**Electrical and Computer Engineering, Purdue University Northwest**
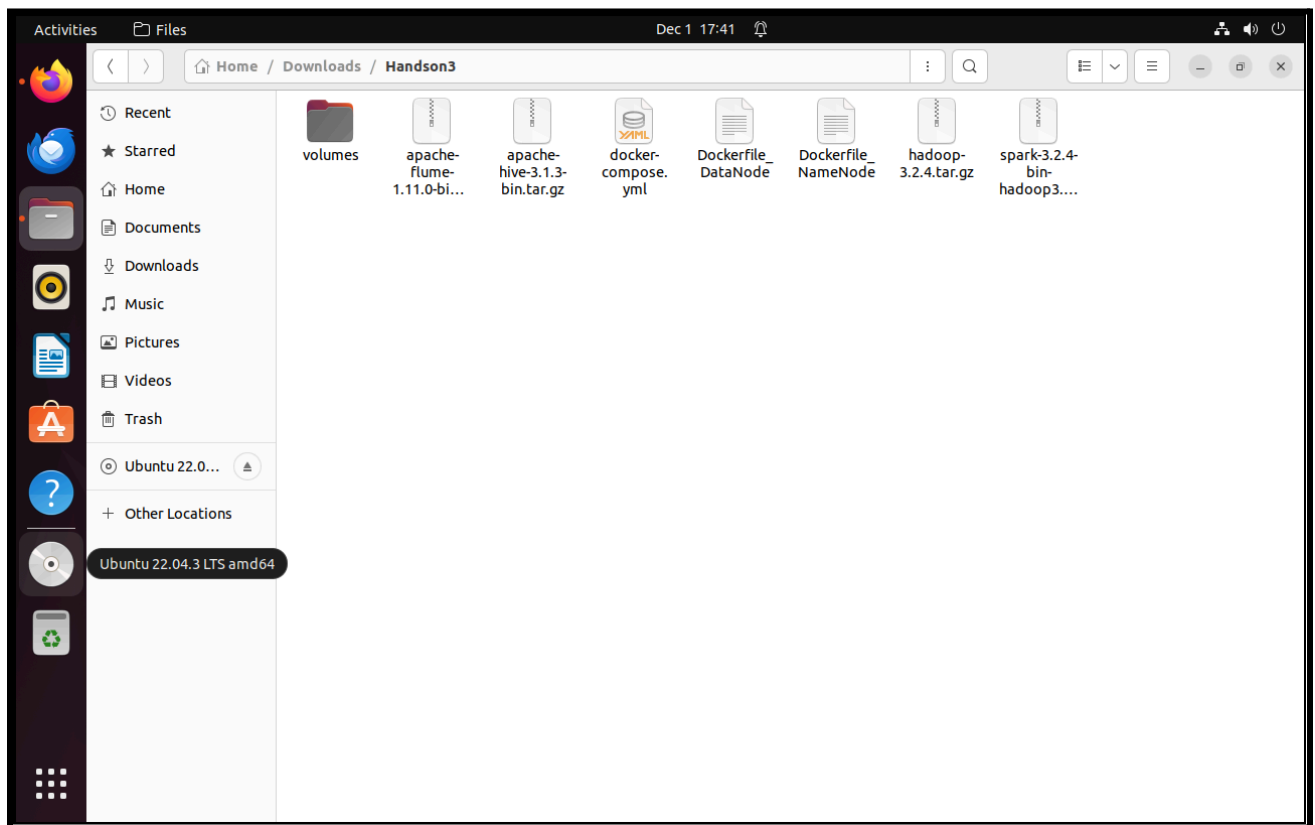**Big Data (ECE49500/ECE59500)**
**Assignment 3**
**Name - Deeksha Hareesha Kulal**

Task 1 [10 points] Hadoop, Spark, Hive, and Flume set-up. You must have to add relevant screenshots in the report for your work to get full credits for the task.

1. Download Handson3.zip file and extract the Handson3 folder from it.

2. Open the terminal to check-in to Handson3 folder and execute command: sudo docker compose up -d

3. Execute the command: sudo docker images

```
deeksha@deeksha-virtual-machine:~/Downloads/Handson3$ sudo docker images
REPOSITORY            TAG       IMAGE ID       CREATED         SIZE
handson3-datanode1    latest    5e5bb204c377   2 minutes ago   2.4GB
handson3-datanode2    latest    5a2e1c63b336   2 minutes ago   2.4GB
handson3-datanode3    latest    5f6958fe1467   2 minutes ago   2.4GB
handson3-namenode     latest    525f5092e297   3 minutes ago   1.85GB
```

4. Open terminal tabs to access containers' shells. Run following command at each tab to access a container's shell. Replace with the actual container's name. sudo docker exec -it $(sudo docker ps --filter "name=" --format "{{.ID}}") /bin/bash

```
                               hadoop@957b41e021c8: ~

deeksha@deeksha-virtual-machine:~/Downloads/Handson3$ sudo docker exec -it $(sudo docker ps --filter "name=namenode" --format "{{.ID
}}") /bin/bash
root@957b41e021c8:/# /volumes/setup.sh
root@957b41e021c8:/# su - hadoop
hadoop@957b41e021c8:~$ hdfs namenode -format
WARNING: /usr/local/hadoop/logs does not exist. Creating.
2023-12-01 22:50:37,631 INFO namenode.NameNode: STARTUP_MSG:
/************************************************************
STARTUP_MSG: Starting NameNode
STARTUP_MSG:   host = namenode/10.9.0.5
STARTUP_MSG:   args = [-format]
STARTUP_MSG:   version = 3.2.4
STARTUP_MSG:   classpath = /usr/local/hadoop/etc/hadoop:/usr/local/hadoop/share/hadoop/common/lib/httpclient-4.5.13.jar:/usr/local/h
adoop/share/hadoop/common/lib/kerby-util-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/kerb-identity-1.0.1.jar:/usr/local/hado
op/share/hadoop/common/lib/jackson-core-asl-1.9.13.jar:/usr/local/hadoop/share/hadoop/common/lib/failureaccess-1.0.jar:/usr/local/ha
doop/share/hadoop/common/lib/commons-compress-1.21.jar:/usr/local/hadoop/share/hadoop/common/lib/kerb-core-1.0.1.jar:/usr/local/hado
op/share/hadoop/common/lib/jsp-api-2.1.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-cli-1.2.jar:/usr/local/hadoop/share/had
oop/common/lib/kerb-common-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-math3-3.1.1.jar:/usr/local/hadoop/share/hadoo
p/common/lib/kerb-admin-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-servlet-9.4.43.v20210629.jar:/usr/local/hadoop/sha
re/hadoop/common/lib/metrics-core-3.2.4.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-text-1.4.jar:/usr/local/hadoop/share/h
adoop/common/lib/htrace-core4-4.1.0-incubating.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-lang3-3.7.jar:/usr/local/hadoop
/share/hadoop/common/lib/commons-net-3.6.jar:/usr/local/hadoop/share/hadoop/common/lib/kerby-asn1-1.0.1.jar:/usr/local/hadoop/share/
hadoop/common/lib/jackson-jaxrs-1.9.13.jar:/usr/local/hadoop/share/hadoop/common/lib/jackson-databind-2.10.5.1.jar:/usr/local/hadoop
/share/hadoop/common/lib/jackson-xc-1.9.13.jar:/usr/local/hadoop/share/hadoop/common/lib/nimbus-jose-jwt-9.8.1.jar:/usr/local/hadoop
/share/hadoop/common/lib/jackson-core-2.10.5.jar:/usr/local/hadoop/share/hadoop/common/lib/jsr311-api-1.1.1.jar:/usr/local/hadoop/sh
are/hadoop/common/lib/json-smart-2.4.7.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-logging-1.1.3.jar:/usr/local/hadoop/sha
re/hadoop/common/lib/reload4j-1.2.18.3.jar:/usr/local/hadoop/share/hadoop/common/lib/asm-5.0.4.jar:/usr/local/hadoop/share/hadoop/co
mmon/lib/jersey-server-1.19.jar:/usr/local/hadoop/share/hadoop/common/lib/javax.servlet-api-3.1.0.jar:/usr/local/hadoop/share/hadoop
/common/lib/curator-recipes-2.13.0.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-configuration2-2.1.1.jar:/usr/local/hadoop/
share/hadoop/common/lib/kerb-server-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/kerb-crypto-1.0.1.jar:/usr/local/hadoop/shar
e/hadoop/common/lib/animal-sniffer-annotations-1.17.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-io-2.8.0.jar:/usr/local/ha
doop/share/hadoop/common/lib/jackson-mapper-asl-1.9.13.jar:/usr/local/hadoop/share/hadoop/common/lib/protobuf-java-2.5.0.jar:/usr/lo
cal/hadoop/share/hadoop/common/lib/kerb-client-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/stax2-api-4.2.1.jar:/usr/local/ha
doop/share/hadoop/common/lib/javax.activation-api-1.2.0.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-codec-1.11.jar:/usr/lo
cal/hadoop/share/hadoop/common/lib/guava-27.0-jre.jar:/usr/local/hadoop/share/hadoop/common/lib/jersey-core-1.19.jar:/usr/local/hado
op/share/hadoop/common/lib/jsch-0.1.55.jar:/usr/local/hadoop/share/hadoop/common/lib/re2j-1.1.jar:/usr/local/hadoop/share/hadoop/com
mon/lib/jersey-servlet-1.19.jar:/usr/local/hadoop/share/hadoop/common/lib/jcip-annotations-1.0-1.jar:/usr/local/hadoop/share/hadoop/
common/lib/slf4j-reload4j-1.7.35.jar:/usr/local/hadoop/share/hadoop/common/lib/jersey-json-1.19.jar:/usr/local/hadoop/share/hadoop/c
ommon/lib/jul-to-slf4j-1.7.35.jar:/usr/local/hadoop/share/hadoop/common/lib/kerby-pkix-1.0.1.jar:/usr/local/hadoop/share/hadoop/comm
on/lib/kerby-config-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/jaxb-impl-2.2.3-1.jar:/usr/local/hadoop/share/hadoop/common/
```

5. At namenode container, execute the following commands: /volumes/setup.sh (Note: Run this command when you access the container first time) su - hadoop hdfs namenode -format (Note: Run this command only when you access the container first time) /volumes/start.sh

```
2023-12-01 22:50:38,172 INFO namenode.NameNode: Caching file names occurring more than 10 times
2023-12-01 22:50:38,176 INFO snapshot.SnapshotManager: Loaded config captureOpenFiles: false, skipCaptureAccessTimeOnlyChange: false
, snapshotDiffAllowSnapRootDescendant: true, maxSnapshotLimit: 65536
2023-12-01 22:50:38,178 INFO snapshot.SnapshotManager: SkipList is disabled
2023-12-01 22:50:38,181 INFO util.GSet: Computing capacity for map cachedBlocks
2023-12-01 22:50:38,181 INFO util.GSet: VM type       = 64-bit
2023-12-01 22:50:38,181 INFO util.GSet: 0.25% max memory 1.7 GB = 4.4 MB
2023-12-01 22:50:38,181 INFO util.GSet: capacity      = 2^19 = 524288 entries
2023-12-01 22:50:38,186 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.window.num.buckets = 10
2023-12-01 22:50:38,186 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.users = 10
2023-12-01 22:50:38,186 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.windows.minutes = 1,5,25
2023-12-01 22:50:38,189 INFO namenode.FSNamesystem: Retry cache on namenode is enabled
2023-12-01 22:50:38,189 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry expiry time is 600
000 millis
2023-12-01 22:50:38,191 INFO util.GSet: Computing capacity for map NameNodeRetryCache
2023-12-01 22:50:38,191 INFO util.GSet: VM type       = 64-bit
2023-12-01 22:50:38,192 INFO util.GSet: 0.029999999329447746% max memory 1.7 GB = 539.6 KB
2023-12-01 22:50:38,192 INFO util.GSet: capacity      = 2^16 = 65536 entries
2023-12-01 22:50:38,214 INFO namenode.FSImage: Allocated new BlockPoolId: BP-1923238440-10.9.0.5-1701471038209
2023-12-01 22:50:38,225 INFO common.Storage: Storage directory /usr/local/hadoop/data/nameNode has been successfully formatted.
2023-12-01 22:50:38,249 INFO namenode.FSImageFormatProtobuf: Saving image file /usr/local/hadoop/data/nameNode/current/fsimage.ckpt_
0000000000000000000 using no compression
2023-12-01 22:50:38,308 INFO namenode.FSImageFormatProtobuf: Image file /usr/local/hadoop/data/nameNode/current/fsimage.ckpt_0000000
000000000000 of size 401 bytes saved in 0 seconds .
2023-12-01 22:50:38,316 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
2023-12-01 22:50:38,331 INFO namenode.FSNamesystem: Stopping services started for active state
2023-12-01 22:50:38,332 INFO namenode.FSNamesystem: Stopping services started for standby state
2023-12-01 22:50:38,335 INFO namenode.FSImage: FSImageSaver clean checkpoint: txid=0 when meet shutdown.
2023-12-01 22:50:38,335 INFO namenode.NameNode: SHUTDOWN_MSG:
/************************************************************
SHUTDOWN_MSG: Shutting down NameNode at namenode/10.9.0.5
************************************************************/
hadoop@957b41e021c8:~$ /volumes/start.sh
248 ResourceManager
154 NameNode
302 Jps
hadoop@957b41e021c8:~$
```

6. At each datanode container, execute the following commands: /volumes/setup.sh (Note: Run this command only when you access the container first time) su - hadoop /volumes/start.sh
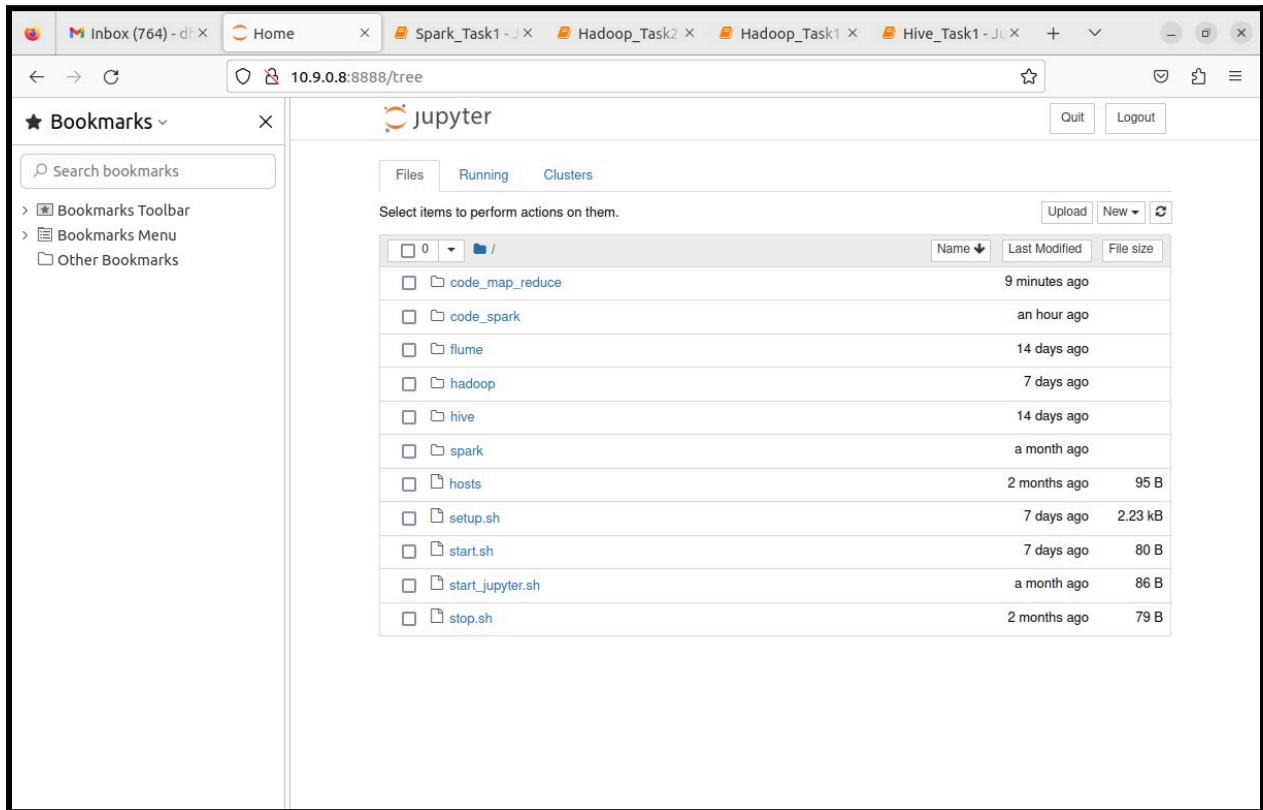
```
deeksha@deeksha-virtual-machine:~/Downloads/Handson3$ sudo docker exec -it $(sudo docker ps --filter "name=datanode1" --format "{{.I
D}}") /bin/bash
[sudo] password for deeksha:
root@927ba7da9af9:/# /volumes/setup.sh
root@927ba7da9af9:/# su - hadoop
hadoop@927ba7da9af9:~$ /volumes/start.sh
WARNING: /usr/local/hadoop/logs does not exist. Creating.
203 NodeManager
108 DataNode
303 Jps
hadoop@927ba7da9af9:~$
```

)

7. Go to any datanode's shell, (e.g. datanode1) and execute the following commands: cd /volumes pip3 install findspark chmod +x start_jupyter.sh /volumes/start_jupyter.sh

8. Copy the Jupyter URL and paste it into a browser in your host machine or VM. Replace the IP 127.0.0.1 with the IP of the datanode. It will bring the Jupyter Notebook. Go to code_map_reduce directory, and complete Task1.pynb file and submit the Notebook file.



**Task 2 [15 points] Go to code_map_reduce directory and complete the tasks given in Hadoop_Task1.ipynb file and submit the Notebook file.**
File submitted

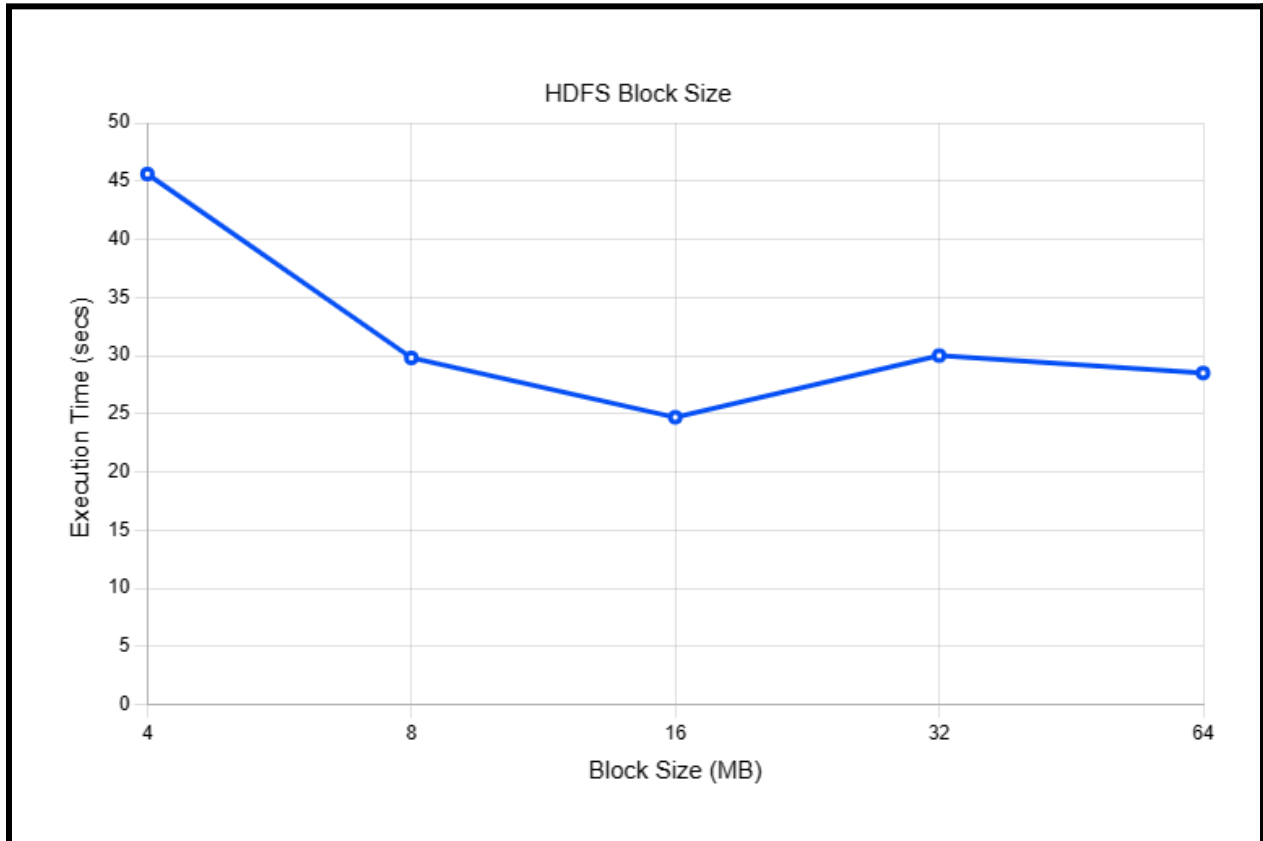**Task 3 [15 points] Complete the tasks given in Hadoop_Task2.ipynb file, submit the Notebook file, and provide the comparison charts for execution times for following cases:**
File submitted
**a. Number of workers = 3, HDFS block size = 4MB, 8MB, 16MB, 32MB, and 64MB (File used - books64.txt)**

| HDFS Block Size | Execution time |
|---|---|
| 4MB | **45.6 secs** |
| 8MB | **29.8 secs** |
| 16MB | **24.7 secs** |

| | |
|---|---|
| 32MB | **30 secs** |
| 64MB | **28.5 secs** |



**b. HDFS block size = 64MB, number of workers = 1, 2, and 3 (Optional for ECE 49500) (File used - books32.txt hence block size of 32MB taken)**

| No of Workers | Execution Time |
|---|---|
| 1 | **24.8 secs** |
| 2 | **22.5 secs** |
| 3 | **21.5 secs** |

Varying the Number of Workers - HFDS Block size (64MB)

**c. HDFS block size = 64MB, number of workers = 3, number of reducers = 1, 2, and 3 (File used - books64.txt)**

| No of Reducers | Execution Time |
|---|---|
| 1 | **30.5 secs** |
| 2 | **29.5 secs** |
| 3 | **29.6 secs** |

Varying the Number of Reducers - HFDS Block size (64MB)

**Task 4 [10 points] Complete the tasks given in Hive_Task1.ipynb file and submit the Notebook file.**
File submitted

**Task 5 [10 points] Complete the tasks given in Spark_Task1.ipynb file, submit the Notebook file, and provide the comparison charts for execution times the following cases:**
File submitted
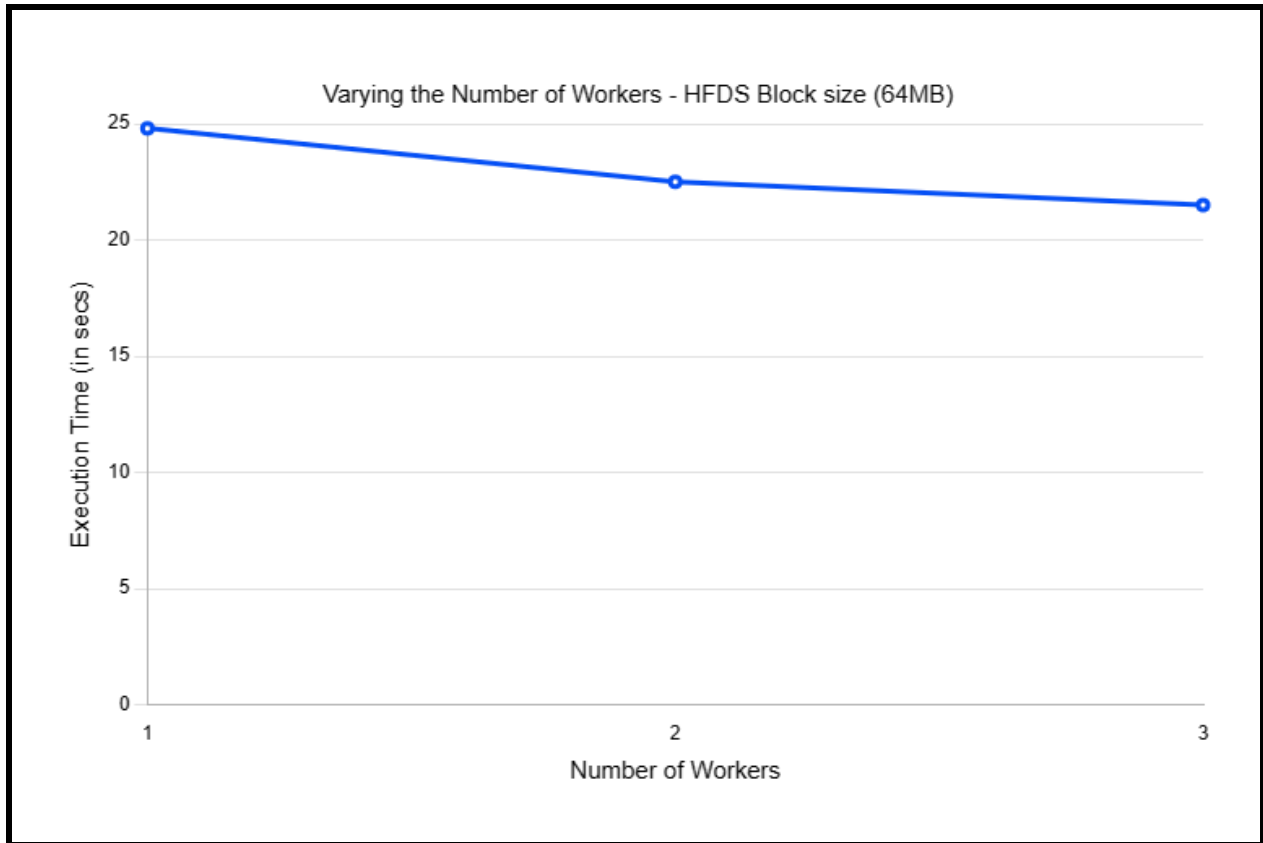**a. Number of workers = 3, HDFS block size = 4MB, 8MB, 16MB, 32MB, and 64MB (File used - books32.txt hence block size of 64 MB not considered)**

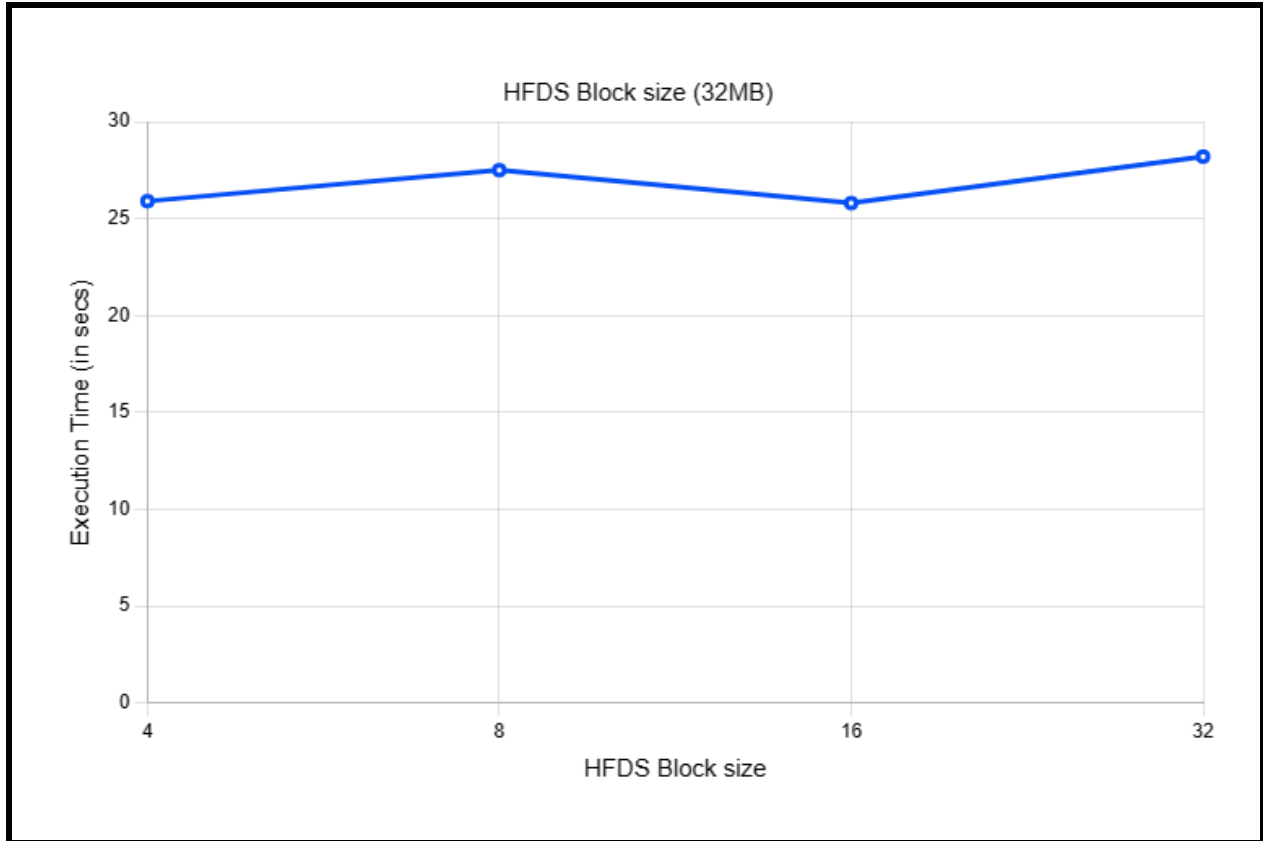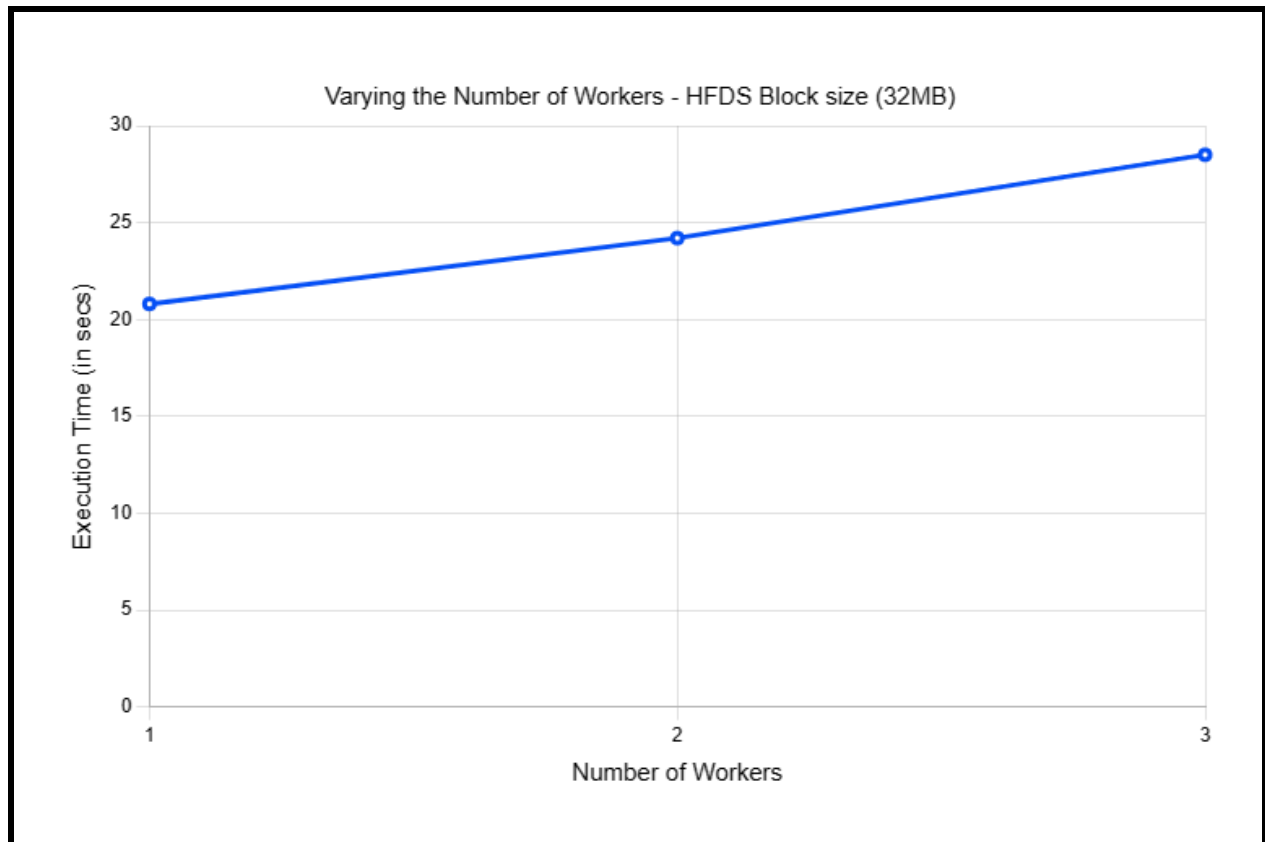| HDFS Block Size | Execution time |
|---|---|
| 4MB | **25.9 secs** |
| 8MB | **27.5 secs** |
| 16MB | **25.8 secs** |
| 32MB | **28.2 secs** |

**b. HDFS block size = 64MB, number of workers = 1, 2, and 3 (Optional for ECE 49500) (HDFS block size of 32MB used with books32.txt file)**

| No of Workers | Execution Time |
|---|---|
| 1 | **20.8 secs** |
| 2 | **24.2 secs** |
| 3 | **28.5 secs** |

Varying the Number of Workers - HFDS Block size (32MB)

**Task 6 [10 points] Go to code_spark directory and complete the tasks given in Spark_Task2.ipynb file and submit the Notebook file.**

File submitted