadoop_usr@ip-172-31-30-70:~$ cd work
hadoop_usr@ip-172-31-30-70:~/work$ wget https://archive.org/download/stackexchange/stackoverflow.com-Posts.7z
--2022-11-17 02:39:32-- https://archive.org/download/stackexchange/stackoverflow.com-Posts.7z
Resolving archive.org (archive.org)... 207.241.224.2
Connecting to archive.org (archive.org)|207.241.224.2|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://ia600107.us.archive.org/27/items/stackexchange/stackoverflow.com-Posts.7z [following]
--2022-11-17 02:39:32-- https://ia600107.us.archive.org/27/items/stackexchange/stackoverflow.com-Posts.7z
Resolving ia600107.us.archive.org (ia600107.us.archive.org)... 207.241.227.247
Connecting to ia600107.us.archive.org (ia600107.us.archive.org)|207.241.227.247|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 19430749009 (18G) [application/x-7z-compressed]
Saving to: 'stackoverflow.com-Posts.7z'

stackoverflow.com-Posts.7z 20%[=====>] 3.67G 8.37MB/s eta 44m 21s
```

The zip file was downloaded into the ubuntu machine using **wget** command



File transferred to HDFS into the location `/DWZ/data` using `put` command in HDFS

[illegible]

## Step 2:

### 3. Build a data warehouse (location:/DWZ) with data partitioned based on Creation Date and then Post Type.

## Install Hive in hadoo using following commands

1. Download 'apache-hive-1.2.1-bin' file from apache mirrors <https://archive.apache.org/dist/hive/hive-1.2.1/>, copy 'apache-hive-1.2.1-bin' file into this '/home/hadoop/work' directory and extract the tar files into same directory.

Using terminal:tar -xvzf apache-hive-1.2.1-bin

2. Open the '~/.bashrc' file using command: `gedit ~/.bashrc` in the terminal add below lines at the end of the document.

```
export HIVE_HOME=/home/hadoop/work/apache-hive-1.2.1-bin
```

```
export PATH=$HIVE_HOME/bin:$PATH
```

3. Copy mysql-connector-java-5.1.38 to /home/hadoop/work/apache-hive-1.2.1-bin/lib

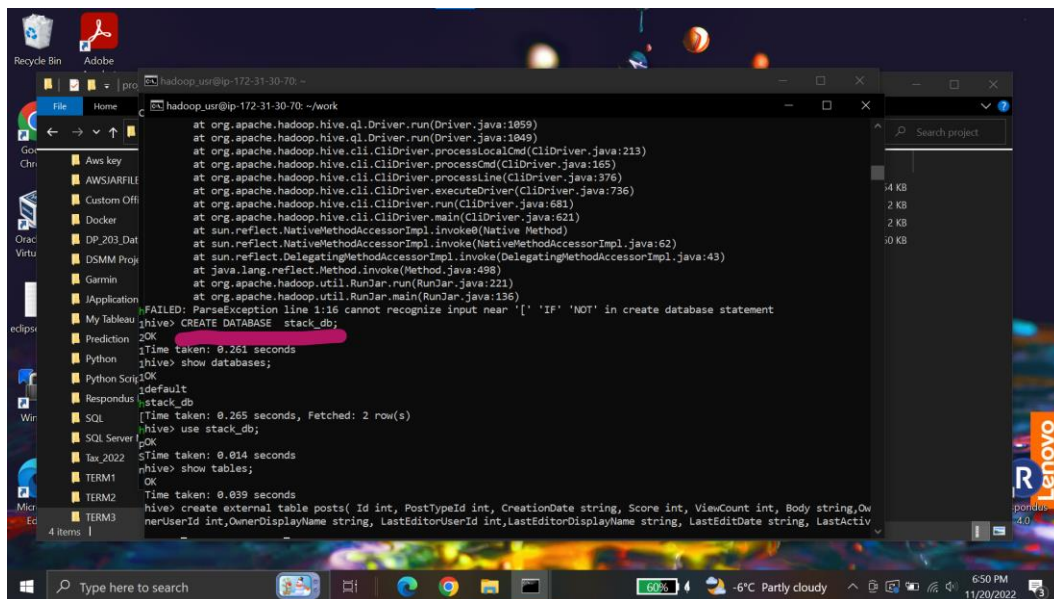
Copy /home/hadoop/work/db-derby-10.11.1.1-bin/lib/derbyclient.jar to /home/hadoop/work/apache-hive-1.2.1-bin/lib

4. The following configuration files are required for hive to be run in different modes.

- 1.Hive-site.xml (Main configuration file)
- 2.hive-site.xml\_local (For this mode, copy paste this script into Main configuration file)
- 3.Hive-site.xml\_derby (For this mode, copy paste this script into Main configuration file)
- 4.hive-site.xml\_mysql (For this mode, copy paste this script into Main configuration file)

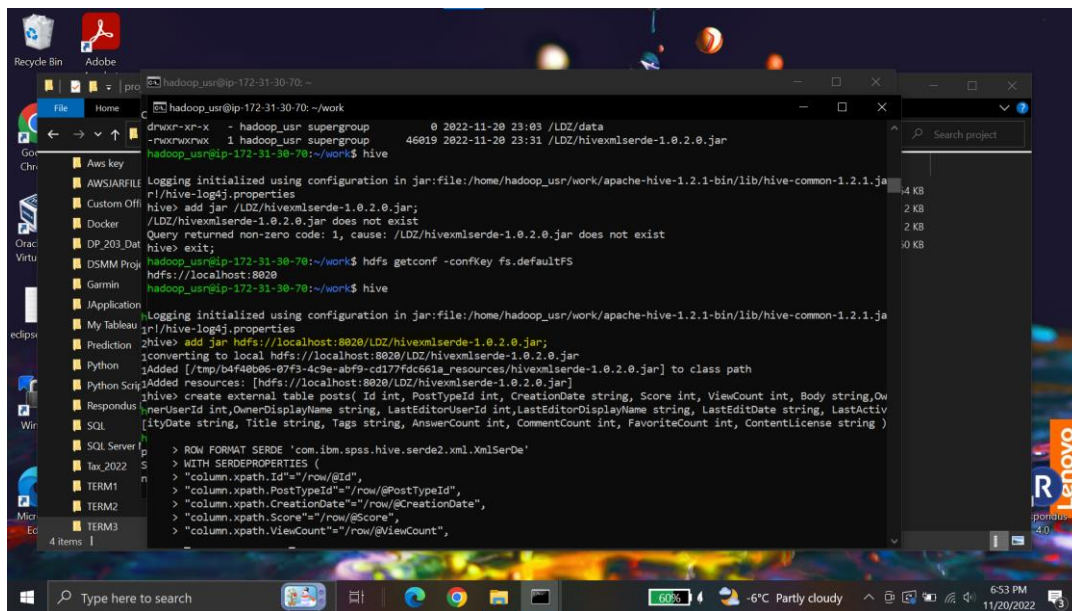
5. Open hive-site.xml from /home/hadoop/work/apache-hive-1.2.1-bin/conf. Edit configuration with below properties.

Once Hive is installed, add xml create database and table in Hadoop

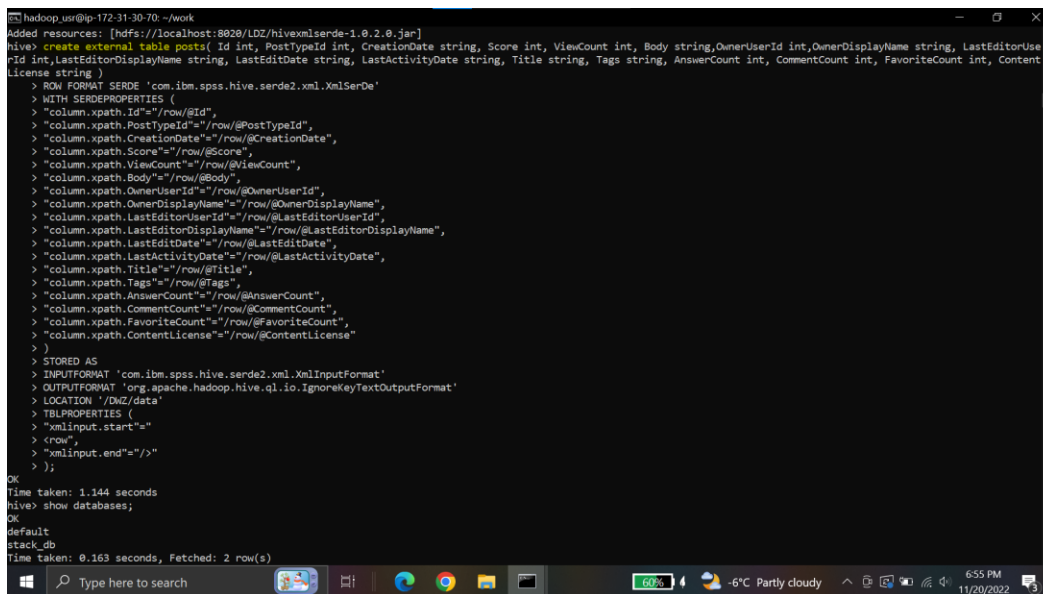


```
at org.apache.hadoop.hive.q1.Driver.run(Driver.java:1859)
at org.apache.hadoop.hive.q1.Driver.run(Driver.java:1849)
at org.apache.hadoop.hive.cli.CliDriver.processLocalCmd(CliDriver.java:213)
at org.apache.hadoop.hive.cli.CliDriver.processCmd(CliDriver.java:165)
at org.apache.hadoop.hive.cli.CliDriver.processLine(CliDriver.java:376)
at org.apache.hadoop.hive.cli.CliDriver.executeDriver(CliDriver.java:736)
at org.apache.hadoop.hive.cli.CliDriver.run(CliDriver.java:681)
at org.apache.hadoop.hive.cli.CliDriver.main(CliDriver.java:621)
at sun.reflect.NativeMethodAccessorImpl.invoke(Native Method)
at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:62)
at java.lang.reflect.Method.invoke(Method.java:498)
at org.apache.hadoop.util.RunJar.run(RunJar.java:221)
at org.apache.hadoop.util.RunJar.main(RunJar.java:136)
hive> CREATE DATABASE stack_db;
Time taken: 0.261 seconds
hive> show databases;
default
stack_db
Time taken: 0.265 seconds, Fetched: 2 row(s)
hive> use stack_db;
Time taken: 0.014 seconds
hive> show tables;
OK
Time taken: 0.039 seconds
hive> create external table posts( Id int, PostTypeId int, CreationDate string, Score int, ViewCount int, Body string, OwnerUserId int, OwnerDisplayName string, LastEditorUserId int, LastEditorDisplayName string, LastEditDate string, LastActiveDate string)
```

We need serializer before creating table in hive, so download the serializer put it in hadoop and add jar file to hive providing the location



## Create External table in hive



create external table posts( Id int, PostTypeId int, CreationDate string, Score int, ViewCount int, Body string, OwnerUserId int, OwnerDisplayName string, LastEditorUserId int, LastEditorDisplayName string, LastEditDate string, LastActivityDate string, Title string, Tags string, AnswerCount int, CommentCount int, FavoriteCount int, ContentLicense string )

ROW FORMAT SERDE 'com.ibm.spss.hive.serde2.xml.XmlSerDe'

WITH SERDEPROPERTIES (

"column.xpath.Id"="/row/@Id",

```

"column.xpath.PostTypeId"="/row/@PostTypeId",
"column.xpath.CreationDate"="/row/@CreationDate",
"column.xpath.Score"="/row/@Score",
"column.xpath.ViewCount"="/row/@ViewCount",
"column.xpath.Body"="/row/@Body",
"column.xpath.OwnerUserId"="/row/@OwnerUserId",
"column.xpath.OwnerDisplayName"="/row/@OwnerDisplayName",
"column.xpath.LastEditorUserId"="/row/@LastEditorUserId",
"column.xpath.LastEditorDisplayName"="/row/@LastEditorDisplayName",
"column.xpath.LastEditDate"="/row/@LastEditDate",
"column.xpath.LastActivityDate"="/row/@LastActivityDate",
"column.xpath.Title"="/row/@Title",
"column.xpath.Tags"="/row/@Tags",
"column.xpath.AnswerCount"="/row/@AnswerCount",
"column.xpath.CommentCount"="/row/@CommentCount",
"column.xpath.FavoriteCount"="/row/@FavoriteCount",
"column.xpath.ContentLicense"="/row/@ContentLicense"
)

```

STORED AS

INPUTFORMAT 'com.ibm.spss.hive.serde2.xml.XmlInputFormat'

OUTPUTFORMAT 'org.apache.hadoop.hive.ql.io.IgnoreKeyTextOutputFormat'

LOCATION '/DWZ/data'

TBLPROPERTIES (

"xmlinput.start"="

<row",

"xmlinput.end"="/>"

);

Load data into posts table

```
hive> LOAD DATA INPATH 'hdfs://localhost:8020/LDZ/data/miniposts.xml' into posts;  
MismatchedTokenException(261+513)  
    at org.apache.hadoop.mapreduce.lib.input.FileInputFormat$2.doOpen(FileInputFormat.java:617)  
    at org.apache.hadoop.mapreduce.TaskIOContext.openTask(TaskIOContext.java:115)  
    at org.apache.hadoop.mapreduce.TaskIOContext.openInputs(TaskIOContext.java:180)  
    at org.apache.hadoop.mapred.MapRunner.run(MapRunner.java:160)  
    at org.apache.hadoop.mapred.YarnRunner.run(YarnRunner.java:119)  
    at org.apache.hadoop.mapred.Main.main(Main.java:262)  
    at org.apache.hadoop.mapred.Driver.compile(Driver.java:166)  
    at org.apache.hadoop.mapred.Driver.compile(Driver.java:396)  
    at org.apache.hadoop.mapred.Driver.compileInternal(Driver.java:1122)  
    at org.apache.hadoop.mapred.Driver.runInterval(Driver.java:1178)  
    at org.apache.hadoop.mapred.Driver.run(Driver.java:1059)  
    at org.apache.hadoop.mapred.Driver.run(Driver.java:1049)  
    at org.apache.hadoop.mapred.CliDriver.processLocalCmd(CliDriver.java:213)  
    at org.apache.hadoop.mapred.CliDriver.processCmd(CliDriver.java:165)  
    at org.apache.hadoop.mapred.CliDriver.processLine(CliDriver.java:376)  
    at org.apache.hadoop.mapred.CliDriver.executeDriver(CliDriver.java:736)  
    at org.apache.hadoop.mapred.CliDriver.run(CliDriver.java:681)  
    at org.apache.hadoop.mapred.CliDriver.main(CliDriver.java:621)  
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)  
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)  
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)  
    at java.lang.reflect.Method.invoke(Method.java:498)  
    at org.apache.hadoop.util.RunJar.run(RunJar.java:221)  
    at org.apache.hadoop.util.RunJar.main(RunJar.java:136)  
FAILED: ParseException line 1:17 mismatched input 'hdfs' expecting string literal near 'INPATH' in load statement  
hive> LOAD DATA INPATH 'hdfs://localhost:8020/LDZ/data/miniposts.xml' into posts;  
FAILED: ParseException line 1:69 missing TABLE at 'posts' near '{EOF}'  
hive> show tables;  
OK  
Time taken: 0.024 seconds  
hive> show databases;  
OK  
default  
stack_00  
Time taken: 0.008 seconds, Fetched: 2 row(s)  
hive> use default;  
OK  
Time taken: 0.025 seconds
```

### Show few data in hive in table posts

```

C:\hdopoc>ip-172.31.31.215 -work\hdopoc.2.6.0\ht\hdopoc
posts
Time taken: 0.017 seconds, Fetched: 1 row(s)
hive> select * from posts;
+-----+
| x |
+-----+
1      2008-07-31T21:41:52.667 781      67914  <p>I want to assign the decimal variable &quot;trans&quot; to the double variable &quot;this.Opacity&quot;.</p>

```

```

hive> select id, posttypeid, creationdate from posts limit 20;
OK
4      1      2008-07-31T21:42:52.667
6      1      2008-07-31T22:08:08.620
7      2      2008-07-31T22:17:57.883
9      1      2008-07-31T23:40:59.743
11     1      2008-07-31T23:55:37.967
12     2      2008-07-31T23:56:41.303
13     1      2008-08-01T00:42:38.903
14     1      2008-08-01T00:59:11.177
4      1      2008-07-31T21:42:52.667
6      1      2008-07-31T22:08:08.620
7      2      2008-07-31T22:17:57.883
9      1      2008-07-31T23:40:59.743
11     1      2008-07-31T23:55:37.967
12     2      2008-07-31T23:56:41.303
13     1      2008-08-01T00:42:38.903
14     1      2008-08-01T00:59:11.177
16     1      2008-08-01T04:59:33.643
17     1      2008-08-01T05:09:55.993
18     2      2008-08-01T05:12:44.193
19     1      2008-08-01T05:21:22.257
Time taken: 0.087 seconds, Fetched: 20 row(s)
hive> _

```

### Step 3:

#### Querying Database:

What are the top 10 most answered questions in Stack Overflow posts for a particular creation date?

Query:

Select \* from posts where

creationdate ='2008-07-31T21:42:52.667'

ORDER BY answercount limit 10;

;





## **Conclusion**

Hence the required architecture for the analysis of the stackoverflow data was built using

- Hadoop
- Hive

On Hive, mysql mode was used and the data was loaded in the external table and queries were performed to answer the given questions.