# CS570: Week 2 Homework 1 : MapReduce on Ubuntu

## Md Shahadat Hossain Khan (19609)

---------------------------------------------------------------------------------------------------------------------

MapReduce on Ubuntu - IDE
      Step 1: Create a Ubuntu virtual Machine on GCP
      Step 2: Hadoop: Setting up a Single Node Cluster.
      Step 3: MapReduce Tutorial: Word Count


>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>

## Create a Ubuntu virtual Machine on GCP

> Create project: Click current project drop down menu from top menu
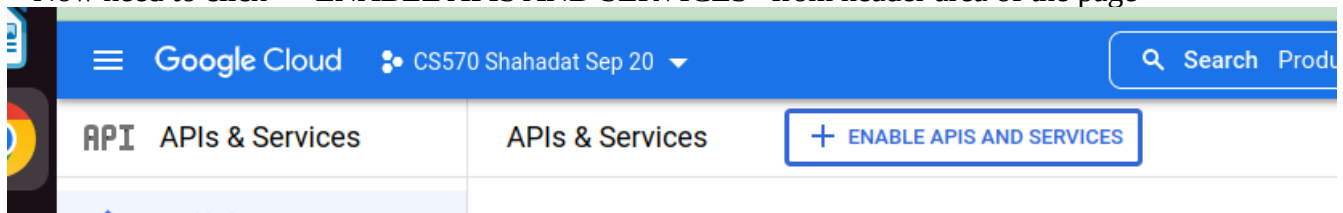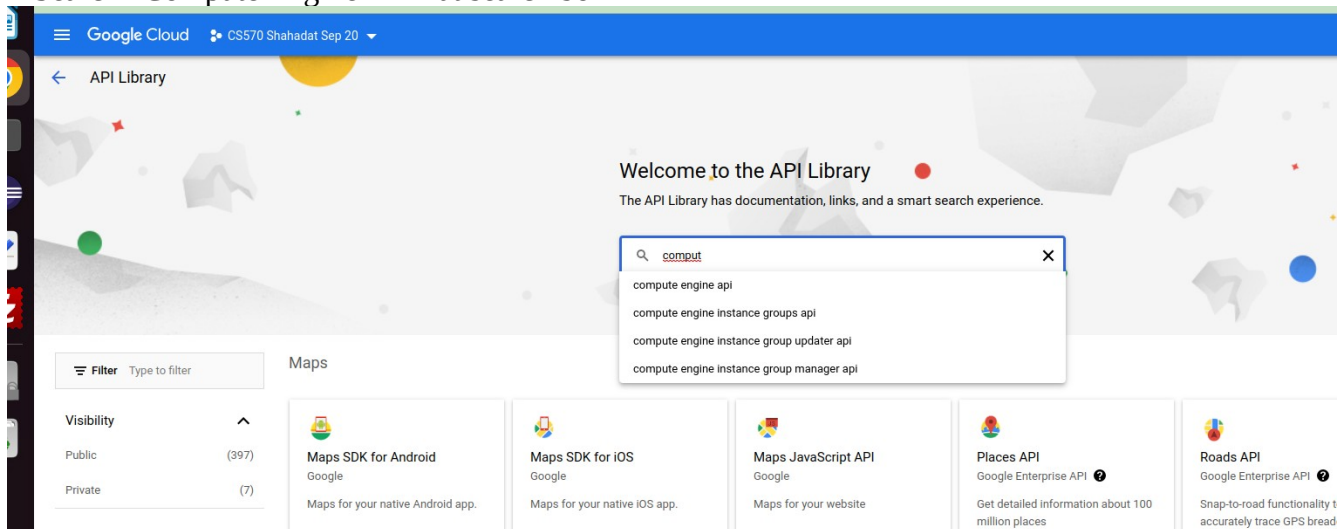


> Click "NEW PROJECT" from top right button of opened pop-up
> Provide a project name and click "Create"
> Now we need to select our new project, to do that we need to click project drop down menu from top menu to open project list pop-up
> Click on newly create project
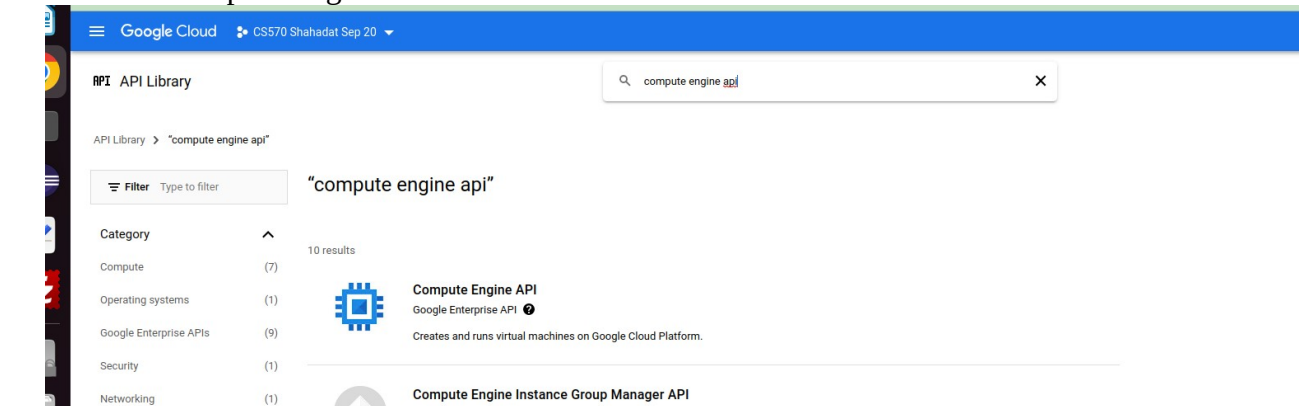> Now we need to enable "Compute Engine API", to do that click "API & Services" from left menu.

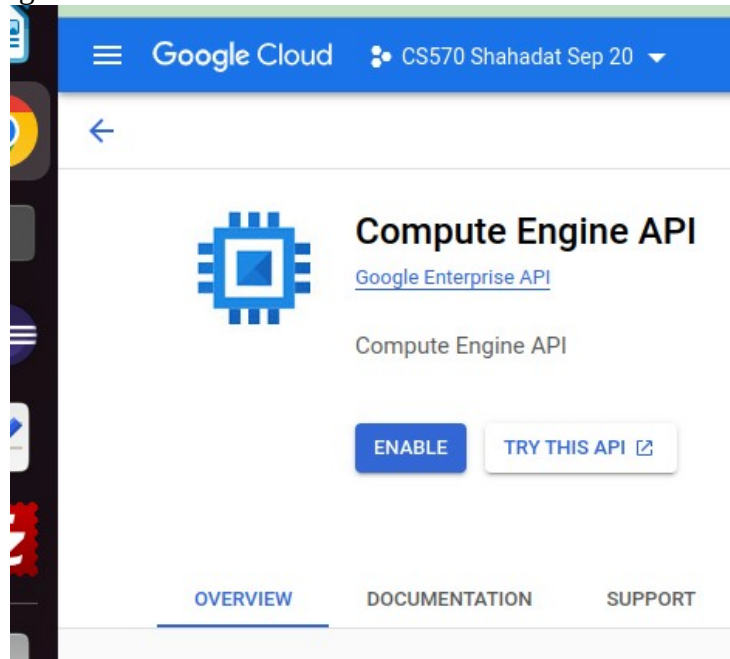> Now need to click "+ ENABLE APIS AND SERVICES" from header area of the page



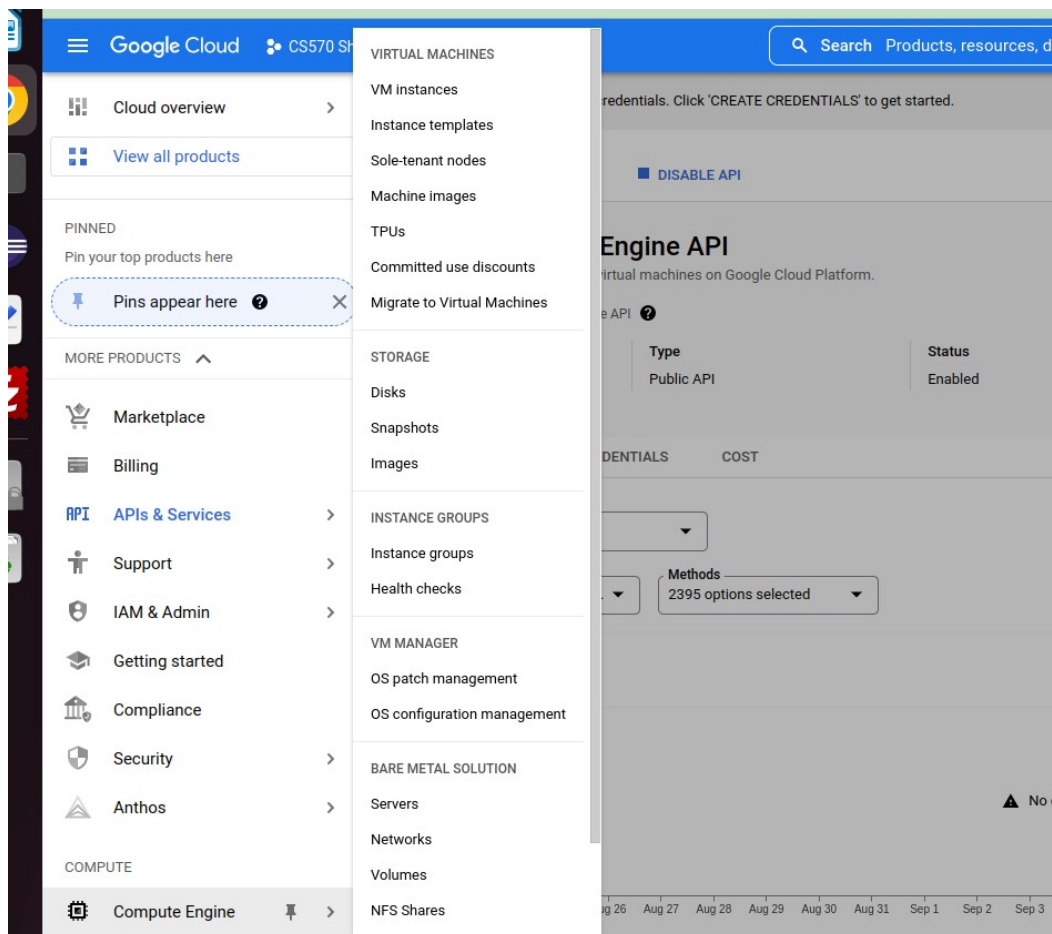> Search "Compute Engine API" at search box



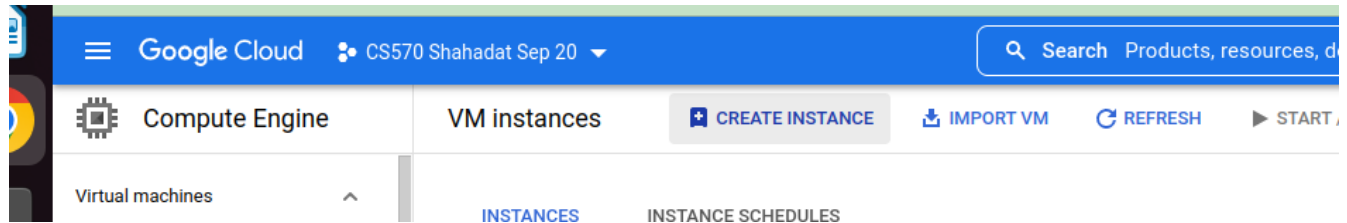> Click on "Compute Engine API" from search result

> Enable API by clicking "Enable" button



> After enabling API open VM instances page by clicking "Compute Engine" from left menu

> Click "Create Instance" from header menu



> Change configuration as image below



> Click "Create" button and wait until its ready to use

## We are going to add our SSH key into new VM

> At list of VM instances page, click on your new VM name to open that VM control panel page



> Now click "Edit" from that control panel
> From Edit page find "SSH" section click "+ ADD ITEM"

> Paste your local SSH public key into text box showed up and click "Save" button at the very bottom part of this page

Note: Please refer my assignment from week 1 and homework 2 to find proper format of SSH key before paste.

## Hadoop: Setting up a Single Node Cluster

> sudo apt-get update
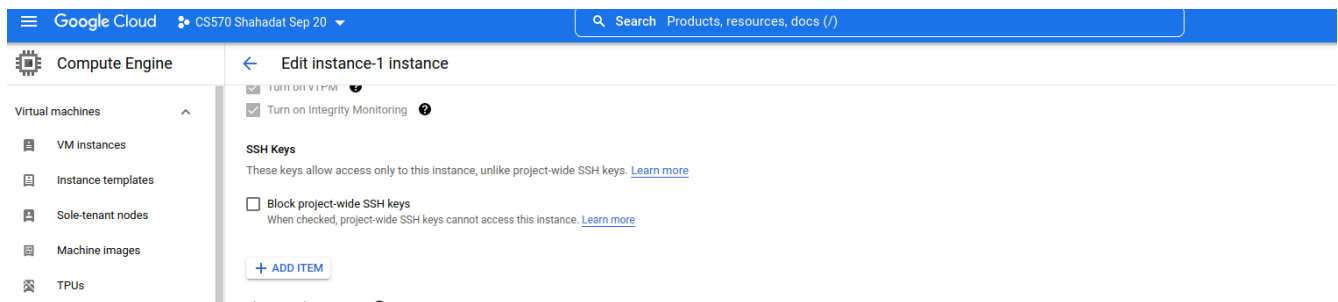> sudo apt-get install wget
> sudo apt-get install pdsh
> echo "ssh" | sudo tee /etc/pdsh/rcmd_default
Add `export PDSH_RCMD_TYPE=ssh into ~/.bashrc`
`You can check rcmd settings by using "pdsh -q -w localhost"`



> sudo apt-get install software-properties-common

At this stage we need to install java but hadoop need java 8.

We can install Java 8 automatically or manually. I did it manually, so I download java 8 from oracle archive site. Here are some link [https://www.oracle.com/java/technologies/javase/javase8u211-later-archive-downloads.html](https://www.oracle.com/java/technologies/javase/javase8u211-later-archive-downloads.html) and untared it and set some configuration as follows -
> cd /usr/lib/jvm/
> sudo tar -xvzf jdk-8u333-linux-x64.tar.gz
> sudo update-alternatives --install "/usr/bin/java" "java" "/usr/lib/jvm/jdk1.8.0_333/bin/java" 0
> sudo update-alternatives --install "/usr/bin/javac" "javac" "/usr/lib/jvm/jdk1.8.0_333/bin/javac" 0
> sudo update-alternatives --set java /usr/lib/jvm/jdk1.8.0_333/bin/java
> sudo update-alternatives --set javac /usr/lib/jvm/jdk1.8.0_333/bin/javac

That's all for my Java installation, now we need to set java path (if it is not set already) by modifying /etc/environment file, where we need append following line at the end.
*JAVA_HOME="/usr/lib/jvm/jdk1.8.0_333"*
Also we need to change PATH variable into this file. We need to prepend ${JAVA_HOME}

```
e.g.
PATH="/usr/lib/jvm/jdk1.8.0_333:/usr/local/sbin:/usr/local/bin:/usr/sbin:/
usr/bin:/sbin:/bin:/usr/games:/usr/local/games:/snap/bin"
JAVA_HOME="/usr/lib/jvm/jdk1.8.0_333"
```

### *Download Hadoop*

Now we can access our VM through local ssh.

Visit [https://dlcdn.apache.org/hadoop/common/stable/](https://dlcdn.apache.org/hadoop/common/stable/) or [https://www.apache.org/dyn/closer.cgi/hadoop/common/](https://www.apache.org/dyn/closer.cgi/hadoop/common/) or [https://dlcdn.apache.org/hadoop/common/](https://dlcdn.apache.org/hadoop/common/)

And using wget like -

wget [https://dlcdn.apache.org/hadoop/common/stable/hadoop-3.3.4.tar.gz](https://dlcdn.apache.org/hadoop/common/stable/hadoop-3.3.4.tar.gz)

After untar the hadoop, we can run hadoop locally. But we need to modify some settings to to run hadoop as single node cluster. Here is the commands that need to apply -

> cd <hadoop untared location>

> vi etc/hadoop/hadoop-env.sh

Add following line at the end of file or find somewhere you found "HADOOP_CLASSPATH" and add there.

export HADOOP_CLASSPATH="${JAVA_HOME}/lib/tools.jar"

Now we are ready to test hadoop!

> cd ~

> mkdir hadoop-namenode

> mkdir hadoop-datanode

> cd <hadoop untared location>

> vi etc/hadoop/core-site.xml

```
<configuration>
    <property>
        <name>fs.defaultFS</name>
        <value>hdfs://localhost:9000</value>
    </property>
</configuration>
```
> vi etc/hadoop/hdfs-site.xml

```
<configuration>
    <property>
        <name>dfs.replication</name>
        <value>1</value>
    </property>
<property>
        <name>dfs.permission</name>
        <value>false</value>
    </property>
<property>
        <name>dfs.namenode.name.dir</name>
        <value>file:///home/mkhan/hadoop-namenode</value>
    </property>
    <property>
        <name>dfs.datanode.data.dir</name>
        <value>file:///home/mkhan/hadoop-datanode</value>
    </property>
</configuration>
```

Now we need to able ssh to the localhost without a passphrase. To achieve this apply following commands -

> ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa

> cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys

> chmod 0600 ~/.ssh/authorized_keys

We all set, now we need to reboot our instance, like stop and start or we can restart.


Now we ready to run hadoop as single node cluster. Here is example execution -

> bin/hdfs namenode -format

> sbin/start-dfs.sh

> wget http://localhost:9870/

> bin/hdfs dfs -mkdir /user

> bin/hdfs dfs -mkdir /user/<username>

```
Note: Please use your login user name of your system/instance. In my case "mkhan",
but I don't research about this username properly
```

> bin/hdfs dfs -mkdir input

> bin/hdfs dfs -put etc/hadoop/*.xml input

> bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-3.3.4.jar grep input output 'dfs[a-z.]+'

```
Note: Many warning may come at screen, we can ignore by this time. I didn't
research on these warning.
```

> bin/hdfs dfs -cat output/*

> sbin/stop-dfs.sh


So, we ran hadoop successfully! Now I'm going to execute a real example like compiling code by Java. Following https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html

> bin/hdfs namenode -format

> sbin/start-dfs.sh

> bin/hdfs dfs -mkdir /user

> bin/hdfs dfs -mkdir /user/mkhan

> bin/hdfs dfs -mkdir /user/mkhan/wordcount

> bin/hdfs dfs -mkdir /user/mkhan/wordcount/input

> bin/hdfs dfs -put ../wordcount/input/* /user/mkhan/wordcount/input

> bin/hadoop fs -ls /user/mkhan/wordcount/input/

> bin/hadoop fs -cat /user/mkhan/wordcount/input/file01

I already create our source file into {user root}/wordcount/src directory. Now I'm going to compiling my code with hadoop and java

> cd ../wordcount/src

> ../../hadoop/bin/hadoop com.sun.tools.javac.Main WordCount.java

Now creating "jar" package to run our program into hadoop.

> jar cf wc.jar WordCount*.class

Here we executing our code by hadoop!

> ../../hadoop/bin/hadoop jar wc.jar WordCount /user/mkhan/wordcount/input /user/mkhan/wordcount/output

We all done! Now displaying output and stop our hadoop instance.

> cd ../../hadoop

> bin/hadoop fs -cat /user/mkhan/wordcount/output/part-r-00000

> bin/hadoop fs -ls /user/mkhan/wordcount/output/

> sbin/stop-dfs.sh

And here is the final result that created by hadoop process...



Thank You