# Report

## Group Members:

- Bhavana Devulapally
- Shusrita Venugopal
- Neha Navarkar

Our submission (attached zip file) contains separate folders for each Part. Following is the list of files submitted:

Part 1:

1. Dockerfile
2. bootstrap.sh
3. core-site.sh
4. hdfs-site.xml
5. mapred-site.xml
6. yarn-site.xml

Part 2:

1. ngrammaper.py
2. ngramreducer.py
3. input.txt

Part 3:

Contains mapper and reducer for each problem.

## Part 1: Setting up Hadoop in Docker

In this project, we explored Docker images and learned how to build them. Docker is a powerful virtualization tool similar to the previous tool, Sandbox, but with additional features for system developers. It combines the functionality of a version control system (VCS) with that of a programming language, such as Java, for virtual machines.
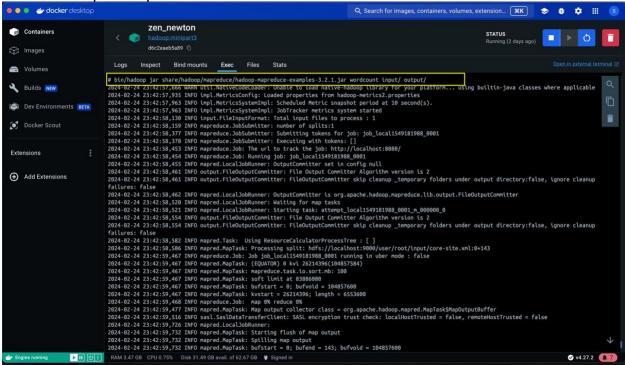
1.To begin, we set up a single node in local Docker by building the Docker image based on the Dockerfile in the "local_demo.zip" file, which contains the necessary code and files for our project.

2. We then setup SSH keys using ssh-keyget command to generate public-private RSA key pair.

3. We configured start-dfs.sh and start-yarn.sh, to initialize the core components of a Hadoop cluster, enabling data storage and resource management for distributed computing tasks. These scripts are typically executed on the master node (or nodes) of the Hadoop cluster to start the necessary services.

4. We configured stop-dfs.sh and stop-yarn.sh, to gracefully stop the core components of a Hadoop cluster, ensuring that data is safely persisted, and resources are released. These scripts are typically executed on the master node (or nodes) of the Hadoop cluster to stop the necessary services.

5. We configured hdfs-site.xml to configure properties related to HDFS, such as replication factor, block size, and data node settings. It contains parameters that govern the behavior of the HDFS components, including the NameNode, DataNode, and Secondary NameNode. Some common properties defined in hdfs-site.xml include:

- dfs.replication: Specifies the default replication factor for file blocks in HDFS.
- dfs.blocksize: Defines the default block size for files in HDFS.
- dfs.namenode.name.dir: Specifies the directory where the NameNode stores its metadata.
- dfs.datanode.data.dir: Defines the directory where DataNode stores HDFS blocks.

6. We configured mapred.xml to configure properties related to the MapReduce framework, including JobTracker and TaskTracker settings. It defines parameters that control job execution, task scheduling, and resource allocation within the Hadoop cluster. Some common properties defined in mapred.xml include:

- mapred.job.tracker: Specifies the hostname and port of the JobTracker.
- mapred.map.tasks: Defines the number of map tasks to run per job.
- mapred.reduce.tasks: Specifies the number of reduce tasks to run per job.
- mapred.tasktracker.map.tasks.maximum: Defines the maximum number of map tasks that can be executed concurrently on a TaskTracker node.
- mapred.tasktracker.reduce.tasks.maximum: Defines the maximum number of reduce tasks that can be executed concurrently on a TaskTracker node.

7. We then executed The bin/hadoop namenode -format command initializes and prepares the Hadoop Distributed File System (HDFS) by creating necessary metadata structures and clearing any existing data. It ensures a clean and consistent starting state for the HDFS, essential for setting up new clusters or reinstalling HDFS.

8. Here we encountered error for java_home on MAC devices. We set JAVA_HOME and PATH in /opt/hadoop/etc/hadoop/hadoop-enc.sh file using the following commands:

- cd /opt/hadoop/etc/hadoop
- nano hadoop-env.sh
- In the last line of Hadoop-env.sh file add export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64

9. Further, we executed ./start-dfs.sh command to start the HDFS daemons i.e. NameNode and DataNode daemons.

10. With JPS command we verified that Hadoop daemons (like NameNode, DataNode, ResourceManager) are up and running. As we know, the NameNode specifically listens on port 9000.

```
# jps
641 DataNode
7573 Jps
473 NameNode
863 SecondaryNameNode
```

Wordcount Example:
bin/hadoop      jar      share/hadoop/mapreduce/hadoop-mapreduce-examples-3.2.1.jar wordcount input/ output/
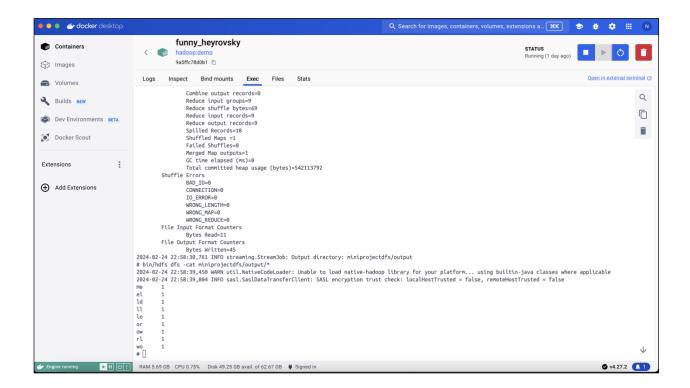


Output:



# Part 2: Hadoop program that produces the n-gram frequencies of the text "Helloworld".

We executed the below command for ngram frequency (n) "2".

bin/hadoop jar ./share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar -files /miniproject1/ngrammapper.py,/miniproject1/ngramreducer.py -input miniprojectdfs/ - output miniprojectdfs/output -mapper "python3 ngrammapper.py 2" -reducer "python3 ngramreducer.py"
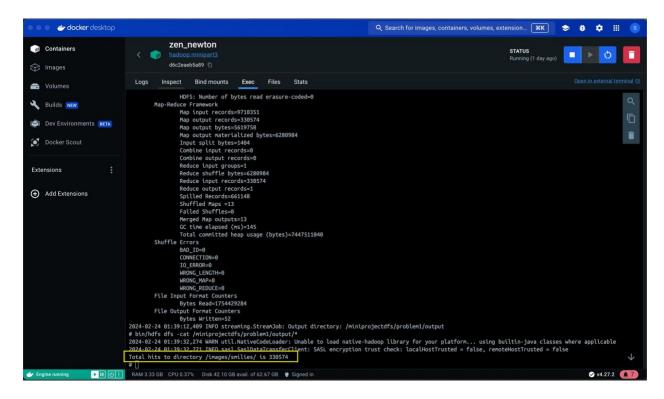
# **Part 3:** Hadoop program to analyze real logs.

In this part, we developed MapReduce programs to analyze a real anonymous log file and answer several questions based on the log data. The log file, named access_log, follows the Common Log Format with additional fields that we ignored for the purpose of these problems.

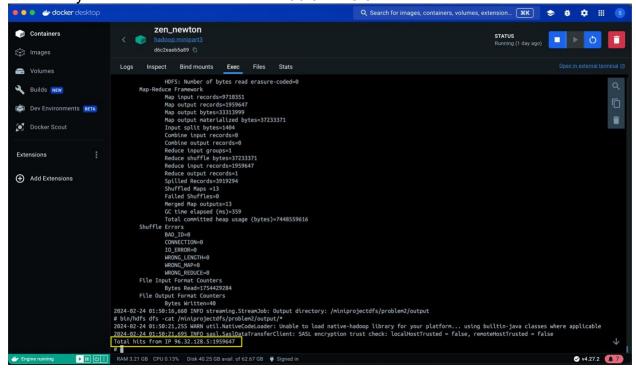Here are the problems we addressed along with the solutions:

Problem 1 –
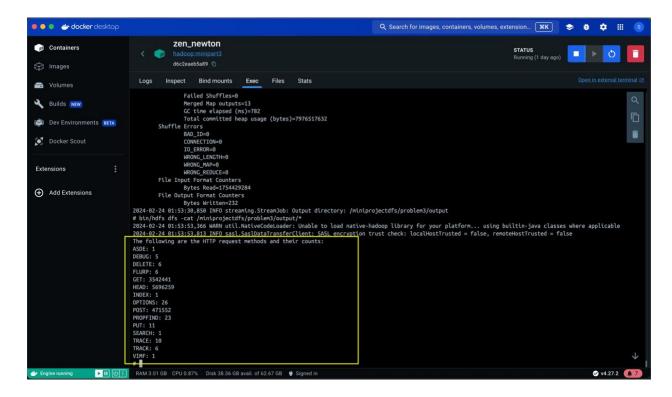How many hits were made to the website directory "/images/smilies/"(including subdirectories and files)?
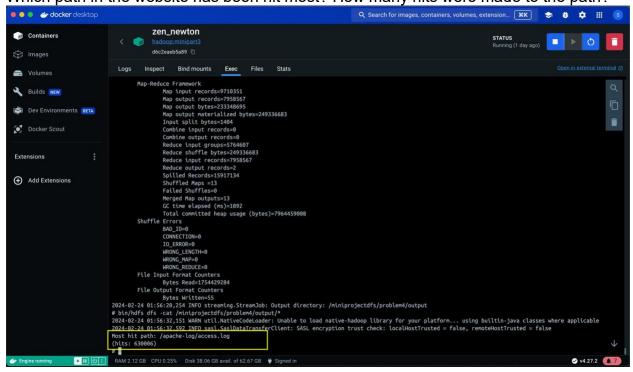
## Problem 2 –
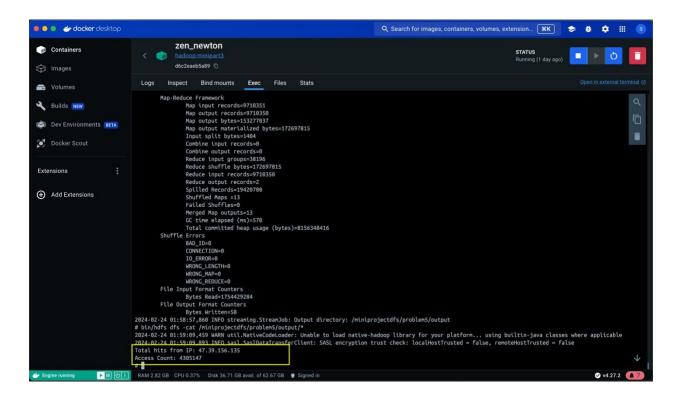## How many hits were made from the IP: 96.32.128.5?



## Problem 3 –
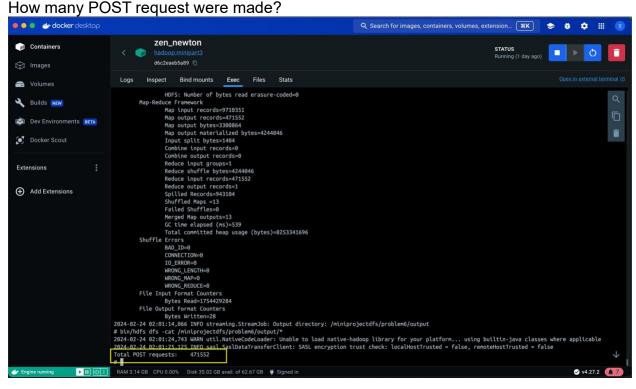## How many HTTP request methods are used in this file? What are they?

## Problem 4 –
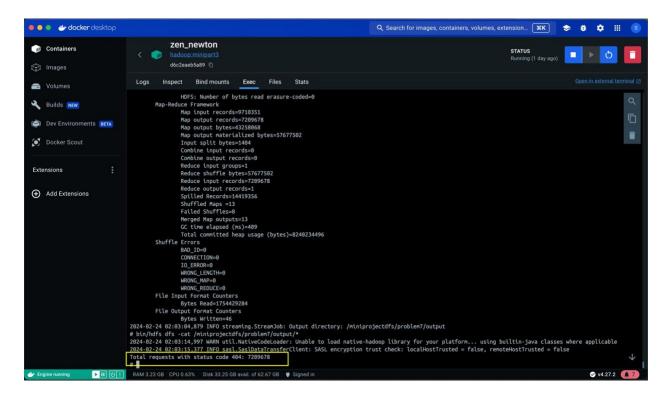Which path in the website has been hit most? How many hits were made to the path?



## Problem 5 –
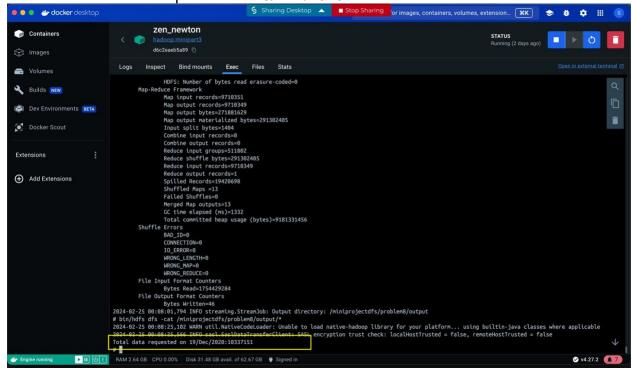Which IP accesses the website most? How many accesses were made by it?

## Problem 6 –
How many POST request were made?



## Problem 7 –
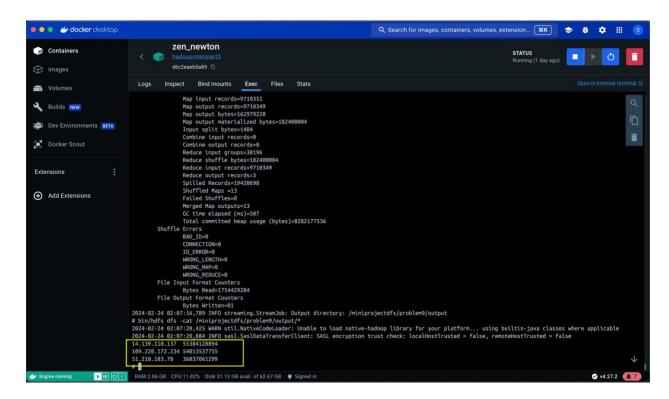How many requests received a 404-status code?

## Problem 8 –
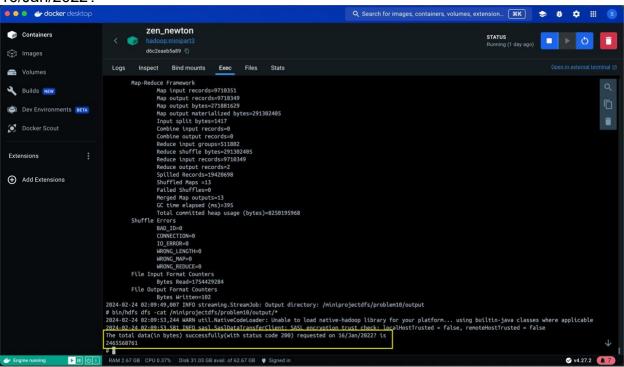How much data was requested on 19/Dec/2020?



## Problem 9 –
List 3 IPs that access the most, and what is the total data flow size of each IP?

## Problem 10 –
How much data (in bytes) was successfully (with status code 200) requested on 16/Jan/2022?

# Distributed mode Hadoop (VM)

We tried to step the project on VM mode but unfortunately, we were unable to connect to the IPs provided. Below is the error we ran into:

```
Last login: Sat Feb 24 16:51:15 on ttys000
[shusritavenugopal@Shusritas-MacBook-Air ~ % ssh student@128.105.146.166
 ssh: connect to host 128.105.146.166 port 22: Network is unreachable
[shusritavenugopal@Shusritas-MacBook-Air ~ % ssh student@128.105.146.166
 ssh: connect to host 128.105.146.166 port 22: Network is unreachable
[shusritavenugopal@Shusritas-MacBook-Air ~ % ssh student@128.105.146.166
 ssh: connect to host 128.105.146.166 port 22: Network is unreachable
 shusritavenugopal@Shusritas-MacBook-Air ~ %
```

# References

Patel, Mehul. "10 Essential Things to Know About Docker - Mehul Patel - Medium." *Medium*,

26 July 2023, medium.com/@nomadicmehul/10-essential-things-to-know-about-docker-

19b827235bfa.

Mehul - Codedamn. "Docker Tutorial 6: Running Ubuntu in Container." *YouTube*, 29 Jan.

2020, www.youtube.com/watch?v=pBTD9bowokA.

Mehul - Codedamn. "Docker Tutorial 11: Dockerfile Introduction." *YouTube*, 8 Mar.

2020, www.youtube.com/watch?v=wGY1WnSPUOk.

*Running Hadoop on Ubuntu Linux (Single-Node Cluster)*. www.michael-

noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster.

*How to Fix the Error JAVAHOME Is Not Set and Could Not Be Found After Hadoop Installation

| Saturn Cloud Blog*. 26 Oct. 2023, https://saturncloud.io/blog/how-to-fix-the-error-

javahome-is-not-set-and-could-not-be-found-after-hadoop-installation/