# Đại học Khoa học Tự nhiên – ĐHQG TP.HCM

----------🙞🙜----------

# Course: DỮ LIỆU LỚN

## Lab 1: A Gentle Introduction to Hadoop

**Class: 21KHMT1**

**Teacher:**

Mrs. Nguyễn Ngọc Thảo

Mr. Bùi Huỳnh Trung Nam

Mr. Đỗ Trọng Lễ

**Student information:**

21127456 – Võ Cao Trí

21127608 – Trần Trung Hiếu
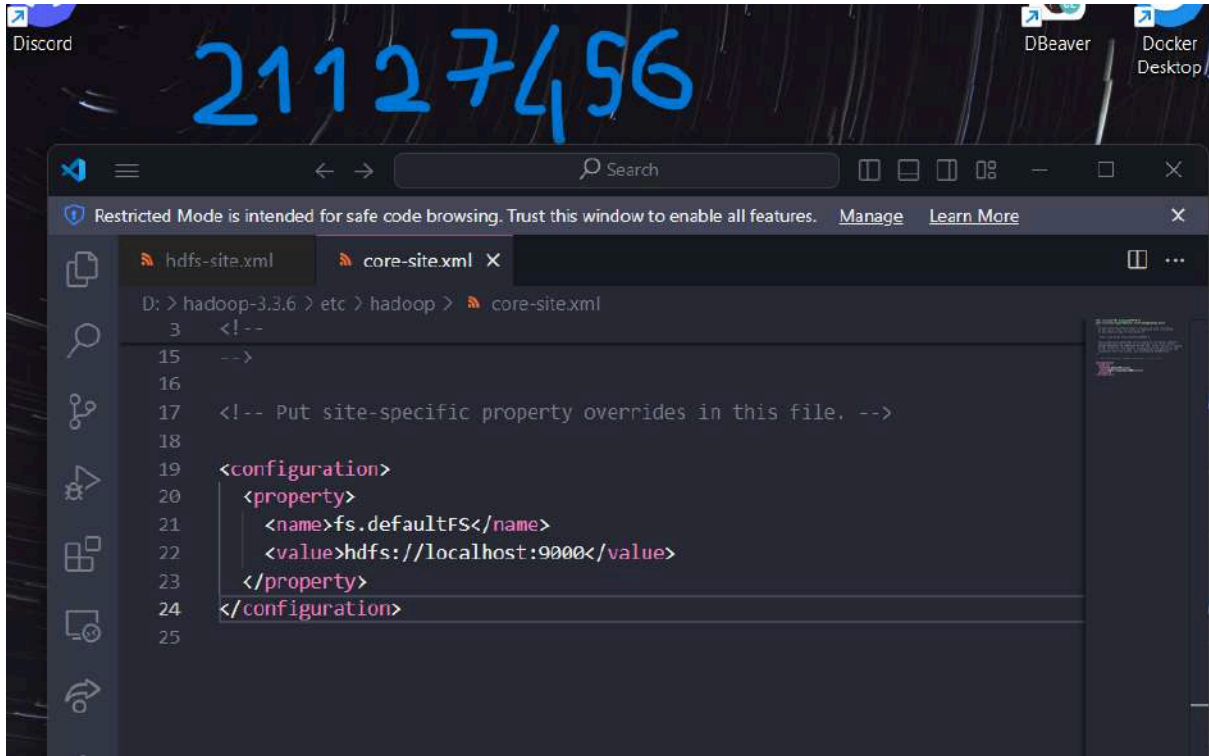
21127668 – Đinh Quang Phong

Assign task:

| NO | Name | Task | Completion level |
| --- | --- | --- | --- |
| 1 | Tri, Hieu, Phong | Setting up Single-node Hadoop Cluster | 100% |
| 2 | Hieu, Tri | Paper Reading | 100% |
| 3 | Phong | Running a warm-up problem: Word Count | 100% |
| 4 | Hieu, Tri, Phong | Write report | 100% |

# 1. Setting up SNC - Single Node Cluster

    a.   The result of Setting up SNC:
- 21127456
  - Set up core-site.xmr



  - Set up hdfs-site.xml

○ Set up mapred-site.xml



```xml
3    <!--
15   -->
16
17   <!-- Put site-specific property overrides in this file. -->
18
19   <configuration>
20     <property>
21       <name>mapreduce.framework.name</name>
22       <value>yarn</value>
23     </property>
24   </configuration>
25
```
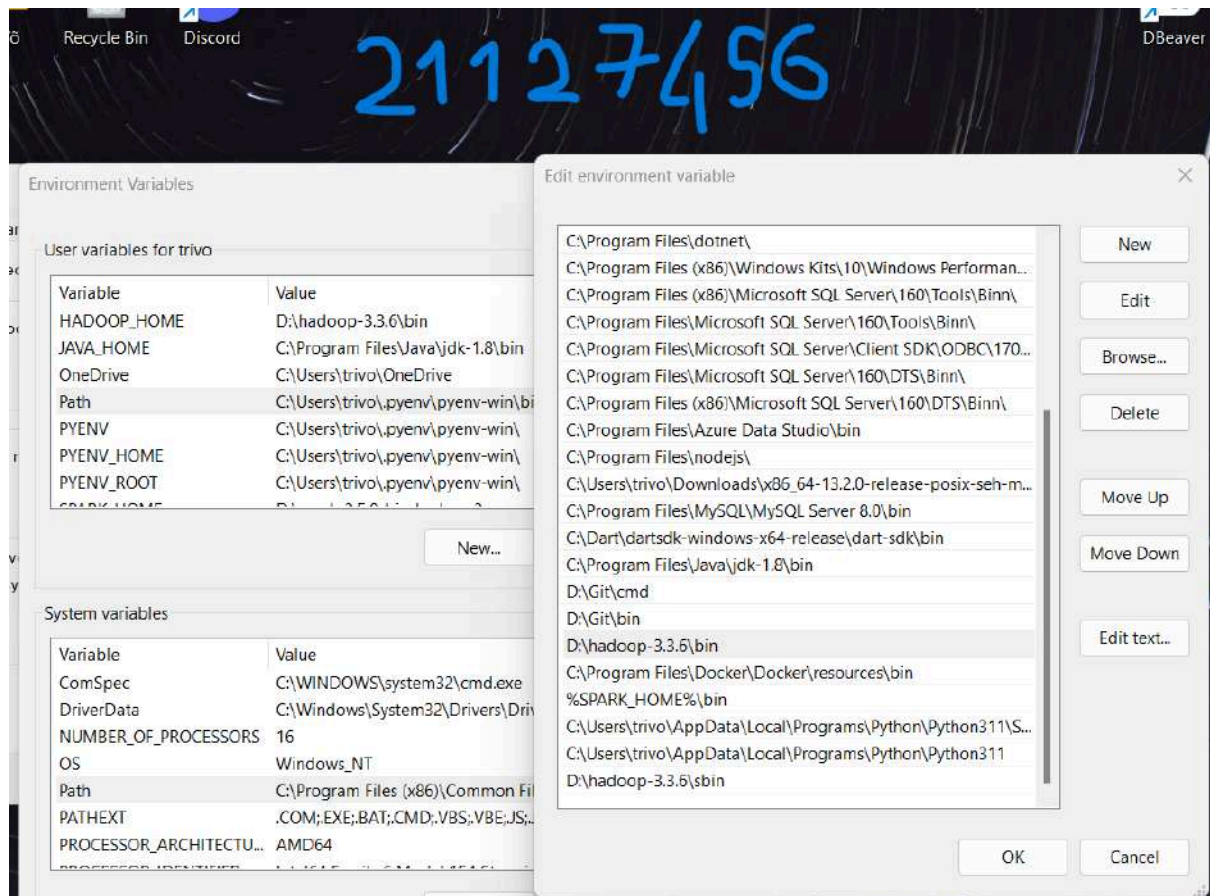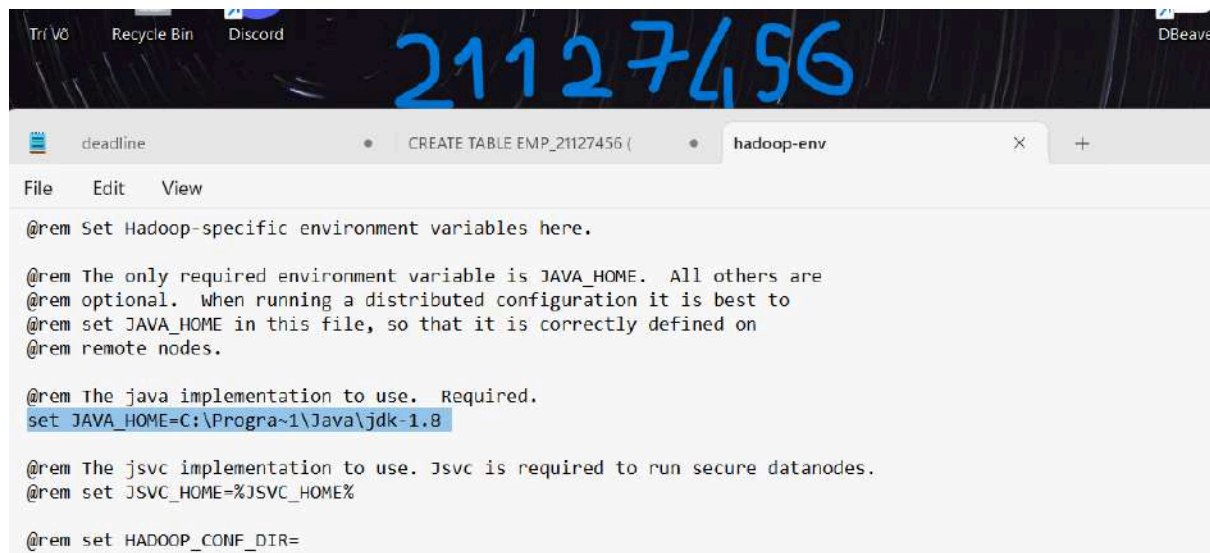
○ Set up yarn-site.xml



```xml
14   -->
15   <configuration>
16     <property>
17       <name>yarn.nodemanager.aux-services</name>
18       <value>mapreduce_shuffle</value>
19     </property>
20     <property>
21       <name>yarn.nodemanager.auxservices.mapreduce.shuffle.class</name>
22     <value>org.apache.hadoop.mapred.ShuffleHandler</value>
23     </property>
24   </configuration>
25
```

○ Set up in Environment Variable

○ Set up JAVA_HOME of hadoop-env.cmd



Result:

○ HDFS namenode format screen

- start-all.cmd screen

○ namenode screen



○ datanode screen

```
68.1.155-1709476691062: 4ms
2024-03-03 21:40:36,978 INFO checker.ThrottledAsyncChecker: Scheduling a check for D:\hadoop-3.3.6\data\datanode
2024-03-03 21:40:36,987 INFO checker.DatasetVolumeChecker: Scheduled health check for volume D:\hadoop-3.3.6\data\datano
de
2024-03-03 21:40:36,990 INFO datanode.VolumeScanner: Now scanning bpid BP-898571220-192.168.1.155-1709476691062 on volum
e D:\hadoop-3.3.6\data\datanode
2024-03-03 21:40:36,992 INFO datanode.VolumeScanner: VolumeScanner(D:\hadoop-3.3.6\data\datanode, DS-9e4b62e0-b263-4cd5-
8797-77ee4ceab5df): finished scanning block pool BP-898571220-192.168.1.155-1709476691062
2024-03-03 21:40:36,995 WARN datanode.DirectoryScanner: dfs.datanode.directoryscan.throttle.limit.ms.per.sec set to valu
e above 1000 ms/sec. Assuming default value of -1
2024-03-03 21:40:36,995 INFO datanode.DirectoryScanner: Periodic Directory Tree Verification scan starting in 7625106ms
with interval of 21600000ms and throttle limit of -1ms/s
2024-03-03 21:40:37,004 INFO datanode.DataNode: Block pool BP-898571220-192.168.1.155-1709476691062 (Datanode Uuid ba947
0d9-bc45-4728-b7a5-d14f4e820978) service to localhost/127.0.0.1:9000 beginning handshake with NN
2024-03-03 21:40:37,009 INFO datanode.VolumeScanner: VolumeScanner(D:\hadoop-3.3.6\data\datanode, DS-9e4b62e0-b263-4cd5-
8797-77ee4ceab5df): no suitable block pools found to scan.  Waiting 1814399981 ms.
2024-03-03 21:40:37,103 INFO datanode.DataNode: Block pool BP-898571220-192.168.1.155-1709476691062 (Datanode Uuid ba947
0d9-bc45-4728-b7a5-d14f4e820978) service to localhost/127.0.0.1:9000 successfully registered with NN
2024-03-03 21:40:37,105 INFO datanode.DataNode: For namenode localhost/127.0.0.1:9000 using BLOCKREPORT_INTERVAL of 2160
0000msecs CACHEREPORT_INTERVAL of 10000msecs Initial delay: 0msecs; heartBeatInterval=3000
2024-03-03 21:40:37,106 INFO datanode.DataNode: Starting IBR Task Handler.
2024-03-03 21:40:37,248 INFO datanode.DataNode: After receiving heartbeat response, updating state of namenode localhost
:9000 to active
2024-03-03 21:40:37,326 INFO datanode.DataNode: Successfully sent block report 0x8be113c163ab0c8b with lease ID 0xf89318
3c42061da8 to namenode: localhost/127.0.0.1:9000,  containing 1 storage report(s), of which we sent 1. The reports had 0
 total blocks and used 1 RPC(s). This took 8 msecs to generate and 68 msecs for RPC and NN processing. Got back one comm
and: FinalizeCommand/5.
2024-03-03 21:40:37,327 INFO datanode.DataNode: Got finalize command for block pool BP-898571220-192.168.1.155-170947669
1062
```
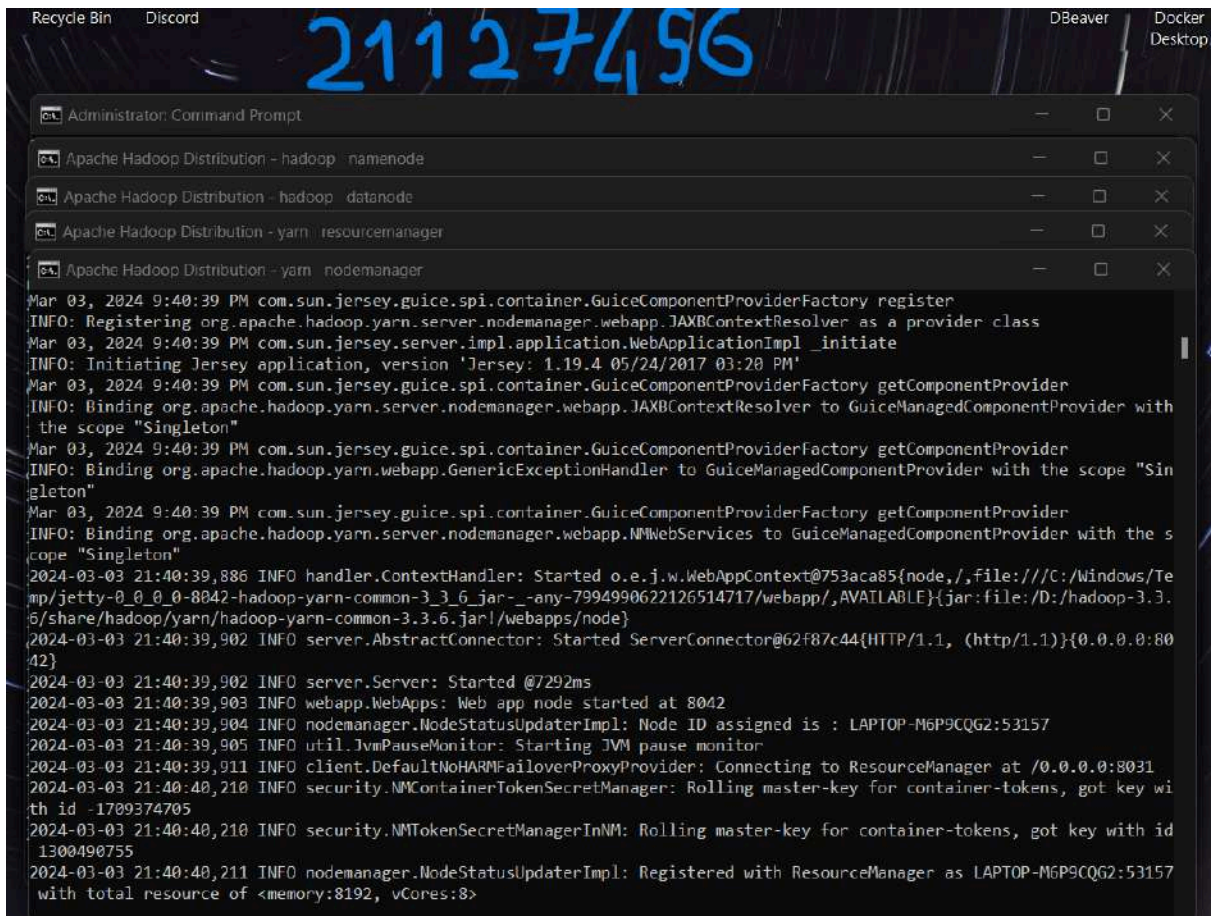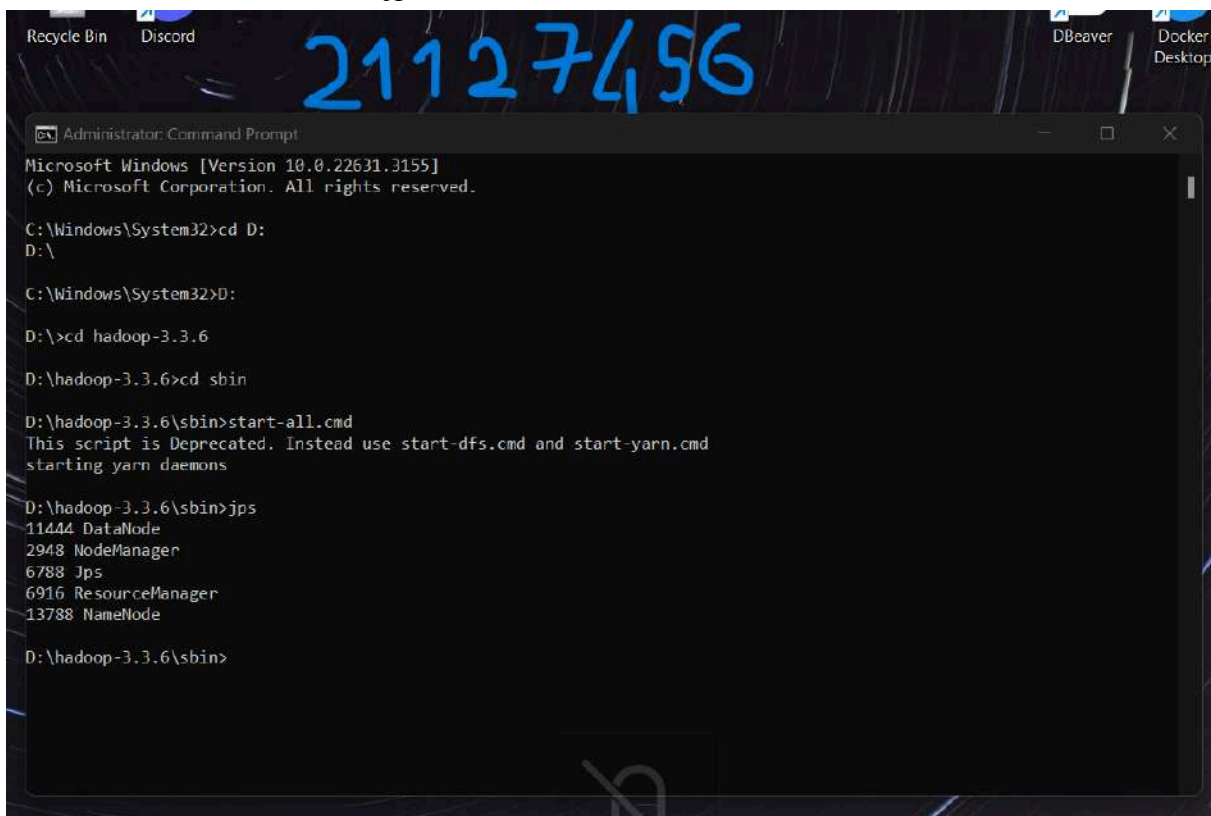
○ resource manager screen

```
2024-03-03 21:40:38,781 INFO pb.RpcServerFactoryPBImpl: Adding protocol org.apache.hadoop.yarn.server.api.ResourceTracke
rPB to the server
2024-03-03 21:40:38,781 INFO ipc.Server: IPC Server listener on 8031: starting
2024-03-03 21:40:38,783 INFO ipc.Server: IPC Server Responder: starting
2024-03-03 21:40:38,787 INFO util.JvmPauseMonitor: Starting JVM pause monitor
2024-03-03 21:40:38,800 INFO ipc.CallQueueManager: Using callQueue: class java.util.concurrent.LinkedBlockingQueue, queu
eCapacity: 5000, scheduler: class org.apache.hadoop.ipc.DefaultRpcScheduler, ipcBackoff: false.
2024-03-03 21:40:38,805 INFO ipc.Server: Listener at 0.0.0.0:8030
2024-03-03 21:40:38,807 INFO ipc.Server: Starting Socket Reader #1 for port 8030
2024-03-03 21:40:38,814 INFO pb.RpcServerFactoryPBImpl: Adding protocol org.apache.hadoop.yarn.api.ApplicationMasterProt
ocolPB to the server
2024-03-03 21:40:38,814 INFO ipc.Server: IPC Server listener on 8030: starting
2024-03-03 21:40:38,814 INFO ipc.Server: IPC Server Responder: starting
2024-03-03 21:40:38,922 INFO ipc.CallQueueManager: Using callQueue: class java.util.concurrent.LinkedBlockingQueue, queu
eCapacity: 5000, scheduler: class org.apache.hadoop.ipc.DefaultRpcScheduler, ipcBackoff: false.
2024-03-03 21:40:38,923 INFO ipc.Server: Listener at 0.0.0.0:8032
2024-03-03 21:40:38,925 INFO ipc.Server: Starting Socket Reader #1 for port 8032
2024-03-03 21:40:38,930 INFO pb.RpcServerFactoryPBImpl: Adding protocol org.apache.hadoop.yarn.api.ApplicationClientProt
ocolPB to the server
2024-03-03 21:40:38,931 INFO ipc.Server: IPC Server Responder: starting
2024-03-03 21:40:38,932 INFO ipc.Server: IPC Server listener on 8032: starting
2024-03-03 21:40:39,313 INFO webproxy.ProxyCA: Created Certificate for OU=YARN-33a70ae1-2d12-4e58-b442-d45c4f57cb0a
2024-03-03 21:40:39,364 INFO recovery.RMStateStore: Storing CA Certificate and Private Key
2024-03-03 21:40:39,364 INFO resourcemanager.ResourceManager: Transitioned to active state
2024-03-03 21:40:40,192 INFO resourcemanager.ResourceTrackerService: NodeManager from node LAPTOP-M6P9CQG2(cmPort: 53157
 httpPort: 8042) registered with capability: <memory:8192, vCores:8>, assigned nodeId LAPTOP-M6P9CQG2:53157
2024-03-03 21:40:40,193 INFO rmnode.RMNodeImpl: LAPTOP-M6P9CQG2:53157 Node Transitioned from NEW to RUNNING
2024-03-03 21:40:40,216 INFO capacity.CapacityScheduler: Added node LAPTOP-M6P9CQG2:53157 clusterResource: <memory:8192,
 vCores:8>
```

○ node manager screen



○ jps screen

- On localhost:9870 of Hadoop File System



- On localhost:8088 is cluster app of hadoop



- 21127608
  - Set up core-site.xml

○ Set up hdfs-site.xml



○ Set up httpfs-site.xml

○   Set up mapred-site.xml



○   Set up yarn-site.xml

○ Set up HADOOP_HOME in Environment Variable



○ Set up JAVA_HOME in Environment Variable

○ Set up path of sbin

○ Set up path of bin



○ Set up JAVA_HOME of hadoop-env.cmd

- Format namenode of hadoop



- Run all file .cmd in sbin folder

○ Yarn nodemanager of sbin



○ Yarn resourcemanager of sbin

- Hadoop datanode of sbin



- Hadoop namenode of sbin

○ Configuration of hadoop system
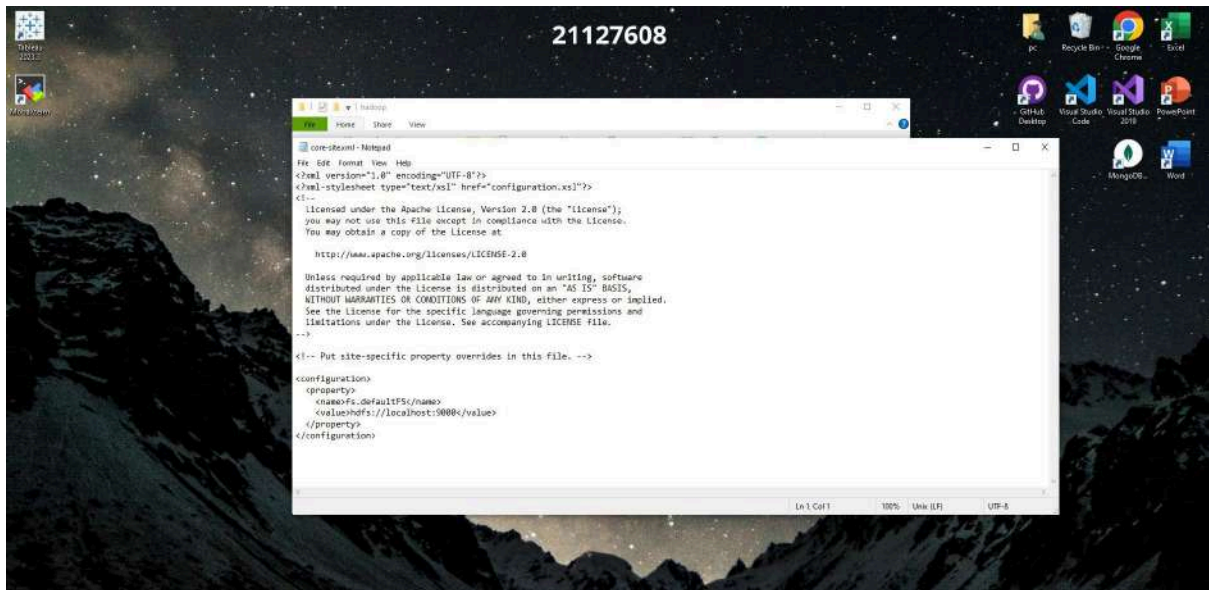


○ On localhost:9870 of Hadoop File System
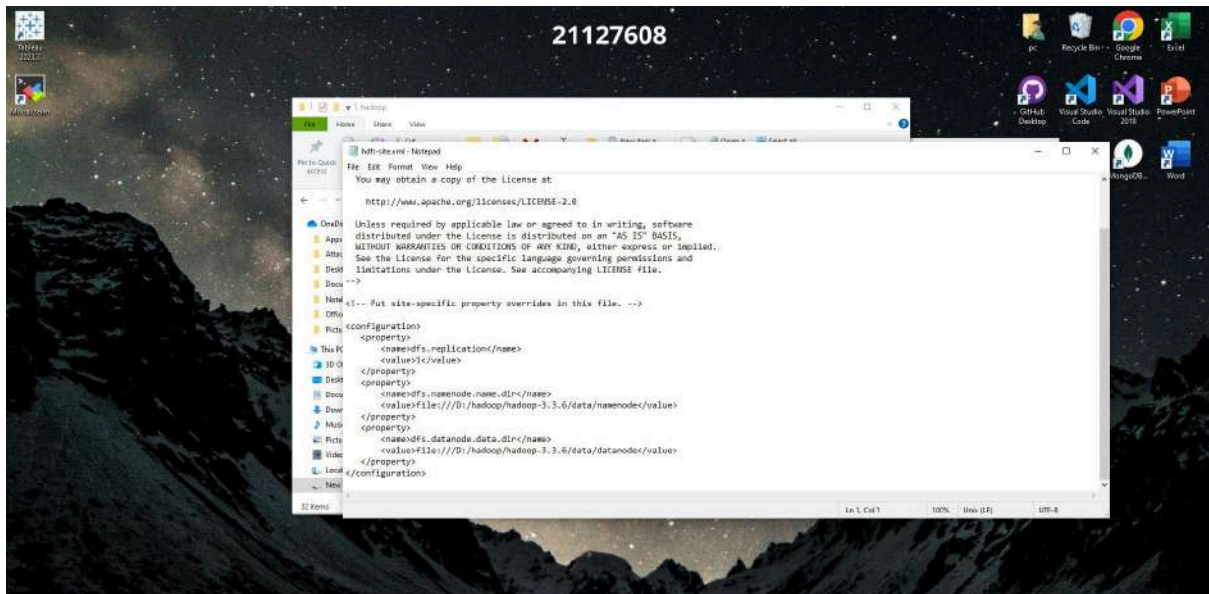
○ On localhost:9870 of Hadoop File System
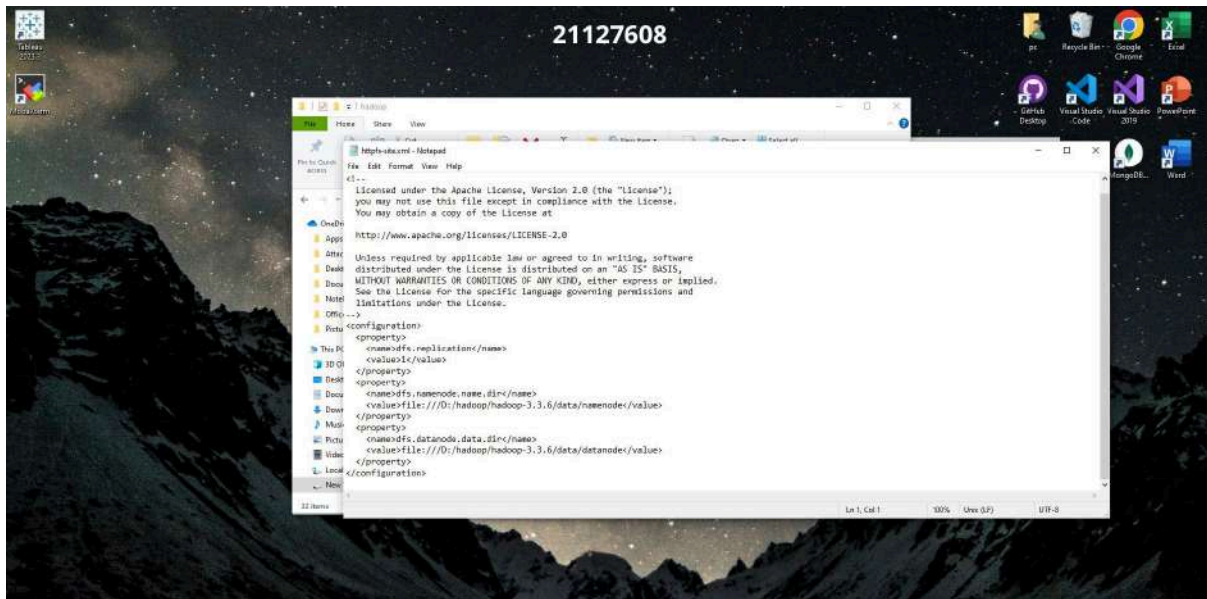


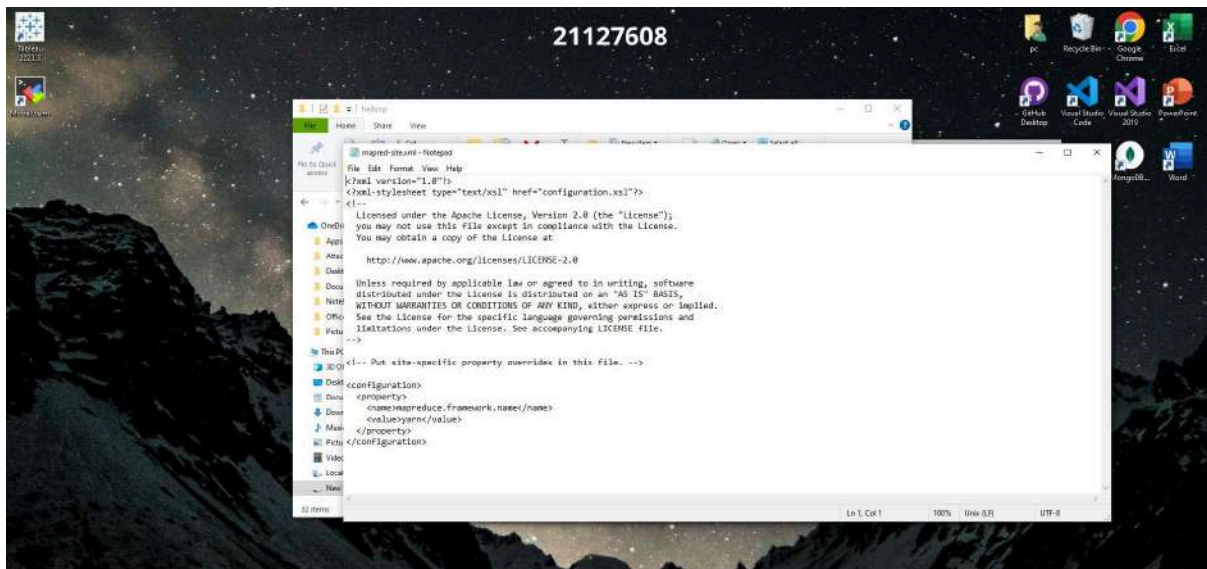○ On localhost:8088 is cluster app of hadoop
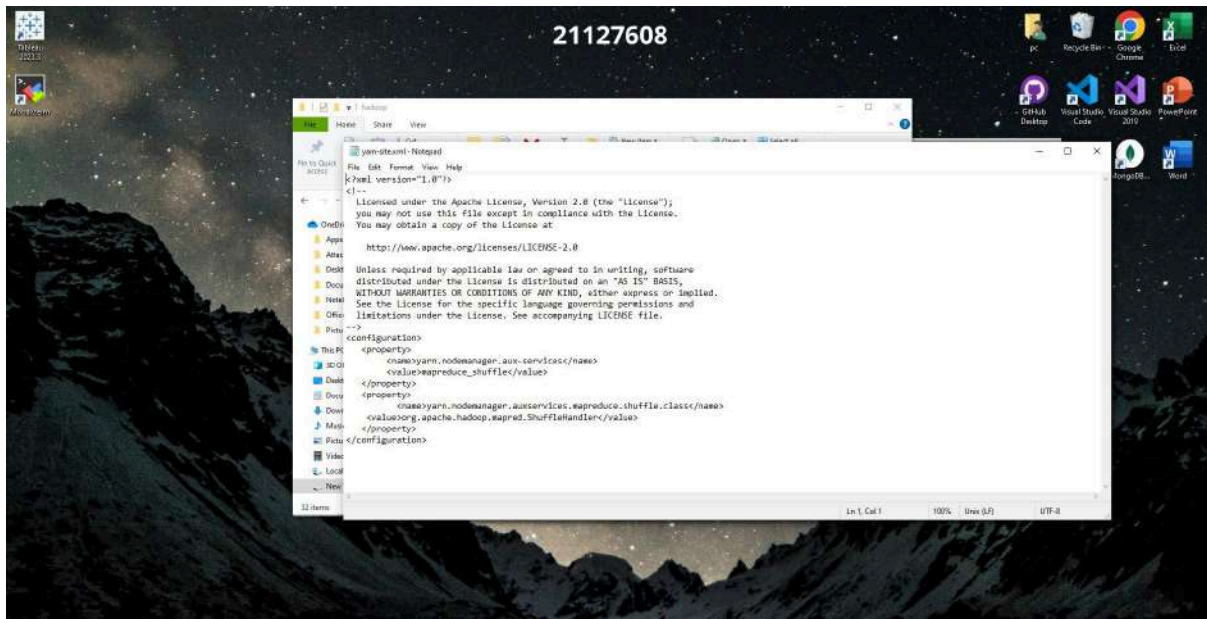
- 21127668
  - Set up core-site.xml
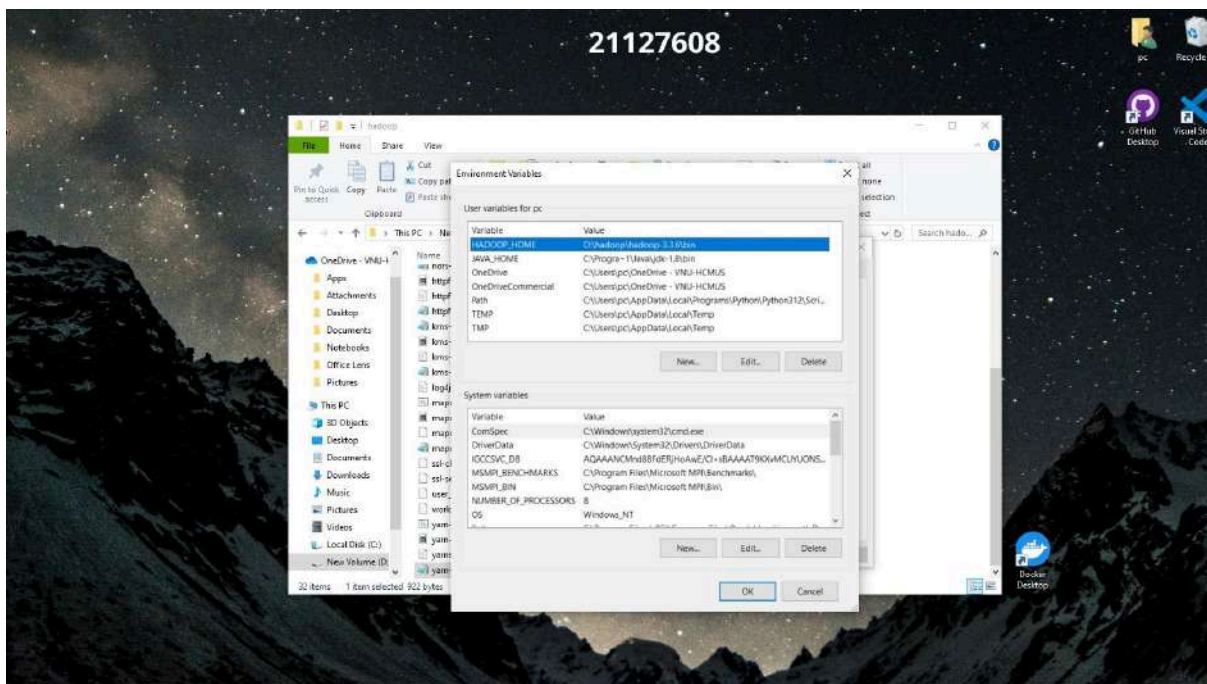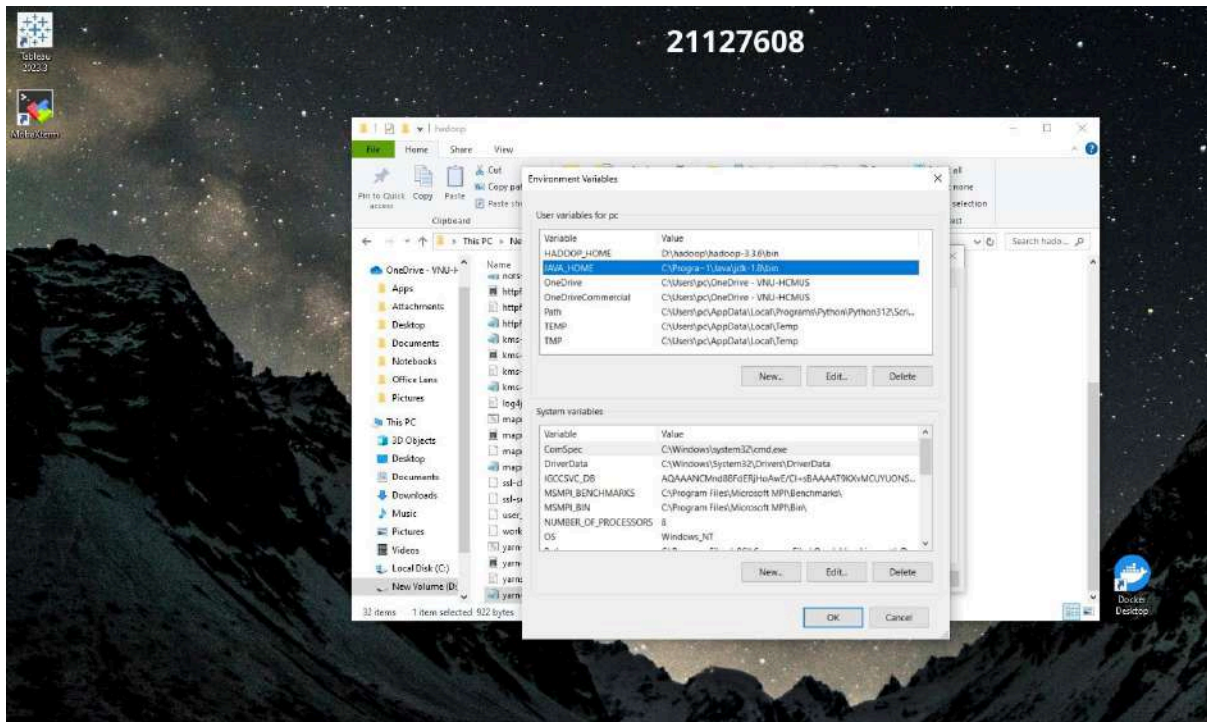
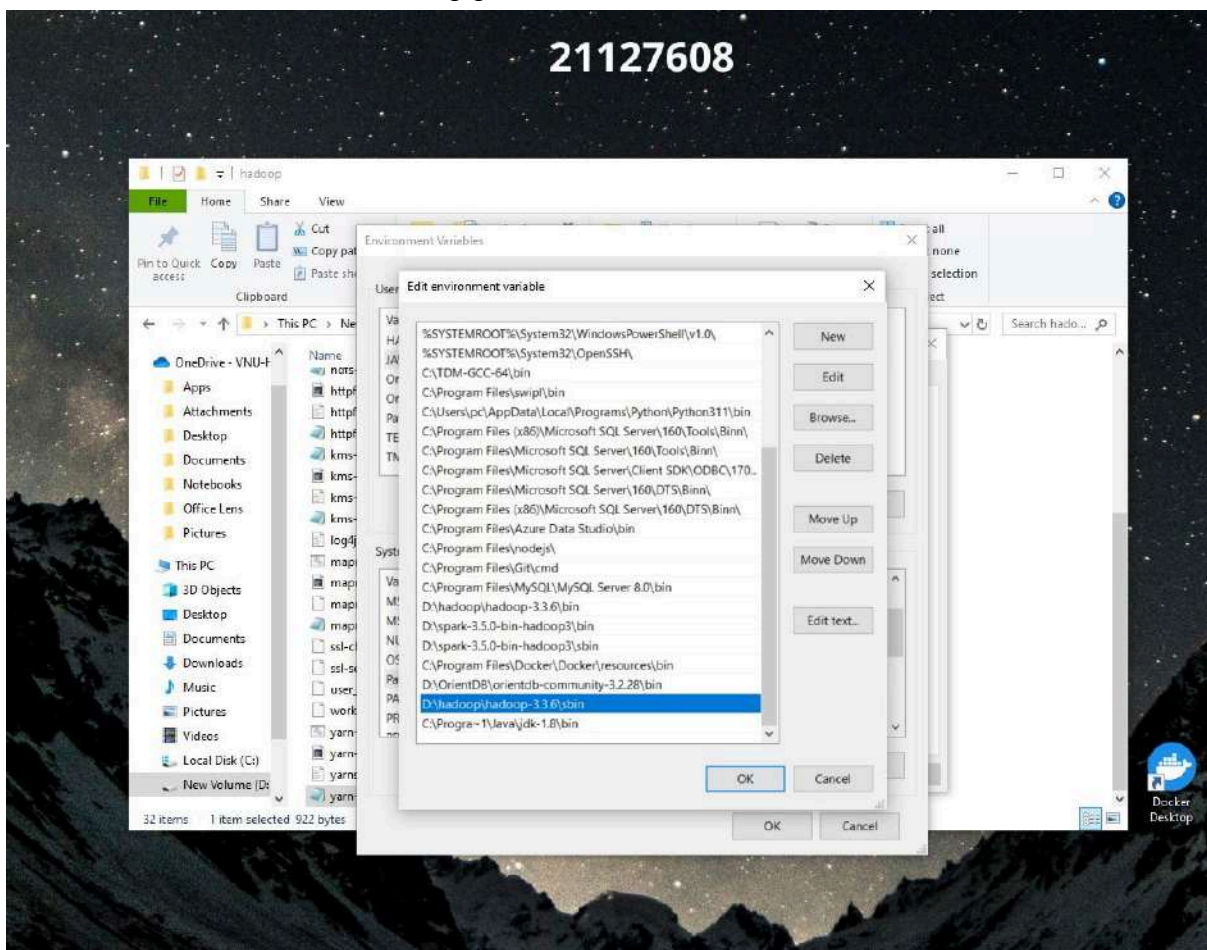○ Set up hdfs-site.xml



○ Set up mapred-site.xml



○ Set up yarn-site.xml

○ Set up in Environment Variable



○ Set up JAVA_HOME of hadoop-env.cmd



Result:

○ HDFS namenode format screen

○ start-all.cmd screen



```
D:\hadoop-3.3.6\sbin>start-all.cmd
This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd
starting yarn daemons

D:\hadoop-3.3.6\sbin>
```

○ namenode screen

○ datanode screen



○ resource manager screen

○ node manager screen



○ jps screen

- On localhost:9870 of Hadoop File System



- On localhost:8088 is cluster app of hadoop



b. The problem and solution

| NO | Problem | Solution |
|---|---|---|
| 1 | Hadoop's bin folder is incomplete, leading to errors during execution and automatic shutdown of the server | Search on stackoverflow and hadoop related websites to add components in bin folder for completeness |
| 2 | The version of Java does not match hadoop leading to | Find the version of Java suitable for Hadoop: Java |

| | ERROR namenode.NameNode: Failed to start namenode. | 1.8.0, The Oracle JDK 8 license changed in April 2019 |
|---|---|---|
| 3 | File yarn.cmd of sbin can not run when executing the start-all.cmd command | Execute the command on the terminal with administrator rights |

# 2. Introduction to MapReduce:

a. How do the input keys-values, the intermediate keys-values, and the output keys-values relate?

- The computation takes a set of input key/value pairs, and produces a set of output key/value pairs. The user of the MapReduce library expresses the computation as two functions: Map and Reduce.
  - Map, written by the user, takes an input pair and produces a set of intermediate key/value pairs. The MapReduce library groups together all intermediate values associated with the same intermediate key I and passes them to the Reduce function.
  - The Reduce function, also written by the user, accepts an intermediate key I and a set of values for that key. It merges together these values to form a possibly smaller set of values. Typically just zero or one output value is produced per Reduce invocation. The intermediate values are supplied to the user's reduce function via an iterator. This allows us to handle lists of values that are too large to fit in memory.
- ➔ In summary, the input keys-values are transformed into intermediate keys-values by the map tasks, and then these intermediate key-value pairs are processed by the reduce tasks to produce the final output keys-values. The relationship between them is characterized by the flow of data through the MapReduce computation pipeline.

b. How does MapReduce deal with node failures?

- The master periodically pings workers to monitor their status. If a worker does not respond for a certain period of time, it is marked as failed by the master.
- Any map tasks that were in progress or completed on the failed worker are reset to their initial idle state and become eligible for rescheduling on other workers. This avoids losing work due to failures.
- Completed reduce tasks do not need to be re-executed since their output is stored in a global file system not local to any worker node. This output is not impacted by worker failures.

- When a map task is re-executed on a new worker after the original worker failed, all reduced tasks are notified so they fetch the updated output rather than potentially outdated data from the failed worker.
- Large scale failures of groups of workers can be tolerated as failed work is simply rescheduled on remaining healthy nodes, allowing jobs to complete even if large portions of the cluster temporarily fail.
- The use of atomic task outputs and failure notifications ensures the overall results meet the expected semantics - they are equivalent to non-faulty sequential execution if tasks are deterministic, or provide reasonable semantics otherwise.

c. What is the meaning and implication of locality? What does it use?

- The workings of the MapReduce framework within a computing environment utilizing the Google File System (GFS), "locality" refers to the principle of executing computational tasks (specifically map tasks in this case) on machines where the required input data is already stored locally, i.e., on the same physical node or within the same network proximity. By leveraging data locality, tasks can read their input directly from local storage rather than requiring data transfer over the network.
- Implication: By processing data locally, MapReduce avoids the need to transfer large datasets across the network, which can be slow and resource-intensive. This leads to:
  - Reduce network traffic: Transferring massive amounts of data across the network can be slow and resource-intensive. Locality minimizes this by processing data on the node where it resides, significantly reducing network bandwidth usage.
  - Improve performance: By avoiding unnecessary data movement, locality leads to faster overall job execution time. This is crucial for large-scale data processing where efficiency is paramount.
  - Increase scalability: By reducing network load, MapReduce can handle larger datasets and workloads more efficiently, contributing to better scalability of the system.
- Usage: The MapReduce scheduler aims to place tasks on nodes with local input replicas. If not possible locally, it tries to schedule tasks on nearby nodes, e.g. same network switch/rack, to improve data locality.

d. Which problem is addressed by introducing a combiner function to the MapReduce model?

- Introducing a combiner function to the MapReduce model addresses the problem of excessive data transfer over the network, especially when there is significant repetition in the intermediate keys produced by each map task.

- In MapReduce, after the map phase, intermediate key-value pairs are shuffled and sorted before being passed to the reduce tasks. However, in scenarios like word counting, where many identical intermediate key-value pairs are generated by map tasks (e.g., <the, 1>), sending all these pairs over the network to a single reduced task can result in substantial network traffic and potentially lead to network congestion and increased processing latency.

- The combiner function acts as a partial aggregation stage before sending data over the network. It reduces the number of key-value pairs with the same key by combining their values (e.g., summing "the" counts from different mappers). This significantly reduces network traffic and improves overall performance.

# 3. Running a warm-up problem: Word Count

- Add dependencies to Maven pom.xml file after creating the Maven project.



- Write Mapper class. In the context of the WordCount program, the input data is a text file. The Mapper takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key/value pairs). Specifically, the Mapper reads a line of text, breaks it into words, and for each word, it emits a key-value pair, where the key is the word and the value is 1.

- Write Reducer class to take the output from the Mapper as input and combine those data tuples into a smaller set of tuples. As the sequence of the name MapReduce implies, the reduce task is always performed after the map job. In the context of the WordCount program, the Reducer takes the input from the Mapper (the key-value pairs), sums up the values for each unique key (the word), and emits a key-value pair, where the key is the word and the value is the total count



- Write WordCount main function, sets up the configuration for the job, specifies the input and output paths, sets the Mapper and Reducer
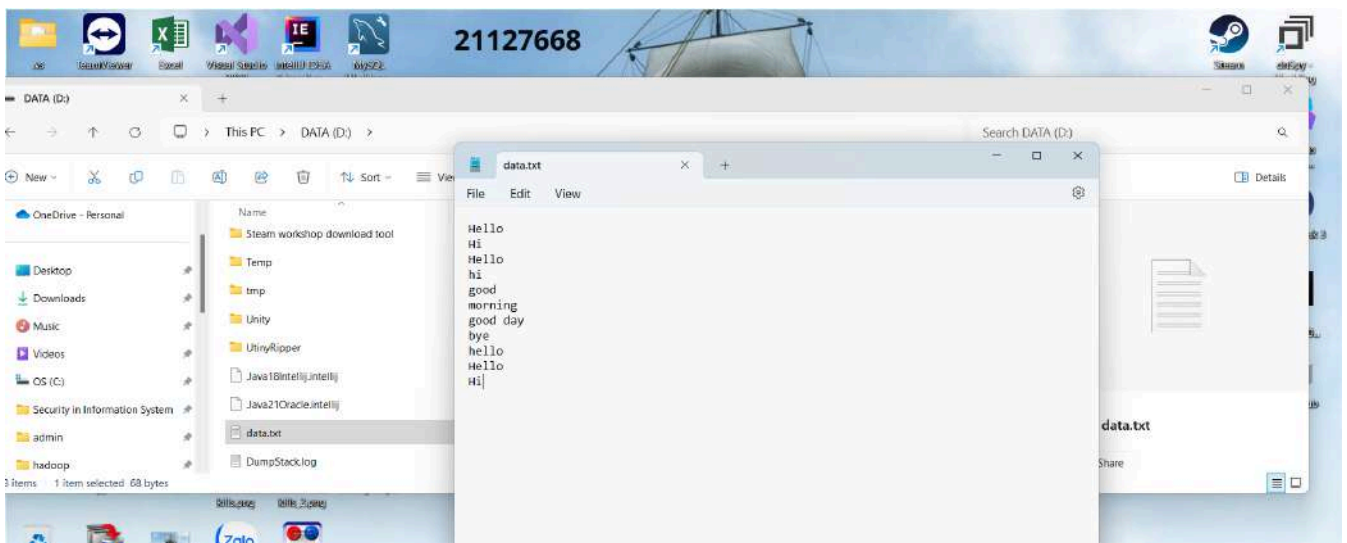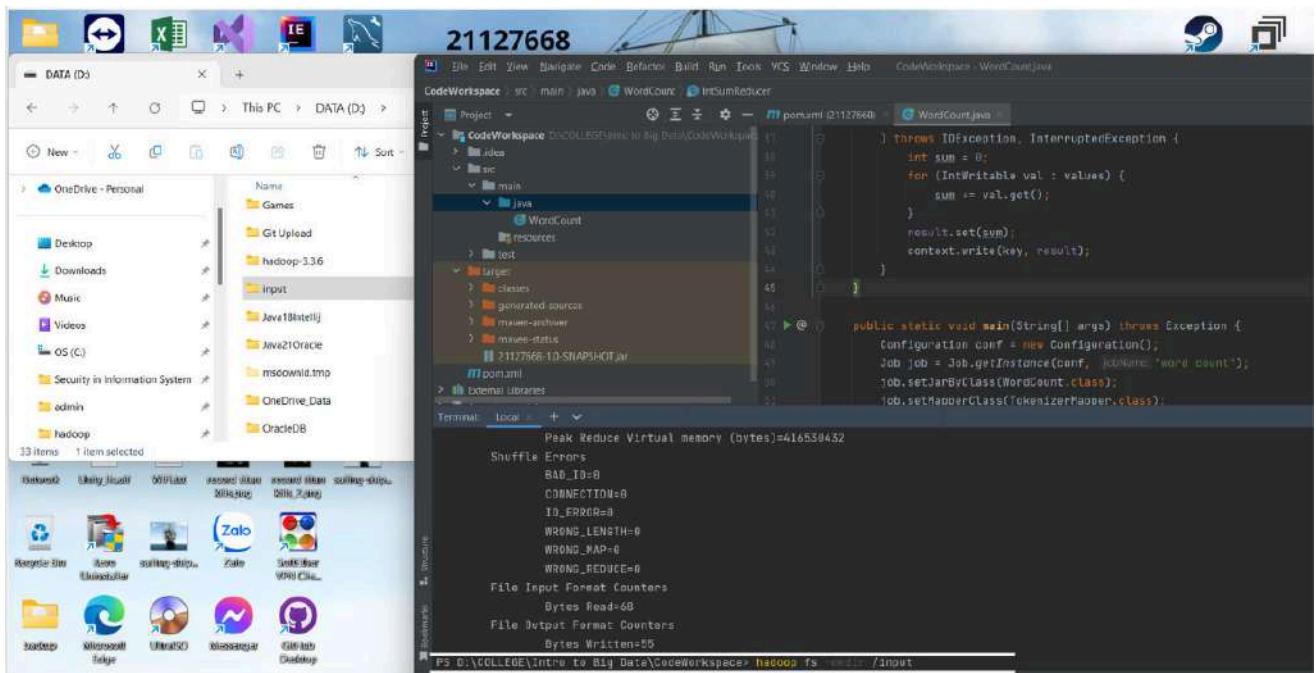
classes, and submits the job to the Hadoop cluster. Once the job is submitted, Hadoop takes care of distributing the data, scheduling and running the map and reduce tasks, and collecting the results.
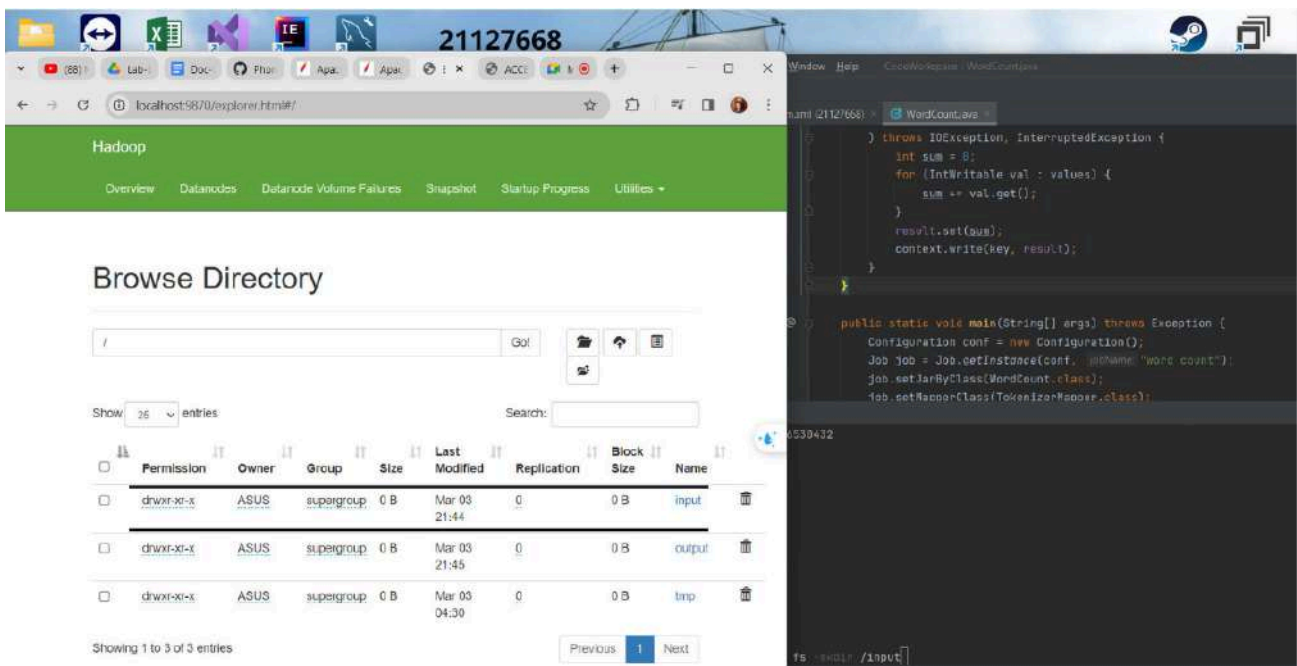


● Create **data.txt** file



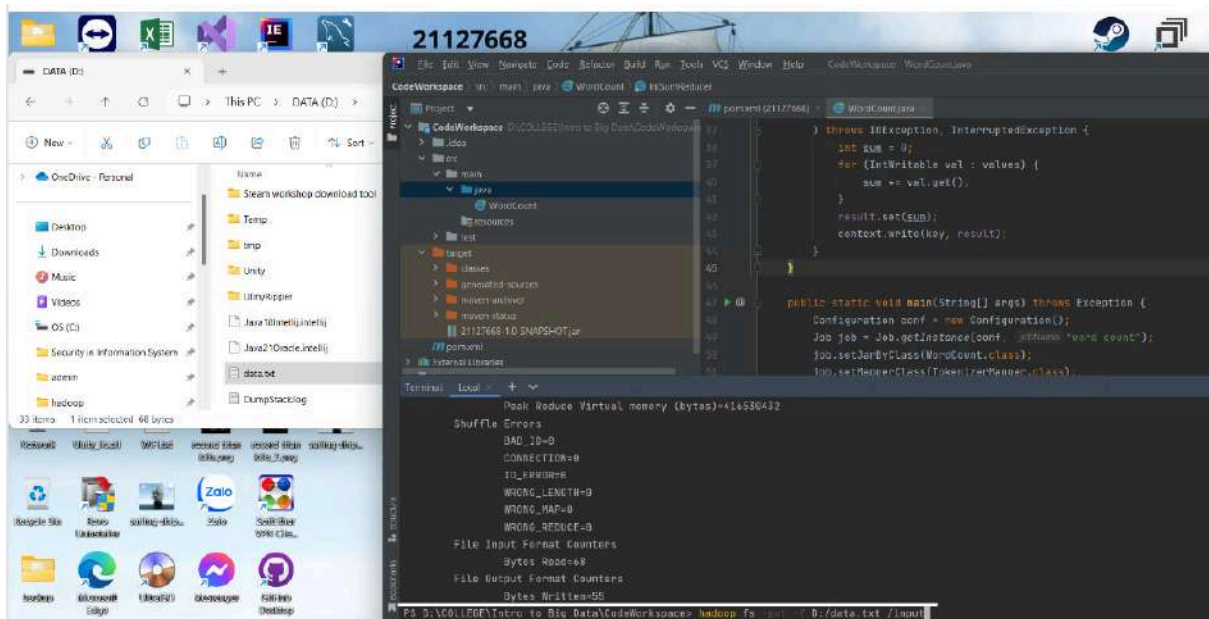● Create **input** folder in the Hadoop File System

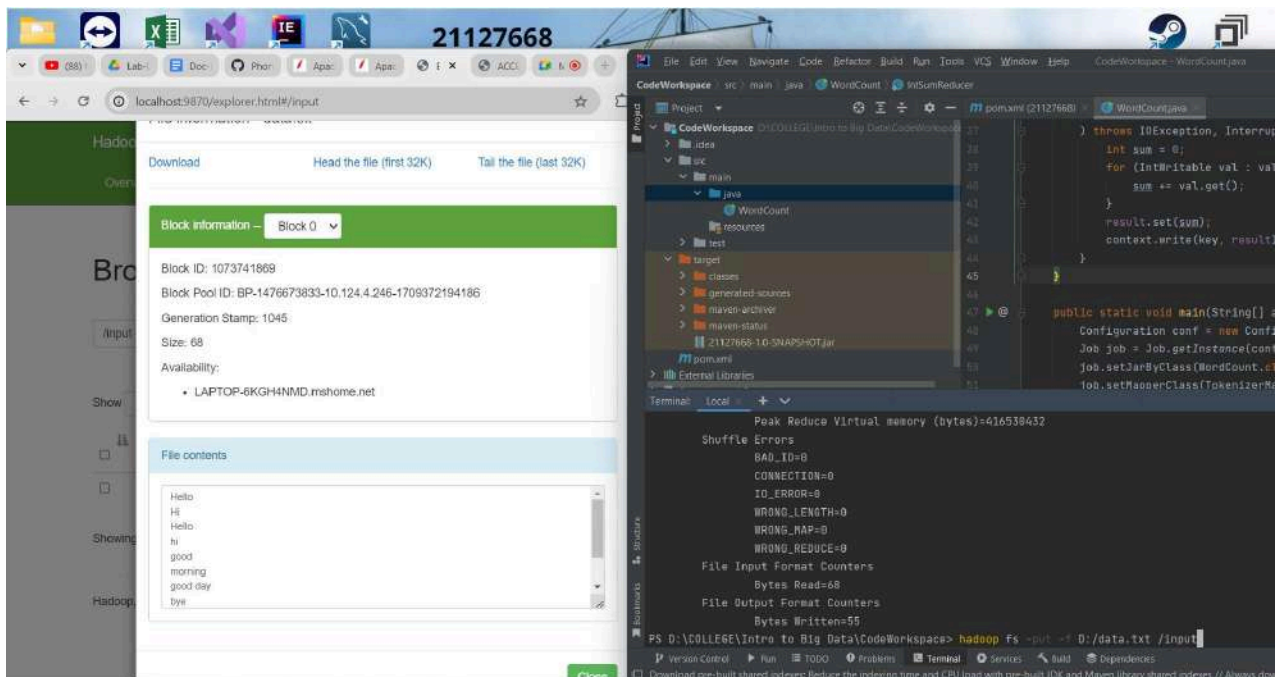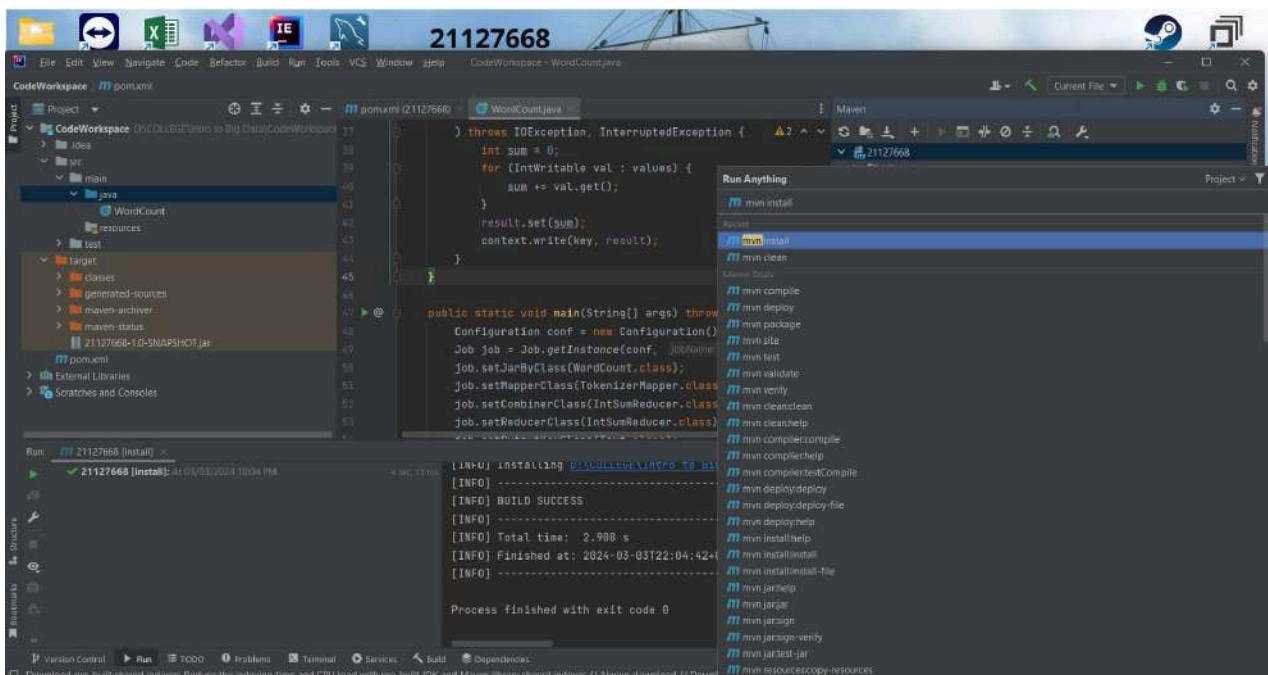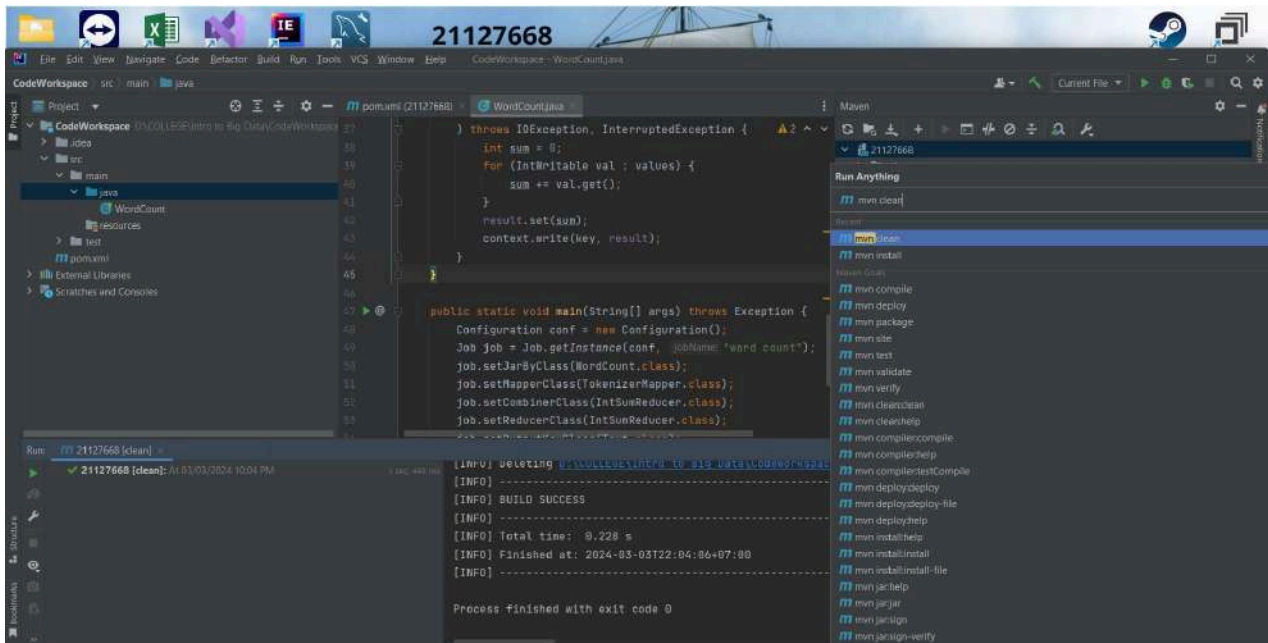On **localhost:9870**



- Put **data.txt** into input folder in the Hadoop File System
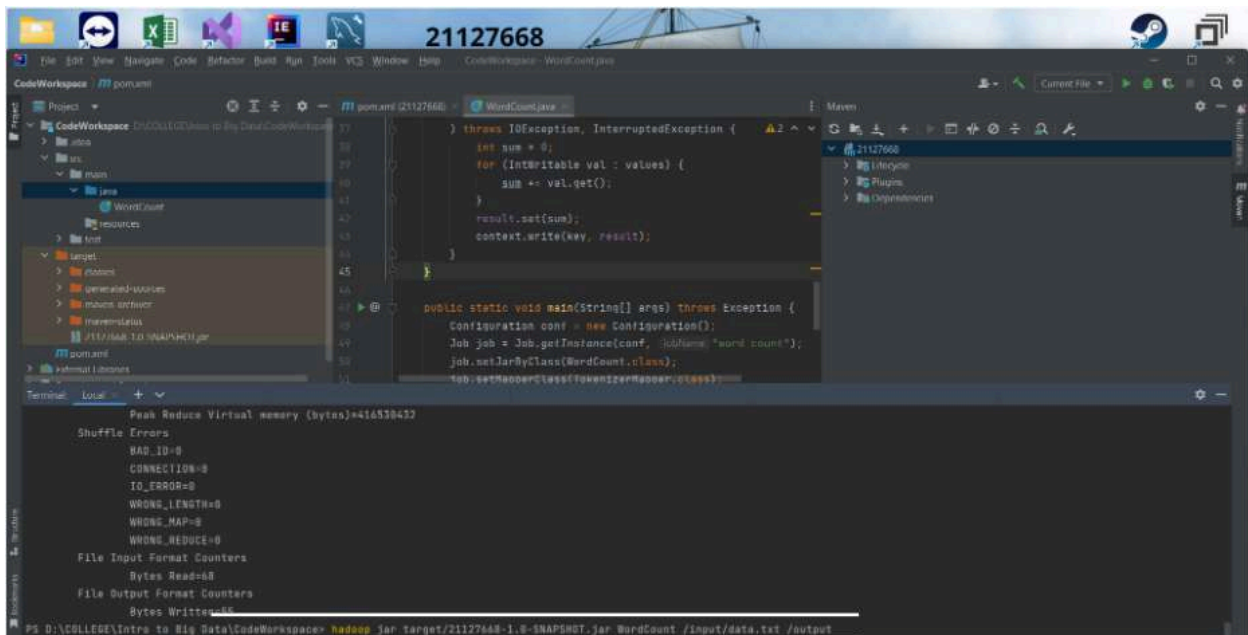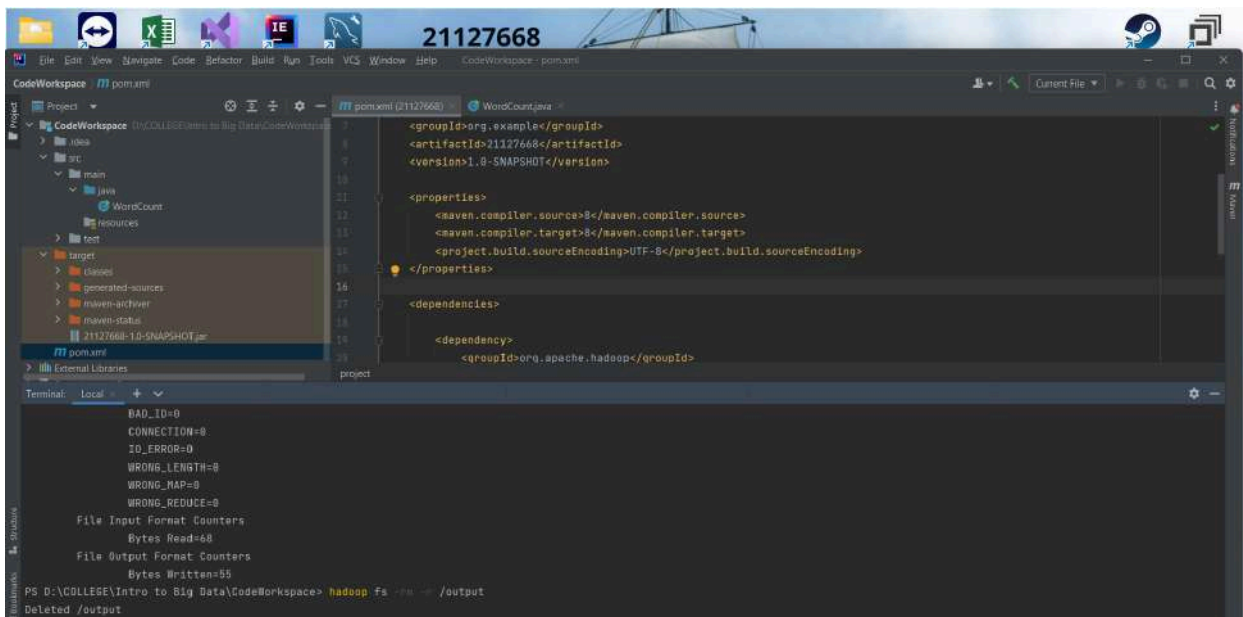
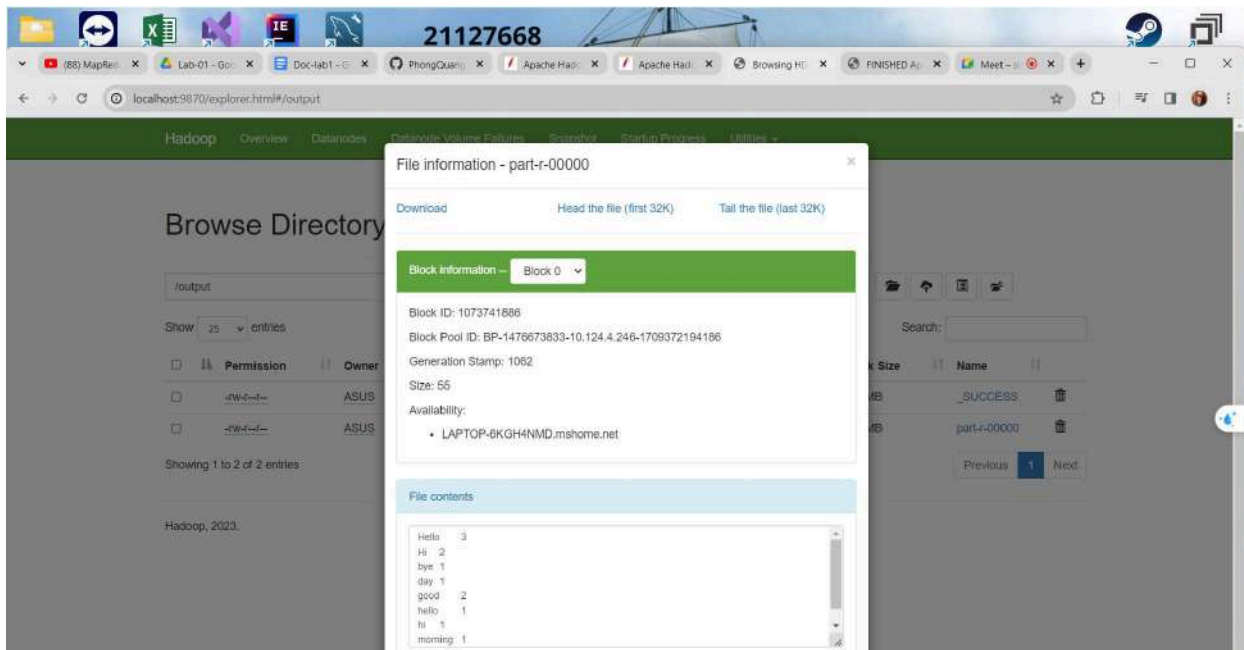On **localhost:9870**

● Create **jar file** of maven





● Run the **jar file** using hadoop jar command

- Eventually an output folder will be created and on **localhost:9870** at **part-r-00000** file which represents the final result
  If the output folder is already exist then we need to delete it

# 4. Reference

[1]https://www.youtube.com/watch?v=knAS0w-jiUk

[2]https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/ClusterSetup.html

[3]https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/SecureMode.html

[4]https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html

[5]https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html#Example%3A_WordCount_v1.0

[6]https://research.google/pubs/mapreduce-simplified-data-processing-on-large-clusters/

[7]https://storage.googleapis.com/gweb-research2023-media/pubtools/pdf/16cb30b4b92fd4989b8619a61752a2387c6dd474.pdf

[8]https://www.slideshare.net/oom65/ hadoop-security-architecture.

[9]Word Count

[10]Install Hadoop