**Project Report**

**Course Name: Software and Data Engineering (SDE)**

**Course Code: CSL7090**

# Exploring BigData: CDH Platform Architecture, Local Machine EDA

**Submitted To: -**

Dr. Sumit Kalra,

Assistant Professor,

Dept. of CSE,

IIT Jodhpur.

**Submitted By: -**

Bhawna Bhoria (M22MA003)

Jash Patel (M22CS061)

# Introduction

As part of this project , we will be observing and analyzing processing techniques across CDH, local machine platform to comprehend intricate systems and address data processing challenges in the context of expanding data volumes.
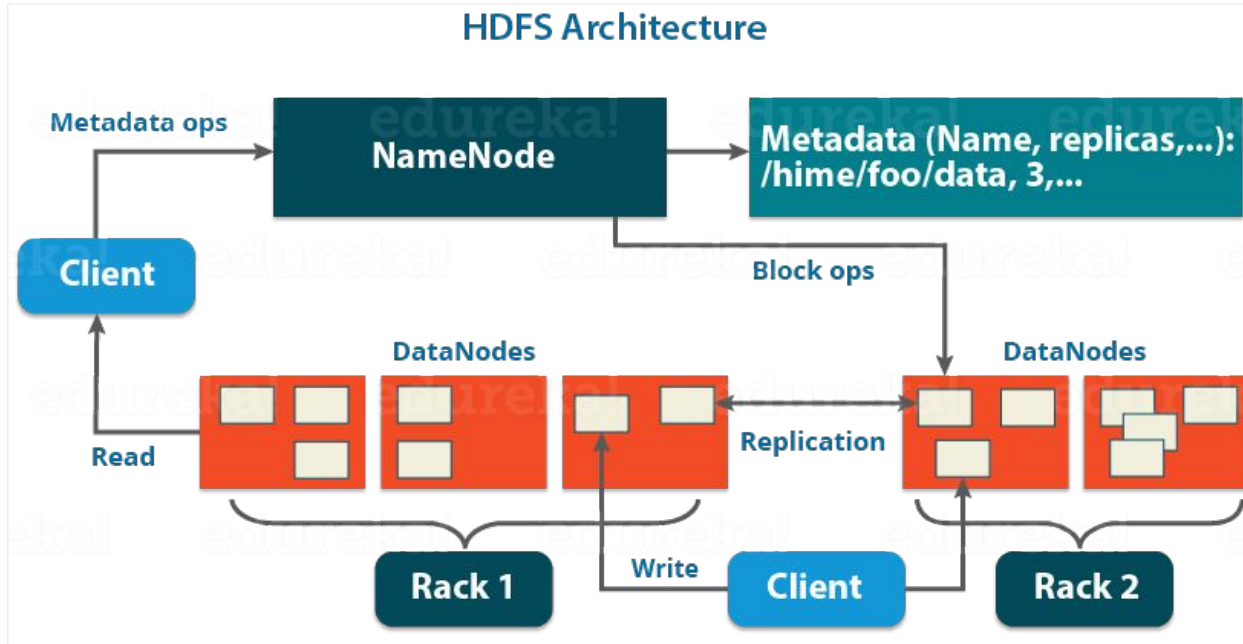
Platform 1:-

CDH - Architecture, Storage System and Job Scheduling.

Platform 2:-

EDA - Exploratory Data Analysis on Local Machine.

# Platform 1 : Cloudera Distribution of Hadoop



The customised distribution employs the Namenode and YARN Resource Manager to oversee and allocate resources in a Hadoop cluster. The functionalities of the Job Tracker and Task Tracker have been substituted by the ResourceManager and NodeManagers, providing a resource management architecture that is more adaptable and capable of handling larger workloads.

# Processes Running on CDH Architecture:-

jps output :-

```
[[root@quickstart /]# jps
880 JournalNode
2000 NodeManager
7292
6163 Bootstrap
2713 HMaster
5165 HistoryServer
5073 Bootstrap
5494 HRegionServer
1673 Bootstrap
674 DataNode
1321 SecondaryNameNode
```

```
1081 NameNode
563 QuorumPeerMain
1802 JobHistoryServer
7183 Bootstrap
2281 ResourceManager
9166 Jps
3953 RunJar
3346 ThriftServer
7443
3030 RESTServer
3616 RunJar
```
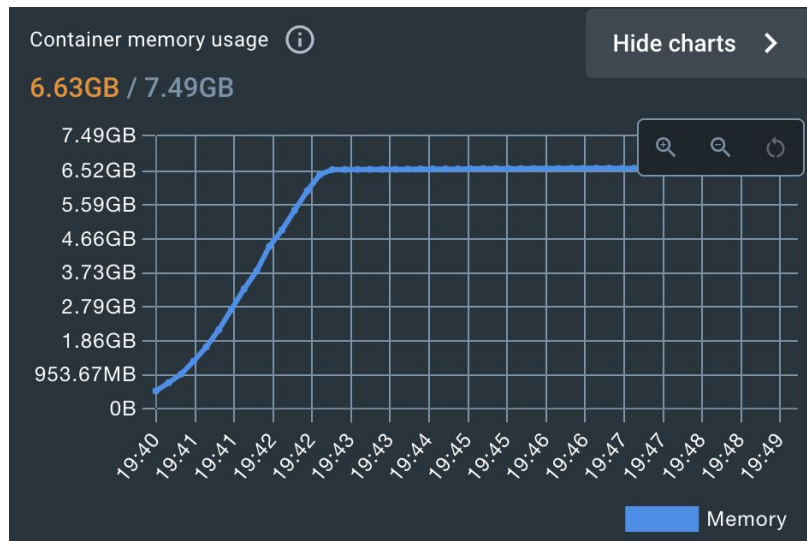
Checking Namenode :-

```
[[root@quickstart /]# hdfs getconf —confKey fs.defaultFS
hdfs://quickstart.cloudera:8020
```

# Laptop configurations and CDH Image Size

Physical Memory on device: 16GB

CDH Distribution Memory Occupance : 8 GB

# CDH Ecosystem : Cluster and HDFS Storage

# Taking File to hdfs

```
[[root@quickstart /]#
[root@quickstart /]# hdfs dfs -put /user/sde_project/taxi_zone_lookup.csv /user/cloudera/
```

The Hadoop Distributed File System (HDFS) is renowned for its fault-tolerant nature. It allows clients to upload files from NameNodes, which then direct them to DataNodes for storage information. The client then uploads the file to the specified DataNode, completing the write request. In a read operation, the client seeks metadata from NameNodes, which then identifies the DataNode server where the file resides in the cluster and randomly selects one.
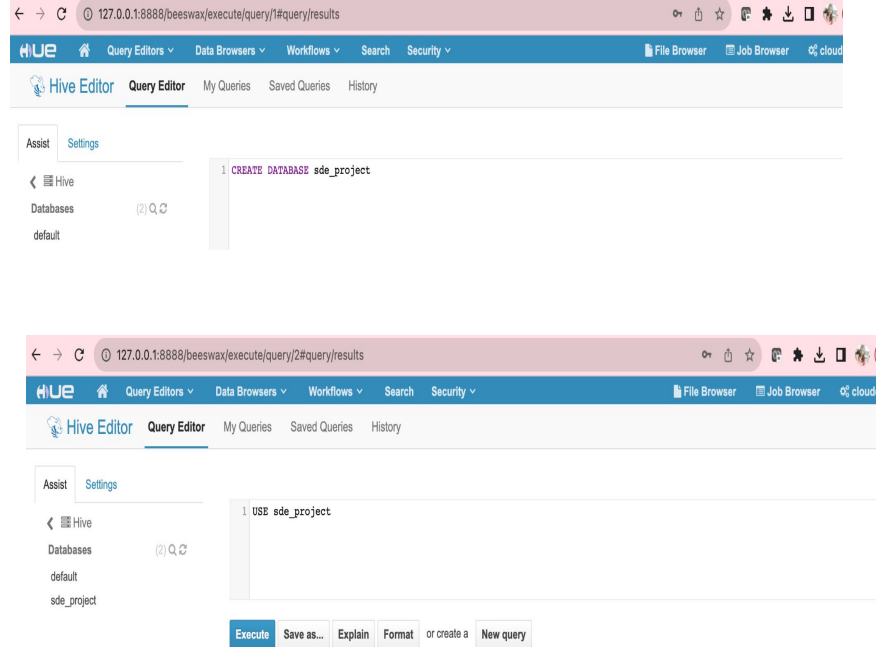
File placed at HDFS storage system into distributed blocks of 64 bytes.

# HIVE

Hive, a Hadoop-based data warehouse tool, efficiently processes queries and generates results for extensive datasets. It executes statements, similar to MySQL, and takes charge of user tasks. The compiler retrieves metadata, compiles the task, selects the optimal strategy, and delivers results. The pseudocode outlines Hive's data transformation, loading, and extraction process.

# Querying on Hive:-

❮ ▤ sde_project

**Tables**    (1) 🔍 ⟳

▦ yellow_taxi_data

```
1
2   -- Create a Hive table for the location data
3   CREATE TABLE IF NOT EXISTS sde_project.yellow_taxi_data (
4       LocationID INT,
5       Borough STRING,
6       Zone STRING,
7       service_zone STRING
8   )
9   ROW FORMAT DELIMITED
10  FIELDS TERMINATED BY ','
11  STORED AS TEXTFILE
12  TBLPROPERTIES("skip.header.line.count"="1")
```

**Execute**   Save as...   Explain   Format   or create a   New query

```
1   LOAD DATA INPATH 'hdfs:///user/cloudera/taxi_zone_lookup.csv' INTO TABLE sde_project.yellow_taxi_cab;
```

## Querying on Hive:-



Map Reduce job Runs in background on YARN.

# What happens behind the curtains?- YARN

Application Master Instance gets created and State is ACCEPTED.

# Logs are stored on HistoryServer and Jobs is moved to FINISHED Applications.

# CDH Oozie Workflows

Workflow to create HDFS directories and files (action name : fs-12c7)

```xml
<workflow-app name="My_Workflow" xmlns="uri:oozie:workflow:0.5">
    <start to="fs-12c7"/>
    <kill name="Kill">
        <message>Action failed, error message[${wf:errorMessage(wf:lastErrorNode())}]</message>
    </kill>
    <action name="fs-12c7">
        <fs>
            <touchz path='${nameNode}/user/cloudera/test'/>
        </fs>
        <ok to="End"/>
        <error to="Kill"/>
    </action>
    <end name="End"/>
</workflow-app>
```

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000001-
231122121616918-oozie-
oozi-W

**VARIABLES**

| Name | Value |
| --- | --- |
| dryrun | False |
| hue-id-w | 3 |
| jobTracker | localhost:8032 |
| mapreduce.job.user.name | cloudera |
| nameNode | hdfs://quickstart.cloudera:8020 |
| oozie.use.system.libpath | True |
| oozie.wf.application.path | hdfs://quickstart.cloudera:8020/user/hue/oozie/workspaces/hue-oozie-1700655912.35 |
| security_enabled | False |
| user.name | cloudera |

Back

Workflow to run shell script checking availability of file on hdfs

```xml
<workflow-app name="Test" xmlns="uri:oozie:workflow:0.5">
    <start to="shell-869a"/>
    <kill name="Kill">
        <message>Action failed, error message[${wf:errorMessage(wf:lastErrorNode())}]</message>
    </kill>
    <action name="shell-869a">
        <shell xmlns="uri:oozie:shell-action:0.1">
            <job-tracker>${jobTracker}</job-tracker>
            <name-node>${nameNode}</name-node>
            <exec>/user/cloudera/file_avail_check.sh</exec>
                <capture-output/>
        </shell>
        <ok to="End"/>
        <error to="Kill"/>
    </action>
    <end name="End"/>
</workflow-app>
```

# Platform 2 : EDA using standard libraries and dataprep

The main objective is to tune down the CPU Utilisation while performing Exploratory Data Analysis on Local Machine.

1. Traditional Method : CPU utilization of EDA using standard libraries

2. EDA Library Method : CPU Utilisation after data profiling using "dataprep"



CPU Utilization has come down significantly after using the profiling method of "dataprep".

# Performance comparison on EDA

| CPU Core | EDA with dataprep in % | EDA with standard lib in % |
|----------|------------------------|----------------------------|
| Core 1 | 7 | 68.4 |
| Core 2 | 38 | 69.4 |
| Core 3 | 7.9 | 59.6 |
| Core 4 | 10.8 | 60.6 |
| Core 5 | 5.1 | 35.6 |
| Core 6 | 72.3 | 100 |
| Core 7 | 6.1 | 88.2 |
| Core 8 | 5.1 | 53.5 |
| Average | 19.0375 | 66.9125 |

# Conclusion

The project explores big data using CDH platform architecture and local machine Exploratory Data Analysis (EDA). It highlights the need for robust big data ecosystems due to the surge in global data storage capacity. Practical insights into CDH, local machine EDA are presented, showcasing their processing techniques. The project provides a holistic view of big data tools and platforms, laying the foundation for future data engineering and analytics endeavors.

# References

https://blog.clairvoyantsoft.com/cloduera-quickstart-vm-using-docker-on-mac-2308acd196f2

https://cwiki.apache.org/confluence/display/Hive/LanguageManual+DDL

https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9463522&tag=1

https://github.com/sfu-db/dataprep

# Thank You