# ITECH2302 Big Data Management Laboratory - Hadoop

## Objectives:

- Installation of Hadoop/Spark environment
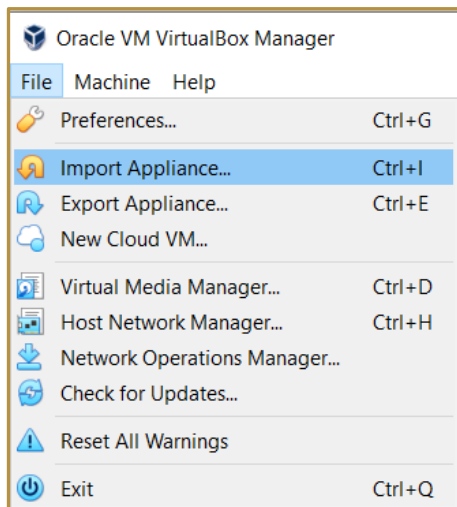- Introduction to Hadoop
- Review questions and activities

# Activity 1
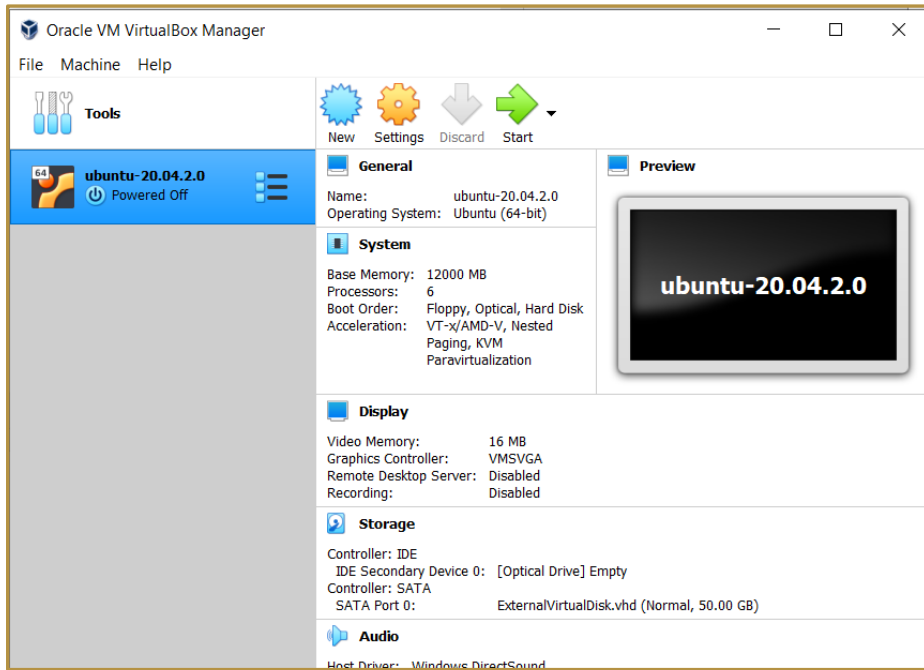
# Installation of Hadoop/Spark environment

1. Download and install the latest version of VirtualBox for your Operating System:
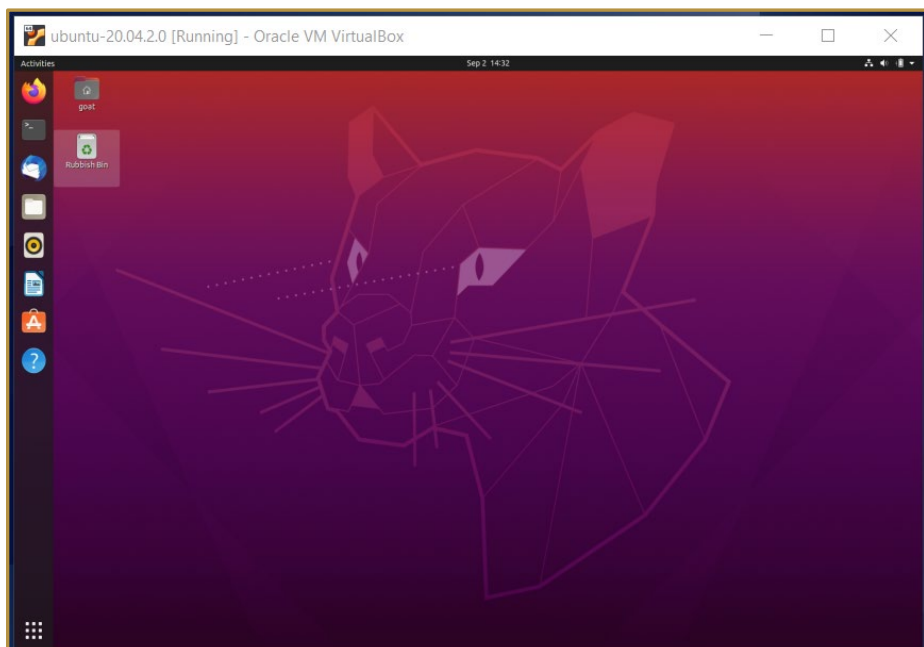


2. Download the OVA file for Hadoop/Spark from the FedUni server (refer to Announcements for details).

3. Load the OVA file into Virtualbox as an Appliance:

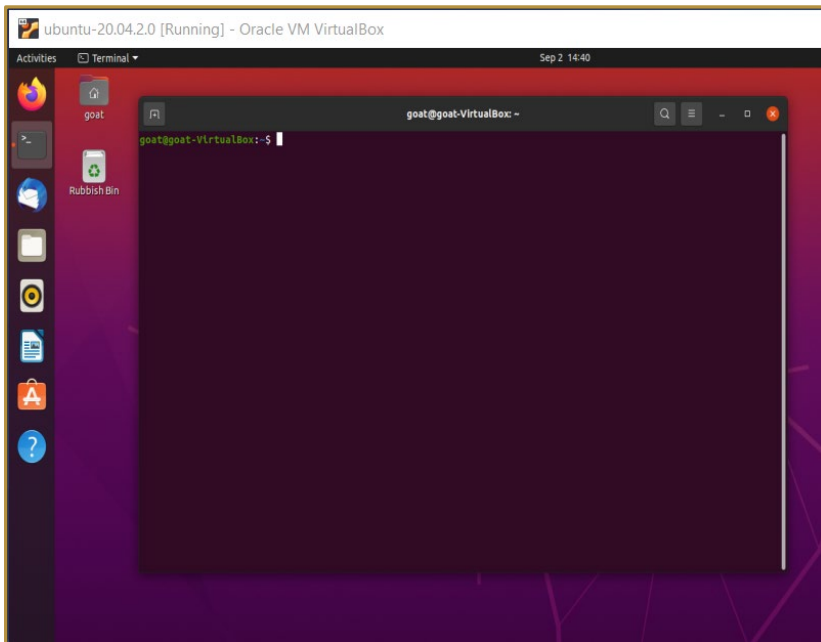4. Open the ubuntu operating system by selecting it and clicking the Start icon



5. Use the provided login details if required (*username*: goat, *password*: goat).
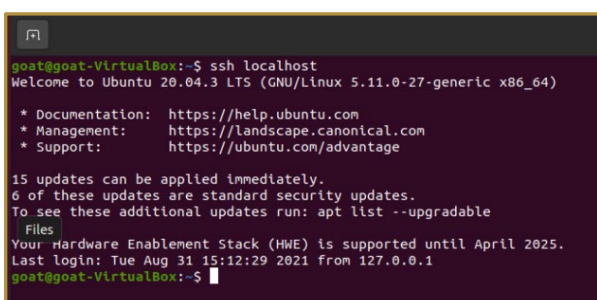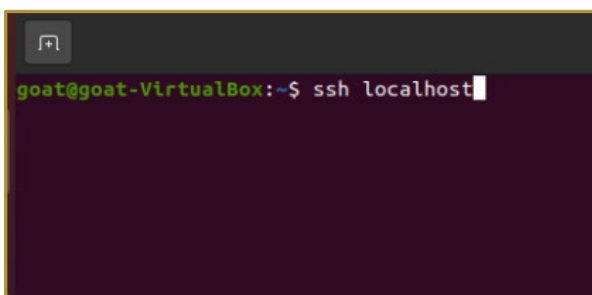
# Activity 2

## Introduction to Hadoop

1. Open a terminal with the ubuntu operating system





2. Familiarise yourself with writing command within the terminal:

3. The following commands should be written within the terminal:

ssh localhost
hdfs namenode -format

```
goat@goat-VirtualBox: ~                                    Q  ≡  _  □  ✕

les: false, skipCaptureAccessTimeOnlyChange: false, snapshotDiffAllowSnapRootDesce
ndant: true, maxSnapshotLimit: 65536
2021-08-26 19:50:07,789 INFO snapshot.SnapshotManager: SkipList is disabled
2021-08-26 19:50:07,819 INFO util.GSet: Computing capacity for map cachedBlocks
2021-08-26 19:50:07,819 INFO util.GSet: VM type       = 64-bit
2021-08-26 19:50:07,819 INFO util.GSet: 0.25% max memory 2.9 GB = 7.3 MB
2021-08-26 19:50:07,819 INFO util.GSet: capacity      = 2^20 = 1048576 entries
2021-08-26 19:50:07,931 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.wind
ow.num.buckets = 10
2021-08-26 19:50:07,931 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.
users = 10
2021-08-26 19:50:07,931 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.wind
ows.minutes = 1,5,25
2021-08-26 19:50:07,939 INFO namenode.FSNamesystem: Retry cache on namenode is ena
bled
2021-08-26 19:50:07,940 INFO namenode.FSNamesystem: Retry cache will use 0.03 of t
otal heap and retry cache entry expiry time is 600000 millis
2021-08-26 19:50:07,970 INFO util.GSet: Computing capacity for map NameNodeRetryCa
che
2021-08-26 19:50:07,970 INFO util.GSet: VM type       = 64-bit
2021-08-26 19:50:07,970 INFO util.GSet: 0.029999999329447746% max memory 2.9 GB =
898.3 KB
2021-08-26 19:50:07,970 INFO util.GSet: capacity      = 2^17 = 131072 entries
Re-format filesystem in Storage Directory root= /home/goat/hadoopdata/hdfs/namenod
e; location= null ? (Y or N)
```

Click Y

start-dfs.sh

```
goat@goat-VirtualBox:~$ start-dfs.sh
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [goat-VirtualBox]
```

start-yarn.sh

```
goat@goat-VirtualBox:~$ start-yarn.sh
Starting resourcemanager
Starting nodemanagers
```

jps

```
goat@goat-VirtualBox:~$ jps
10232 ResourceManager
4411 SecondaryNameNode
10749 Jps
3261 NameNode
10382 NodeManager
```
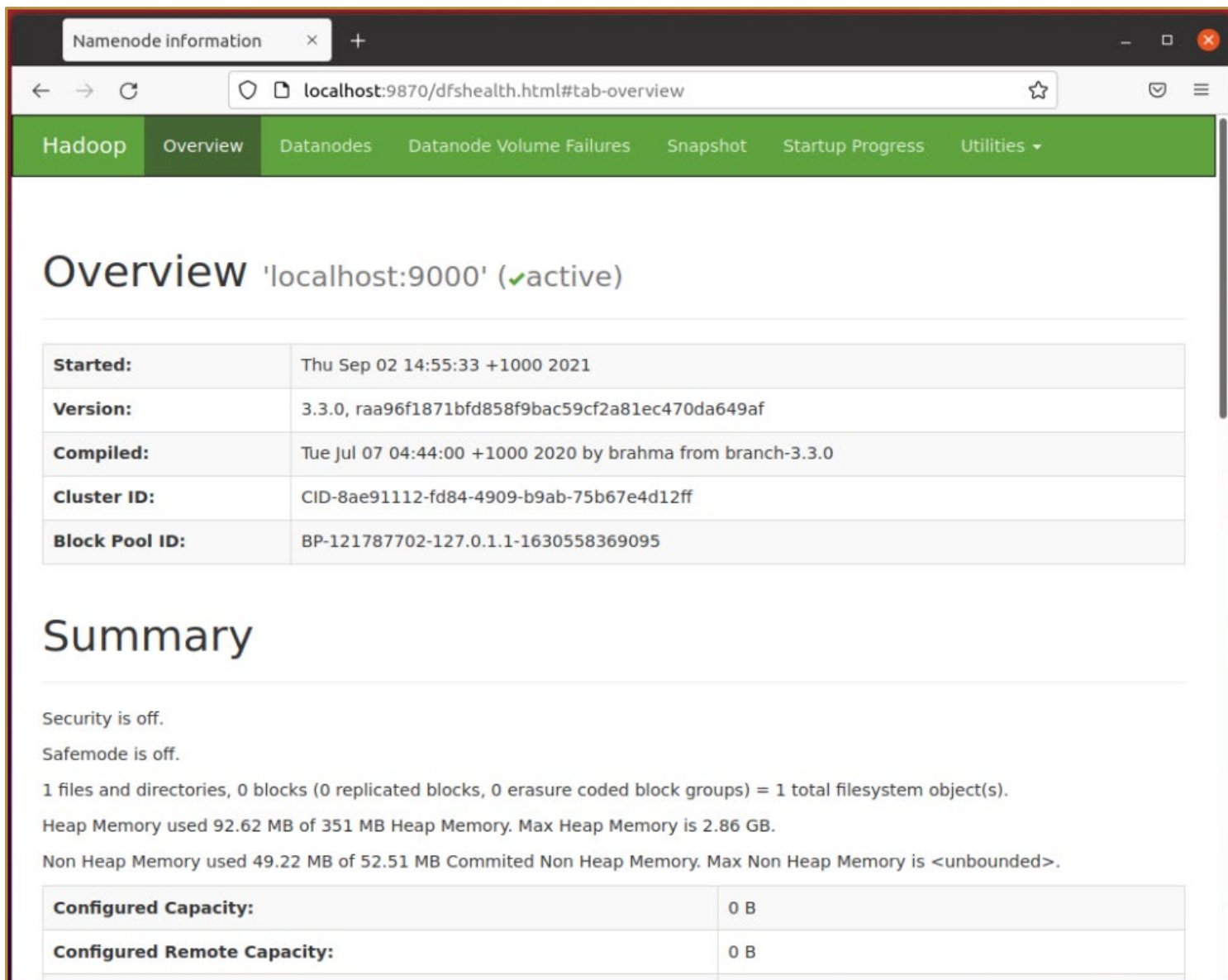
Open the following URL's in Firefox:

http://localhost:9870
http://localhost:8088

The pages should look like the following:

**Browser window 1 — All Applications (localhost:8088/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 |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |

**Cluster Nodes Metrics**

| Active Nodes | Decommissioning Nodes | Decommissioned Nodes |
|---|---|---|
| 1 | 0 | 0 |

**Scheduler Metrics**

| Scheduler Type | Scheduling Resource Type | Minimu |
|---|---|---|
| Capacity Scheduler | [memory-mb (unit=Mi), vcores] | <memory:1024, vCores:1> |

Show 20 entries

| ID | User | Name | Application Type | Application Tags | Queue | Application Priority | StartTime | LaunchTime | FinishTime |
|---|---|---|---|---|---|---|---|---|---|

Showing 0 to 0 of 0 entries

**Browser window 2 — All Applications (localhost:8088/cluster)**

Logged in as: dr.who

# All Applications

| ning | Memory Used | Memory Total | Memory Reserved | VCores Used | VCores Total | VCores Reserved |
|---|---|---|---|---|---|---|
| | 0 B | 8 GB | 0 B | 0 | 8 | 0 |

| Lost Nodes | Unhealthy Nodes | Rebooted Nodes | Shutdown Nodes |
|---|---|---|---|
| 0 | 0 | 0 | 0 |

| imum Allocation | Maximum Allocation | Maximum Cluster Application Priority |
|---|---|---|
| :1> | <memory:8192, vCores:4> | 0 |

Search:

| e | State | FinalStatus | Running Containers | Allocated CPU VCores | Allocated Memory MB | Reserved CPU VCores | Reserved Memory MB | % of Queue | % of Cluster | Progress | Tracking UI | Blacklisted Nodes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

No data available in table

First    Previous    Next    Last

4.  At the end of your session you can issue the following commands:

    stop-dfs.sh
    stop-yarn.sh

    Or:
    stop-all.sh

# Activity 3

There are many great resources for the whole ecosystem, covering a broad set of topics:

- https://www.edureka.co/blog/hadoop-ecosystem