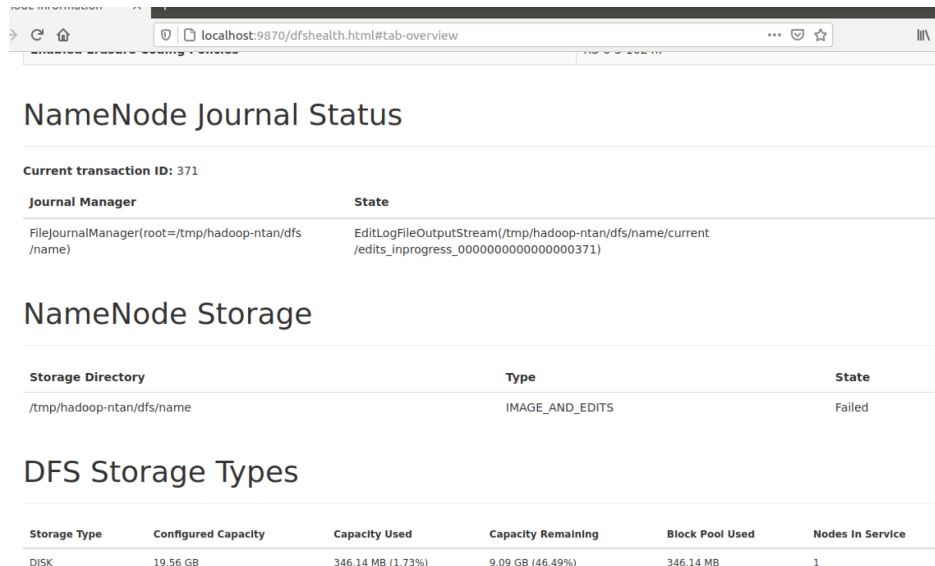# Course: Big Data

## *Lab 01*

# Set up Hadoop

***Fill answers of the questions below in the given tables.***
***Your screenshots must contain commands for required operations.***

## Question 1:

Set up Hadoop following this tutorial and then take a screenshot of the hdfs-site (localhost:9870) which shows information for NameNode Journal Status, NameNode Storage, DFS Storage Types (Overview page).
For example,



| Your screenshot goes here |
| --- |

Enabled Erasure Coding Policies                                     RS-6-3-1024k

## NameNode Journal Status

**Current transaction ID:** 1

| Journal Manager | State |
|---|---|
| FileJournalManager(root=/tmp/hadoop-vltanh/dfs/name) | EditLogFileOutputStream(/tmp/hadoop-vltanh/dfs/name/current/edits_inprogress_0000000000000000001) |

## NameNode Storage

| Storage Directory | Type | State |
|---|---|---|
| /tmp/hadoop-vltanh/dfs/name | IMAGE_AND_EDITS | Active |

## DFS Storage Types

| Storage Type | Configured Capacity | Capacity Used | Capacity Remaining | Block Pool Used | Nodes In Service |
|---|---|---|---|---|---|
| DISK | 279.79 GB | 24 KB (0%) | 107.82 GB (38.53%) | 24 KB | 1 |

Hadoop, 2021.

# Question 2:

Set up Hadoop following this tutorial and then take a screenshot of the yarn-site (localhost:8088) which shows information for nodes (http://localhost:8088/cluster/nodes ).
For example,

localhost:8088/cluster/nodes

Logged in as: dr.who

## Nodes of the cluster

**Cluster**
- About
- Nodes
- Node Labels
- Applications
  - NEW
  - NEW_SAVING
  - SUBMITTED
  - ACCEPTED
  - RUNNING
  - FINISHED
  - FAILED
  - KILLED
- Scheduler

▸ Tools

### Cluster Metrics

| Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Memory Used | Memory Total | Memory Reserved | VCores Used | VCores Total | VCores Reserved |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 B | 8 GB | 0 B | 0 | 8 | 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 |
|---|---|---|---|---|
| Capacity Scheduler | [memory-mb (unit=Mi), vcores] | <memory:1024, vCores:1> | <memory:8192, vCores:4> | 0 |

Show 20 entries                                           Search:

| Node Labels | Rack ⇕ | Node State ⇕ | Node Address ⇕ | Node HTTP Address ⇕ | Last health-update ⇕ | Health-report ⇕ | Containers ⇕ | Allocation Tags ⇕ | Mem Used ⇕ | Mem Avail ⇕ | VCores Used ⇕ | VCores Avail ⇕ | Version ⇕ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | /default-rack | RUNNING | ubuntu:34741 | ubuntu:8042 | Tue May 11 12:30:58 +0700 2021 |  | 0 |  | 0 B | 8 GB | 0 | 8 | 3.2.1 |

Showing 1 to 1 of 1 entries                    First  Previous  1  Next  Last

*Your screenshot goes here*

# Question 3:

Create a folder, named **testdata**, using HDFS commands and then show list of files/folders in the default folder (assume that you had already created the folder /user/<username> following this tutorial).
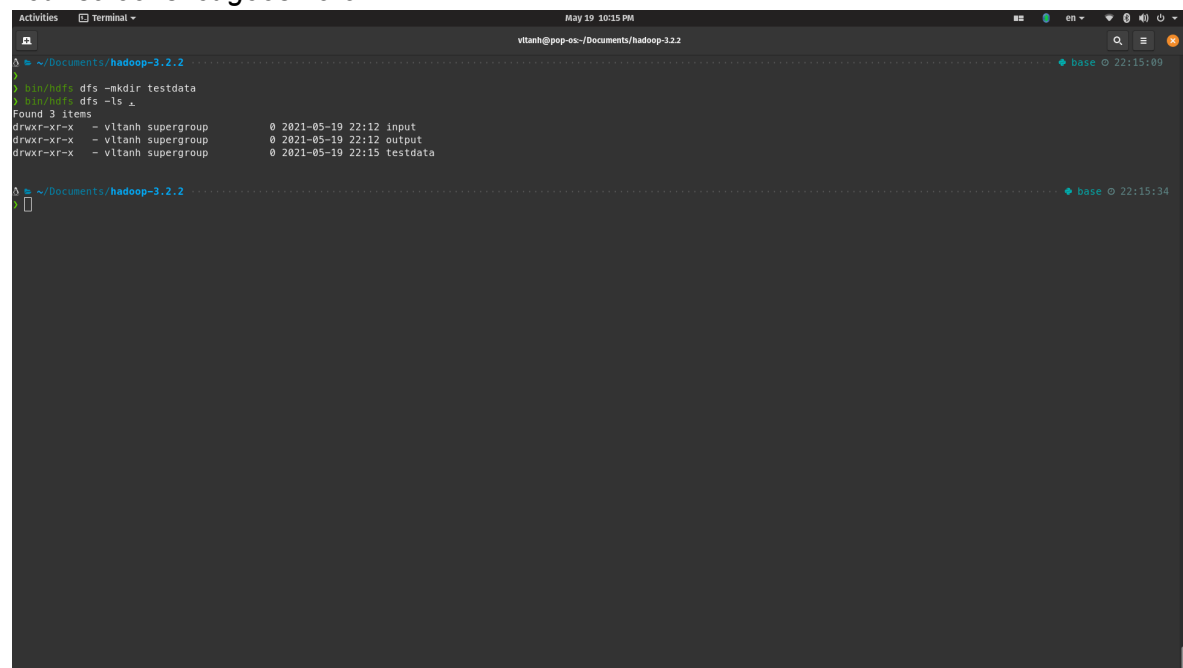
For example,



---

*Your screenshot goes here*

# Question 4:

Use HDFS commands to copy all xml files from **\<Hadoop folder\>/etc/hadoop** to folder **testdata/** and then show the size of each file in **testdata/**.
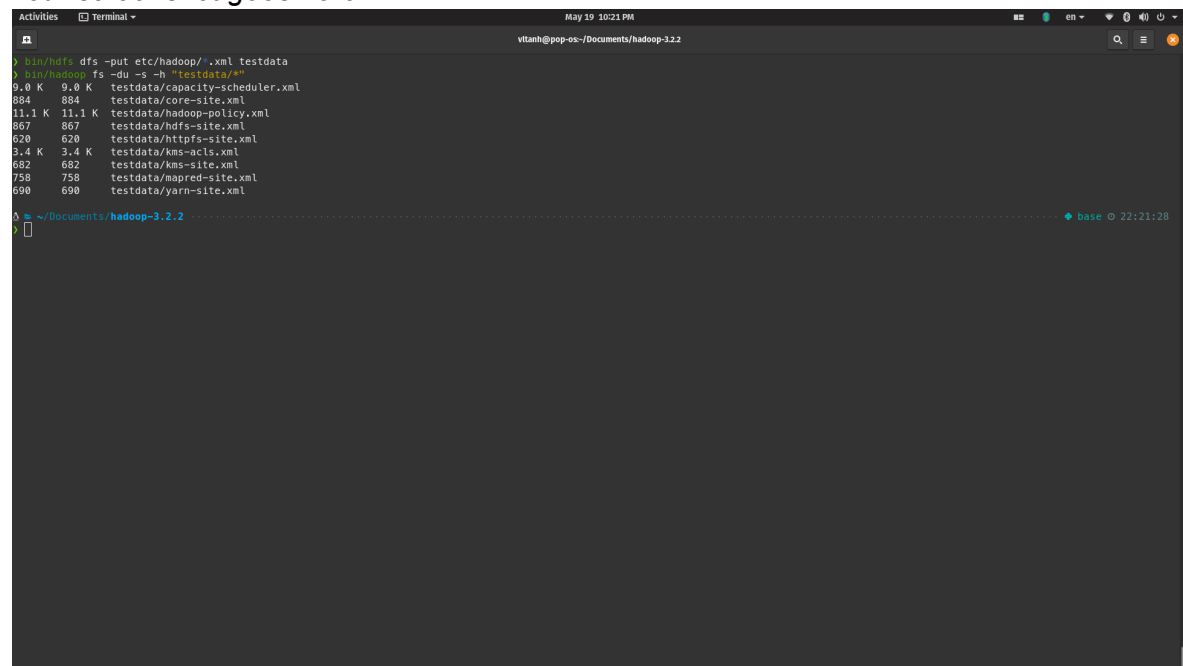
For example,



---

*Your screenshot goes here*

# Question 5:

Create a text file, named **data.txt**, in your local machine.
Write your fullname in **data.txt**.
Upload **data.txt** to **testdata/**
Create a text file, named **append.txt**, in your local machine.
Write your student number in **append.txt**.
Append the content of **append.txt** to **data.txt**
Display the content of **data.txt**
For example,





*Your screenshot goes here*

# Question 6:

Go to http://localhost:9870/explorer.html#/ to access the web interface for the file system. For example,

## Browse Directory

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | drwxrwxrwt | ntan | supergroup | 0 B | May 05 13:30 | 0 | 0 B | grep-temp-650720987 | 🗑 |
| ☐ | drwxrwxrwt | ntan | supergroup | 0 B | May 05 13:24 | 0 | 0 B | input | 🗑 |
| ☐ | drwxrwxrwt | ntan | supergroup | 0 B | May 10 16:07 | 0 | 0 B | output | 🗑 |
| ☐ | drwxr-xr-x | ntan | supergroup | 0 B | May 11 13:07 | 0 | 0 B | testdata | 🗑 |

Create a folder named **testoutput/** in **/user/<username>** and then your operation is denied. Specify the reason, fix the problem, and then conduct the requirement.

*Reason:* The error is "Permission denied: user=dr.who, access=WRITE, inode="/user/vltanh":vltanh:supergroup:drwxr-xr-x". This means that the "/user/vltanh" directory is owned by "vltanh" with 755 (drwxr-xr-x) permissions and therefore only "vltanh" can write to "/user/vltanh" ("dr.who" cannot).

*Solution*: Although this might not be a recommended operation, give everyone ("dr.who" included) the write permission (to "/user/vltanh") using "bin/hdfs dfs -chmod a+w /user/vltanh" then proceed.

*Your screenshot goes here*

# Submission Notice

- Export your answer file as pdf
- Rename the pdf following the format:

  <student number>_HoTen.pdf

  E.g: **123456_NguyenThanhAn.pdf**

  *If you have not been assigned a student number yet, then use 123456 instead.*
- Careless mistakes in filename, format, question order, etc. are not accepted (0 pts).
- Submission URL: https://forms.gle/zSitTjPTvjUogMBfA
- **Deadline: 17:00, 20/05/2021 (Thu).**