



Department of Computer Science and Engineering (Data Science)

Subject: Machine Learning – IV Laboratory AY:  
2024-25

### Experiment 1

**Name: Jhanvi Parekh**

**Sap: 60009210033**

**Aim:** Implement a word count algorithm using MapReduce to count the number of times each word appears in a document.

**Theory:**

Hadoop can be developed in programming languages like Python and C++. MapReduce Hadoop is a software framework for ease in writing applications of software processing huge amounts of data. MapReduce Word Count is a framework which splits the chunk of data, sorts the map outputs and input to reduce tasks. A File-system stores the output and input of jobs. Re-execution of failed tasks, scheduling them and monitoring them is the task of the framework.

**Introduction:**

MapReduce is a programming model and processing technique designed for processing and generating large datasets that can be parallelized across a distributed cluster of computers. It is commonly used for big data processing tasks and is well-suited for tasks like word counting.

**MapReduce Workflow:**

The MapReduce workflow consists of two main phases: the Map phase and the Reduce phase.

#### 1. Map Phase:

- **Input Splitting:** The input data, such as a large text document, is divided into smaller equal-sized splits.
- **Mapping:** In this phase, the Mapper function processes each split independently. The Mapper takes the input data, extracts the required information, and emits key-value pairs. In the case of word counting, the key is typically a word, and the value is a count of 1.

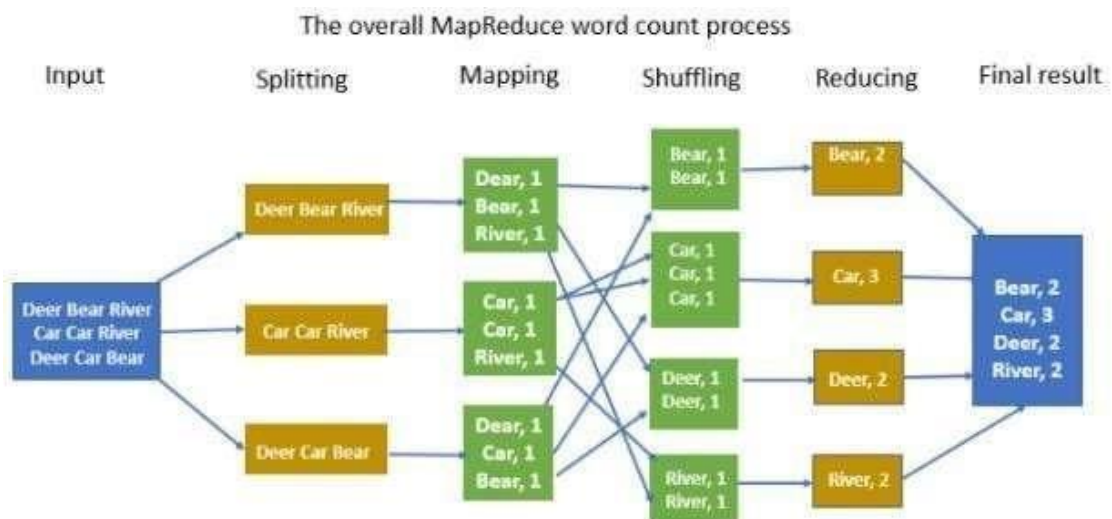
#### 2. Reduce Phase:



- **Shuffling and Sorting:** The MapReduce framework groups and sorts the output keyvalue pairs from the Mappers based on their keys. This ensures that all instances of the same word are grouped together.

Department of Computer Science and Engineering (Data Science)

- **Reducing:** The Reducer function processes the key-value pairs generated by the Map phase. In the word counting task, the Reducer takes the word as a key and a list of counts associated with that word. It then sums up the counts to produce the final count for that word.
- **Final Output:** The Reducers' output consists of key-value pairs where the key is the word and the value is the total count of that word in the input document.



#### Word Count using MapReduce:

The input given is converted into the string. Then it tokenizes them into words as if it need to break them. The mapper will append a single number or digit to each word and mapper outputs are shown above. Once we get the outputs as key-value pairs, once we pass the offset address as input to the mapper, the output of the value would be key-value pairs.

The output is getting into the sorting and shuffling phase. When we sort based on keys, all the keys will come to once a particular place. Sorting on the keys and shuffling the keys is done. A single word will go to a single reducer. Input to the reducer is key-value pairs. Once we pass outputs to reducer as input, the reducer will sum up all the values to keys.



That is, it groups up all the similar keys and output would be the concatenated key-value pair. The reducer will pick the result from the temp path and it will arrive at the final result. When we execute map-reduce, the input and output should be created in HDFS.

Sample Problem:

Input Data: Consider the following input text:

Line 1: Hello world, how are you? Hello, world!

Line 2: I'm doing well, thank you. How about you?

Department of Computer Science and Engineering (Data Science)

Map Phase: In this phase, we process each line and tokenize it into words. For each word found, we emit key-value pairs (word, 1):

"Hello" 1 "world" 1 "how" 1 "are" 1 "you" 1 "Hello" 1 "world" 1 "I'm" 1 "doing" 1 "well" 1 "thank you" 1 "How" 1 "about" 1 "you" 1

Shuffling and Sorting (Intermediate Step): The MapReduce framework groups and sorts the emitted key-value pairs by the key (word):

"Hello" 1, 1 "world" 1, 1 "how" 1 "are" 1 "you" 1, 1 "I'm" 1 "doing" 1 "well" 1 "thank you" 1 "How" 1 "about" 1

Reduce Phase: In this phase, we process each group of key-value pairs, summing up the values for each word to count its occurrences. We emit the word and its total count:

"Hello" 2 "world" 2 "how" 1 "are" 1 "you" 2 "I'm" 1 "doing" 1 "well" 1 "thank you" 1 "How" 1 "about" 1

This illustrates the Word Count process using key-value pairs, where the words serve as keys and the counts as values. The process efficiently counts the occurrences of each word in the input text by distributing the workload across multiple mappers and reducers.

Lab Exercise:

Step 1: Obtain a text document that you want to perform word count on. This can be a plain text file or a dataset in a suitable format.

Step 2: Implement a MapReduce program to count the number of times each word appears in the document. This program should consist of:

- A Mapper function that tokenizes the text, emits key-value pairs (word, 1) for each word.



- A Reducer function that receives key-value pairs, groups them by key (word), and sums up the values to get the word count.

Code with Output:

```
1347 sudo apt update
1348 sudo apt upgrade
1349 sudo apt openjdk-8-jdk
1350 sudo apt install openjdk-8-jdk
1351 wget https://downloads.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz
1352 tar -xvzf hadoop-3.3.6.tar.gz
1353 sudo mv hadoop-3.3.6 /usr/local/hadoop
1354 sudo mv hadoop-3.3.6 /usr/local/hadoop
1355 cd /usr/local/hadoop/etc/hadoop
1356 nano hadoop-env.sh
1357 nano hdfs-site.xml
1358 nano core-site.xml
1359 nano hdfs-site.xml
1360 cp mapred-site.xml.template mapred-site.xml
1361 nano mapred-site.xml
1362 nano yarn-site.xml
1363 nano hadoop-env.sh
1364 nano hdfs-site.xml
1365 nano mapred-site.xml
1366 sudo apt update
1367 sudo apt upgrade -y
1368 sudo apt install openjdk-8-jdk -y
1369 wget https://downloads.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz
1370 tar xzf hadoop-3.3.6.tar.gz
1371 tar xzf hadoop-3.3.6.tar.gz.1
1372 sudo mv hadoop-3.3.6 /usr/local/hadoop
1373 sudo chown -R hdoop:hdoop /usr/local/hadoop
1374 echo 'export JAVA_HOME=$(readlink -f /usr/bin/java | sed "s:bin/java::")' | sudo tee -a /usr/local/hadoop/etc/hadoop/hadoop-env.sh
1375 sudo apt update
1376 sudo apt upgrade -y
1377 sudo apt install openssh-server openssh-client -y
1378 sudo apt install openjdk-11-jdk -y
1379 sudo adduser hdoop
1380 su - hdoop
1381 sudo nano/etc/profile
1382 source ~/.bashrc
1383 sudo apt update
1384 sudo apt upgrade -y
1385 sudo apt install openssh-server openssh-client -y
1386 sudo apt install openjdk-11-jdk -y
1387 sudo adduser hdoop
1388 sudo adduser devansh
1389 su - devansh
1390 sudo apt update && sudo apt install openjdk-8-jdk
1391 java -version
1392 sudo apt install ssh
1393 sudo adduser dhruv209
1394 su - dhruv209
1395 sudo apt update
1396 sudo apt upgrade
1397 tar -xvzf hadoop-3.3.0.tar.gz
```



Administrator: Command Prompt

```
648 Jps
6224 NodeManager
62624 NameNode
9788 ResourceManager
18576 eclipse
18764 DataNode
12700 org.eclipse.equinox.launcher_1.6.700.v20231214-2017.jar

C:\windows\system32>hadoop fs -mkdir /input

C:\windows\system32>hadoop fs -put C:\Users\Admin\OneDrive\Desktop\input.txt \input

C:\windows\system32>hadoop fs -cat /output/part-00000
Additionally, 1
Physical 1
His 1
While 1
AI 5
s 1
Across 1
Also 1
and 3
Application 1
Areas 1
Seen 1
Considerations 1
Continued 1
Deployment 1
Development 2
Development 1
Different 1
Domains 1
Efficacy 1
Emerging 1
Existing 1
Exploration 1
Explore 1
For 1
Frameworks 1
Further 1
Have 1
Improve 1
In 1
Includes 1
Integration 1
Investigation 1
Language 1
Learning 1
Like 1
Made 1
Merit 1
Natural 1
Need 1
Of 2
Processing 1
Refining 1
Reinforcement 1
Research 1
```

Hadoop

Overview

Disks

Disks/Volume Tables

Graphical

Status Progress

Utilities

## Overview 'localhost:9000' (✓active)

Started:	Thu Mar 18 23:09:18 +0530 2023
Version:	3.3.0, 1a28f1871d1385883ac2932a51cc4703a843af
Compiled:	Tue Jul 27 00:14:05 +0530 2020 by beshm from branch-3.3.0
Cluster ID:	CID:6426a12-4604-4872-a780-a6877a76641
Block Pool ID:	BP-1251566485-102-100.1.2-1078987687502

## Summary

Security is off.

SafeMode is off.

1 Size and distribution: C blocks (Duplicated blocks, D denotes deduped block groups) = 1 total filesystem object(s).

Heap Memory used 96.81 MB of 256 MB Heap Memory. Max Heap Memory is 800 MB.

Non-Heap Memory used 53.64 MB of 57.5 MB. Committed Non-Heap Memory. Max Non-Heap Memory is unbounded.

Configured Capacity:	217.25 GB
Configured Remote Capacity:	0 B
DFS Used:	321 B (0%)



localhost:9870/explorer.html#/output

Hadoop Overview

Browse D

/output

Show 25 entries

Permission

Showing 1 to 2 of 2 entries

Hadoop, 2023.

Block information - Block 0

Block ID: 1073741832  
Block Pool ID: BP-542901707-192.168.0.175-1707831928007  
Generation Stamp: 1008  
Size: 601  
Availability:  
• Afreen-Laptop

File contents

Additionally, 1  
Ethical 1  
This 1  
While 1  
XAI 5  
a 1  
across 1  
also 1

Close

SUCCESS

part-00000

Previous 1 Next

Activities Firefox Web Browser Feb 1 15:24

hadoop

Cluster

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Used Resources	Total Resources
0	0	0	0	0	<memory:0 B, vCores:0>	<memory:8 GB, vCores:8>

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes
1	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation
Capacity Scheduler	[memory-mb (unit=Mb), vcores]	<memory:1024, vCores:1>	<memory:8192, vCores:4>

Show 20 entries

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU VCores	Allocated Memory MB
No data available in table														

Showing 0 to 0 of 0 entries





**About the Cluster**

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Used Resources	Total Resources	Reserved Resources	Physical Mem Used %	Physical VCores Used %
0	0	0	0	0	<memory:0 B, vCores:0>	<memory:8 GB, vCores:8>	<memory:0 B, vCores:0>	62	0

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Rebooted Nodes	Shutdown Nodes
1	0	0	0	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation	Maximum Cluster Application Priority	Scheduler Busy %
Capacity Scheduler	[memory-mb (unit=Mi), vcores]	<memory:1024, vCores:1>	<memory:8192, vCores:4>	0	0

Cluster overview

Cluster ID: 1706779227746

ResourceManager state: STARTED

ResourceManager HA state: active

ResourceManager zookeeper connection state: Could not find leader elector. Verify both HA and automatic failover are enabled.

ResourceManager RMStateStore: org.apache.hadoop.yarn.server.resourcemanager.recovery.NullRMStateStore

ResourceManager started on: Thu Feb 01 14:50:27 +0530 2024

ResourceManager version: 3.3.6 from 1be78238728da9266a4f88195058f0f8d012bf9c by ubuntu source checksum d42eb795a5eadb0feb75e44a7f87a9 on 2023-06-18T08:31Z

Hadoop version: 3.3.6 from 1be78238728da9266a4f88195058f0f8d012bf9c by ubuntu source checksum 5652179ad5f76cb287d9c633bb53bbd on 2023-06-18T08:22Z

## Conclusion:

Thus, we successfully implemented a Word Count using MapReduce. This exercise allowed us to harness the power of distributed computing to efficiently count word occurrences in a large dataset. We learned about the MapReduce paradigm and the importance of data preprocessing. This experiment provided valuable insights into handling big data and parallel processing.