School of Science and Engineering

**جـامـعة الأخـويـن**
**AL AKHAWAYN**
**U N I V E R S I T Y**

**Project N°1**

**Intro to Big Data Environment**

**Fall 2023**

**Submitted By :**

Asmaa Dalil (75675)

**Supervised by :**

Dr. Tajjeeddine Rachidi

## Table of Contents:

I.     Introduction
II.    Hadoop Architecture Used
    1.      Hadoop Architecture:

Using clusters of affordable technology, Hadoop is an open-source platform for the distributed archiving and analysis of massive datasets. Its design is made up of a number of essential parts that cooperate to offer fault tolerance, scalability, and effective data processing. Here is a list of the key elements of the Hadoop architecture.

- Hadoop HDFS to store data across slave machines.

- Hadoop YARN for resource management in the Hadoop cluster.

- Hadoop MapReduce to process data in a distributed fashion.

- Zookeeper to ensure synchronization across a cluster.

Under the scope of this project, we are going to use Hadoop Distributed File System (HDFS) and MapReduce model.

    2.      Docker Container
    3.      Hadoop Cluster Setup

Now , we are going to deploy Hadoop cluster by cloning from GitHub  link that includes the docker-hadoop folder in order to install the Hadoop cluster with 5 namenodes.

After we get into the folder, we use the command:

```
$ docker-compose up -d
```



Figure: creation of containers in docker

Then we use  the command *docker ps* in order to pull the running containers

School of Science and Engineering



Figure: the running containers

To get into the namenode container in order to test the 2 simple files as showing in the figure. This figure shows that we tested the wordcount program on the files and we got the results.

```yaml
C: > Windows > System32 > docker-hadoop > docker-compose.yml
60    datanode1:
61      build: ./datanode
62      container_name: datanode1
63      depends_on:
64        - namenode
65      volumes:
66        - hadoop_datanode1:/hadoop/dfs/data
67      env_file:
68        - ./hadoop.env
69
70    datanode2:
71      build: ./datanode
72      container_name: datanode2
73      depends_on:
74        - namenode
75      volumes:
76        - hadoop_datanode2:/hadoop/dfs/data
77      env_file:
78        - ./hadoop.env
79
80    datanode3:
81      build: ./datanode
82      container_name: datanode3
83      depends_on:
84        - namenode
85      volumes:
86        - hadoop_datanode3:/hadoop/dfs/data
87      env_file:
88        - ./hadoop.env
```

Figure: Compose.yml

we will create a directory with the same name in HDFS and copy all of its there. They are now spread over our 5  datanodes in the HDFS:



Figure: wordcount on simple files



Figure: results

## Inverted Index



Figure : copying the 20 files of 1000 words to run inverted index



Figure: running inverted index

School of Science and Engineering

Source code:



```java
import java.io.IOException;
import java.util.StringTokenizer;
import java.util.HashMap;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
//import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
//import org.apache.hadoop.mapreduce.MapContext;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class InvertedIndex {

    /*
    This is the Mapper class. It extends the hadoop's Mapper class.
    This maps input key/value pairs to a set of intermediate(output) key/value pairs.
    Here our input key is a Object and input value is a Text.
    And the output key is a Text and value is an Text. [word<Text> DocID<Text>]<->[aspect 57
    */
    public static class TokenizerMapper
            extends Mapper<Object, Text, Text, Text>{

        /*
        Hadoop supported datatypes. This is a hadoop specific datatype that is used to handle
        numbers and Strings in a hadoop environment. IntWritable and Text are used instead of
```

**Dificulties:**

**I had a problem running the containers on powershell so I had to use Docket ToolBox to run them**

School of Science and Engineering

| Name | Image | Status | CPU (%) | Port(s) | Last started | Actions |
|---|---|---|---|---|---|---|
| welcome-to-docker 0e5c115ce231 | docker/welcome-to-docker:latest | Running | 0% | 8088:80 ↗ | 7 hours ago | |
| docker-hadoop | | Other | 0% | | 32 seconds ago | |
| namenode 646580346dfc | docker-hadoop-namenode | Exited (126) | 0% | 9870:9870 ↗ | 5 hours ago | |
| datanode4 34b2fa26fa2a | docker-hadoop-datanode4 | Exited (126) | 0% | | 5 hours ago | |
| datanode1 60357fcdd434 | docker-hadoop-datanode1 | Exited (126) | 0% | | 5 hours ago | |
| datanode2 3ece792e5f5b | docker-hadoop-datanode2 | Exited (126) | 0% | | 5 hours ago | |
| datanode3 fc4c5ba93e60 | docker-hadoop-datanode3 | Exited (126) | 0% | | 5 hours ago | |

**Containers** Give feedback

Container CPU usage ⓘ
0.00% / 400% (4 cores allocated)

Container memory usage ⓘ
4.68MB / 5.98GB

Show charts ⌄

Search

Only show running containers

References:

https://github.com/big-data-europe/docker-hadoop#supported-hadoop-versions

https://chat.openai.com/

https://www.simplilearn.com/tutorials/hadoop-tutorial/hadoop-architecture#:~:text=Hadoop%20is%20a%20framework%20permitting,management%20in%20the%20Hadoop%20cluster