# Running MapReduce Programs On Single Node Hadoop Cluster - Word Count/Word Frequency

Expt No: 3                                                                     March 06, 2019

Author: Subalakshmi Shanthosi S (186001008)

## Aim

Implementation of MapReduce in Hadoop single node cluster.

## Description

- Apache Hadoop
  - Large Scale,Open Source Software Framework.
  - Supports Three Projects:
    * Hadoop Common.
    * HDFS : Hadoop Distributed File System.
    * MapReduce.

- Hadoop MapReduce
  - Hadoop Programming Model and Software Framework.
  - Computational Processing:
    * Unstructured Data : File system
    * Structured Data : Database
  - MapReduce Layer has job and task tracker nodes.
  - Cluster nodes:
    * Single JobTracker per master.
    * Single TaskTracker per slave.
  - Fundamental Steps:
    * Map Step:
      · Master node slices problem input into several subproblems input.
      · Distributes data slices to worker nodes.
      · Worker nodes processes and hands over the control to master.
    * Reduce Step:
      · Master node takes the answers to the sub problems and combines them in a predefined way to get the output/answer to original problem.

## Software's Used

- Ubuntu 16.04 LTS

- Hadoop 1.0.3

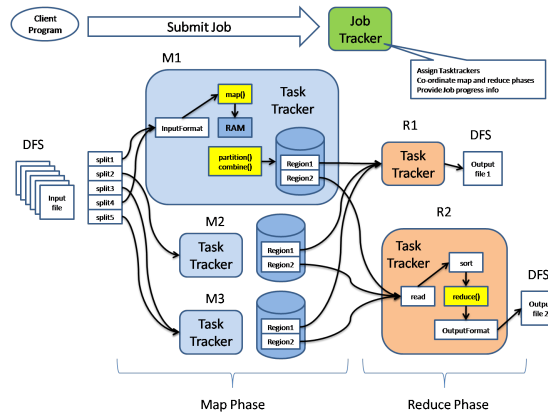# Hadoop MapReduce Architecture

- MapReduce Architecture:



Figure 1: Hadoop MapReduce Architecture Diagram.

# Procedure

1. Launch Ubuntu 16.04 LTS.

2. Login to the OS with sudo permission and install the following packages using apt-get command

   - openssh-server
   - openssh-client
   - java jdk 8
   - javac compiler
   - hadoop 2.7.3

# Output



Figure 2: Install openssh-server,openssh-client in Ubuntu OS.



Figure 3: Setting Java Home environment variable to the specified download path of JDK-1.7.



Figure 4: Adding a dedicated hadoop system user.

Figure 5: Configuring SSH in newly created user.



Figure 6: Disabling IPv6 in the newly created user account.



Figure 7: Disabling IPv6 in the newly created user account.

Figure 8: Installation of Hadoop 2.7.3 in new user login.



Figure 9: Configuring hadoop core-site.xml .
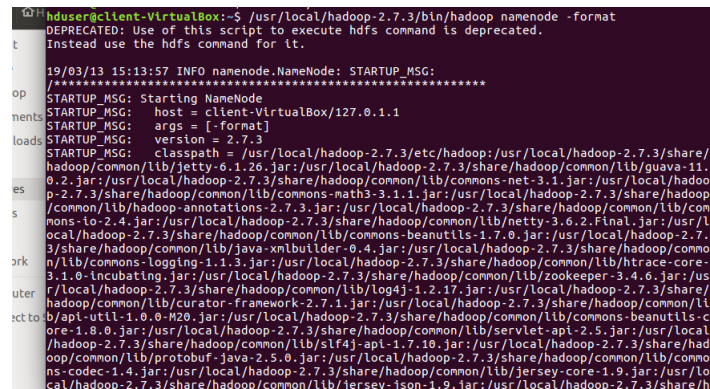


Figure 10: Configuring Hadoop MapReduce.

```
<configuration>

<!-- conf/hdfs-site.xml -->
<property>
  <name>dfs.replication</name>
  <value>1</value>
  <description>Default block replication.
  The actual number of replications can be specified when the file
is created.
  The default is used if replication is not specified in create time.
  </description>
</property>


</configuration>
```

Figure 11: Configuring Hadoop HDFS Site.



Figure 12: Formatting HDFS file system via the NameNode.

Figure 13: Starting hadoop NameNode,Datanode,JobTracker and TaskTracker.

# Result

Thus the hadoop single node cluster is sucessfully created in Ubuntu 16.04 OS version and required packages are installed.