

# Lab 2: MapReduce Programming

---

## 1. Objectives

- Understand the MapReduce concept.
- Get familiar with the Hadoop framework.
- Experience working with small Hadoop cluster using VMs.

## 2. Equipment Needs

- Computers
- Internet

## 3. Experiments

### 3.1 Basics

- Go through the Apache Hadoop introduction to get the general idea about Hadoop:  
<http://hadoop.apache.org/>
- Go through the Apache Hadoop release notes to understand the evolution of Hadoop:  
<http://hadoop.apache.org/releases.html>

### 3.2 Hadoop Single Node Mode.

- Follow the instructions on  
<http://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html>  
to set up Hadoop environment on your own Linux machine.
- Follow the instructions/tutorials from the above link to run simple practice with single node mode.

### 3.3 Hadoop “cluster”

Create two VMs (one master and one slave) on your own computer (ex: Virtualbox) and construct a small hadoop cluster for running the word count program. You can also use Docker containers to perform this assignment.

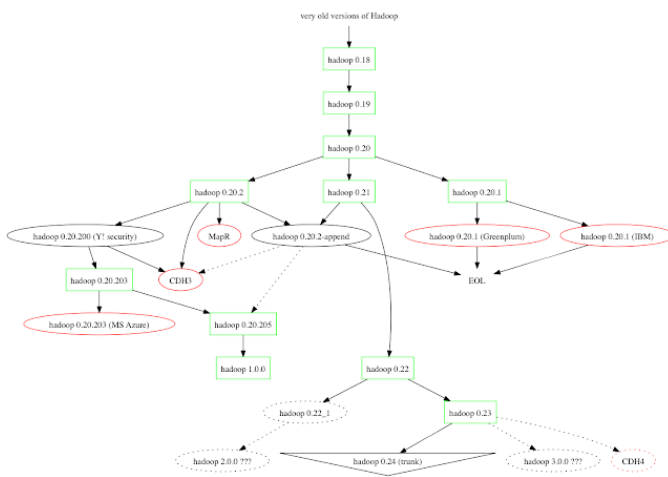
- Create two VMs.
- Configure the VM network so that the VMs can ping and communicate with each other.
- Download and install hadoop and all required tools.
- Configure hadoop configure files for your two-VM cluster
- Run the WordCount Hadoop job on the file (wordCountText.txt) provided by the TA.

## 4. Reports

### (a) What are the differences between Hadoop 0.X, 1.X, 2.X and 3.X?

First of all, the major different between 1.X and 2.X is YARN. The NameNode and JobTracker are replaced by ResourceManager and NodeManager in YARN. The design of YARN is to make the structure more scalable, since it has multi NameNode. (1.X can only support upto 4000 Nodes per cluster) Therefore, Hadoop 1.X is suitable for single purpose system, since it can only support the framework of MapReduce; Hadoop 2.X solve the scalability problem with multiple programming models.

According to the roadmap (<https://wiki.apache.org/hadoop/Roadmap>), we can learn that 2.X is actually not an improvement of 1.X. Instead, they are developed in parallel. 1.X is a branch of 0.20, after 0.21 is released; 2.X is a branch of 0.21. A graph from this website shows the diagram of the version branching. ([https://blogs.apache.org/bigtop/entry/all\\_you\\_wanted\\_to\\_know](https://blogs.apache.org/bigtop/entry/all_you_wanted_to_know))



The structure of Hadoop 3.X is similar with Hadoop 2.X. The main novel feature of 3.X is “Classpath isolation on by default”. Other plans of Hadoop 3.X is still not very clear.

Back to 0.X, from the diagram, we can see that it is the main branching of these versions. When a branch is slower than the next 0.2X, it becomes an independent branch. (eg. 1.X, 2.X)

### (b) What is YARN? Why do we need YARN?

YARN is short for “Yet Another Resource Negotiator”. It provides a more flexible framework for Hadoop to run more than just MapReduce applications by introducing higher-level cluster management.

Firstly, YARN can support more kinds of application. In Hadoop 1.X, users can only submit MapReduce based applications to Job Tracker. The program should follow the framework of MapReduce; however, YARN supports various applications.

Next, YARN provide higher-level cluster management that dynamically improves the utilization of clusters. In Hadoop 1.X, it relies on a single node NameNode to manage the state of an application, whose setting is static; however, in YARN, application manager can be in one of the nodes, managed by NodeManager. Moreover, with ResourceManager, who focus on scheduling and resource management, YARN has a better scalability.

### (c) What is Hadoop streaming?

Hadoop streaming is a utility that allows users to run Map/Reduce jobs with any external executable program/script. By assigning a mapper and reducer with “-mapper <mapper> -reducer <reducer>”, users

Screenshots of the practice of single node mode Hadoop on your own computer. The screenshot should show the output and result of Hadoop execution as well as the files in HDFS.

**“Follow the instructions/tutorials from the above link to run simple practice with single node mode.”**

## 1. Standalone Operation

```
$ mkdir input
$ cp etc/hadoop/*.xml input
$ bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.9.1.jar grep input output 'dfs[a-z.]+'
$ cat output/*
```

Result of Hadoop	Output
<pre> 2018-09-26 03:26:16,926 INFO mapreduce.Job: map: 100% reduce: 100% 2018-09-26 03:26:16,926 INFO mapreduce.Job: Job job_local1480712577_0001 completed successfully 2018-09-26 03:26:16,943 INFO mapreduce.Job: Counters: 36  File System Counters   FILE: Number of bytes read=3434753   FILE: Number of bytes written=0/20370   FILE: Number of read operations=0   FILE: Number of large read operations=0   FILE: Number of write operations=0  Map-Reduce Framework   Map input records=718   Map output records=1   Map output bytes=17   Map output materialized bytes=73   Input split bytes=909   Combine input records=1   Combine output records=1   Reduce input groups=1   Reduce shuffle bytes=73   Reduce input records=1   Reduce output records=1   Spilled Records=2   Shuffled Maps =9   Failed Shuffles=0   Merged Map outputs=0   GC time elapsed (ms)=276 </pre>	<pre> cia@vm1:~\$ ls output/ part-r-000000 _SUCCESS cia@vm1:~\$ cia@vm1:~\$ cia@vm1:~\$ cia@vm1:~\$ cat output/* 1 dfsadmin cia@vm1:~\$ █ </pre>
	Successfully find dfs*.

```
$ bin/hdfs dfs -mkdir input
$ bin/hdfs dfs -put etc/hadoop/*.xml input
$ bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.9.1.jar grep input
output 'dfs[a-z.]+'
$ bin/hdfs dfs -cat output/*
```

[illegible]

Word	Appear times
ashkan	106228
soheil	67768

	55057
frank	47661
in	27613

- (f) Use jps commands on both VMs to show running Hadoop daemons and provide screenshots.

VM1 (master)	VM2 (slave)
<pre> cia@vm1:~\$ jps 6249 Jps cia@vm1:~\$ start-dfs.sh Starting namenodes on [vm1] vm1: starting namenode, logging to /home/cia/hadoop/log s/hadoop-cia-namenode-vm1.out vm2: starting datanode, logging to /home/cia/hadoop/log s/hadoop-cia-datanode-vm2.out Starting secondary namenodes [0.0.0.0] 0.0.0.0: starting secondarynamenode, logging to /home/c ia/hadoop/logs/hadoop-cia-secondarynamenode-vm1.out cia@vm1:~\$ start-yarn.sh starting yarn daemons starting resourcemanager, logging to /home/cia/hadoop/l ogs/yarn-cia-resourcemanager-vm1.out vm2: starting nodemanager, logging to /home/cia/hadoop/l ogs/yarn-cia-nodemanager-vm2.out cia@vm1:~\$ jps 6992 Jps 6649 SecondaryNameNode 6809 ResourceManager 6413 NameNode cia@vm1:~\$ </pre>	<pre> cia@vm2:~\$ jps 3470 Jps cia@vm2:~\$ jps 3706 NodeManager 3828 Jps 3565 DataNode cia@vm2:~\$ </pre>

- (g) Screenshots of configuration files and IP addresses for Master node and Slave node of your small cluster as well as the MapReduce execution result. For each configuration file, please also briefly explain what it does.

VM1 (master)	VM2 (slave)	Description
<pre> root@vm1:~# ifconfig eth0: flags=4096&lt;UP,BROADCAST,MULTICAST&gt; mtu=1500     inet 192.168.50.134 netmask 255.255.255.0     ether 08:00:27:1a:4a:74 txqueuelen 1000 (KBytes)     RX packets 2208 (110.4 KB) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     RX bytes 109632 (106.3 KB)     TX packets 1222 (61.1 KB) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     TX bytes 15360 (15.0 KB) lo: flags=73&lt;UP,LOOPBACK,RUNNING&gt; mtu=65536     inet 127.0.0.1 netmask 255.0.0.0     ether &lt;??&gt; txqueuelen 1000 (KBytes)     RX packets 0 (0.0 B) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     TX packets 0 (0.0 B) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     TX bytes 0 (0.0 B) </pre>	<pre> root@vm2:~# ifconfig eth0: flags=4096&lt;UP,BROADCAST,MULTICAST&gt; mtu=1500     inet 192.168.50.135 netmask 255.255.255.0     ether 08:00:27:1a:4a:74 txqueuelen 1000 (KBytes)     RX packets 2208 (110.4 KB) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     RX bytes 109632 (106.3 KB)     TX packets 1222 (61.1 KB) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     TX bytes 15360 (15.0 KB) lo: flags=73&lt;UP,LOOPBACK,RUNNING&gt; mtu=65536     inet 127.0.0.1 netmask 255.0.0.0     ether &lt;??&gt; txqueuelen 1000 (KBytes)     RX packets 0 (0.0 B) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     TX packets 0 (0.0 B) 0 errors 0 overruns 0 carrier 0     collisons 0 frame 0     TX bytes 0 (0.0 B) </pre>	<b>IP addresses:</b> Master: 192.168.50.134 Slave: 192.168.50.135
<pre> root@vm1:~# cat /etc/hadoop/conf/core-site.xml &lt;configuration&gt;   &lt;property&gt;     &lt;name&gt;fs.defaultFS&lt;/name&gt;     &lt;value&gt;hdfs://vm1:9000/&lt;/value&gt;   &lt;/property&gt; &lt;/configuration&gt; </pre>	<pre> root@vm2:~# cat /etc/hadoop/conf/core-site.xml &lt;configuration&gt;   &lt;property&gt;     &lt;name&gt;fs.defaultFS&lt;/name&gt;     &lt;value&gt;hdfs://vm1:9000/&lt;/value&gt;   &lt;/property&gt; &lt;/configuration&gt; </pre>	<b>core-site.xml:</b> Mainly configure the URI of NameNode, which is the master.
<pre> root@vm1:~# cat /etc/hadoop/conf/hdfs-site.xml &lt;configuration&gt;   &lt;property&gt;     &lt;name&gt;dfs.namenode.name.dir&lt;/name&gt;     &lt;value&gt;/home/cia/hadoop/data/nn&lt;/value&gt;   &lt;/property&gt;   &lt;property&gt;     &lt;name&gt;dfs.datanode.data.dir&lt;/name&gt;     &lt;value&gt;/home/cia/hadoop/data/dn&lt;/value&gt;   &lt;/property&gt;   &lt;property&gt;     &lt;name&gt;dfs.replication&lt;/name&gt;     &lt;value&gt;1&lt;/value&gt;   &lt;/property&gt; &lt;/configuration&gt; </pre>	<pre> root@vm2:~# cat /etc/hadoop/conf/hdfs-site.xml &lt;configuration&gt;   &lt;property&gt;     &lt;name&gt;dfs.namenode.name.dir&lt;/name&gt;     &lt;value&gt;/home/cia/hadoop/data/nn&lt;/value&gt;   &lt;/property&gt;   &lt;property&gt;     &lt;name&gt;dfs.datanode.data.dir&lt;/name&gt;     &lt;value&gt;/home/cia/hadoop/data/dn&lt;/value&gt;   &lt;/property&gt;   &lt;property&gt;     &lt;name&gt;dfs.replication&lt;/name&gt;     &lt;value&gt;1&lt;/value&gt;   &lt;/property&gt; &lt;/configuration&gt; </pre>	<b>hdfs-site.xml:</b> Configure the nameNode and dataNode. Here I assign the directory for nameNode and dataNode for further observation.
<pre> root@vm1:~# cat /etc/hadoop/conf/slaves vm2 </pre>	<pre> root@vm2:~# cat /etc/hadoop/conf/slaves vm2 </pre>	<b>Slaves:</b> Configure the hosts of the slaves.
<pre> root@vm1:~# cat /etc/hadoop/conf/mapred-site.xml &lt;configuration&gt;   &lt;property&gt;     &lt;name&gt;mapreduce.framework.name&lt;/name&gt;     &lt;value&gt;yarn&lt;/value&gt;   &lt;/property&gt; &lt;/configuration&gt; </pre>	<pre> root@vm2:~# cat /etc/hadoop/conf/mapred-site.xml &lt;configuration&gt;   &lt;property&gt;     &lt;name&gt;mapreduce.framework.name&lt;/name&gt;     &lt;value&gt;yarn&lt;/value&gt;   &lt;/property&gt; &lt;/configuration&gt; </pre>	<b>mapred-site.xml:</b> Configure the framework for MapReduce.

The image displays two side-by-side terminal windows showing the configuration of `yarn-site.xml`. The left window shows the initial configuration with the following properties:

```
<configuration>
<!-- Site specific YARN configuration properties -->

<property>
<name>yarn.acl.enabled</name>
<value>false</value>
</property>

<property>
<name>yarn.resourcemanager.hostname</name>
<value>vm11</value>
</property>

<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value>
</property>
</configuration>
-- INSERT --
```

The right window shows the same configuration with an additional property added:

```
<configuration>
<!-- Site specific YARN configuration properties -->

<property>
<name>yarn.acl.enabled</name>
<value>false</value>
</property>

<property>
<name>yarn.resourcemanager.hostname</name>
<value>vm11</value>
</property>

<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value>
</property>
</configuration>
-- INSERT --
```

A large text overlay in the center of the image reads: **yarn-site.xml: Configure the URI/host of yarn.**

## Result

```
$ hadoop jar WordCount.jar WordCountText.txt output
```

```
$ hdfs dfs -cat output/*
```

Result of hadoop	Cat the results
<pre> 10/09/27 02:12:38 INFO mapreduce.jobsubmitter: Submitting tokens for job: job_15380695647_0001 10/09/27 02:12:32 INFO org.apache.hadoop.mapreduce.lib.input: Submitted application application_15380695647_0001 10/09/27 02:12:32 INFO org.mapreduce.job: The url to track the job: http://maprproxy/application_15380695647_0001/ 10/09/27 02:12:32 INFO mapreduce.job: Running job: job_15380695647_0001 10/09/27 02:12:52 INFO mapreduce.job: Job ID: job_15380695647_0001 running in uber mode : false 10/09/27 02:12:52 INFO mapreduce.job: map 0% reduce 0% 10/09/27 02:13:09 INFO mapreduce.job: map 100% reduce 0% 10/09/27 02:13:22 INFO mapreduce.job: map 100% reduce 100% 10/09/27 02:13:24 INFO mapreduce.job: Job completed successfully 10/09/27 02:13:26 INFO mapreduce.job: Counters: 49 File system counters FILE: Number of bytes read=2622214 FILE: Number of bytes written=501250 FILE: Number of read operations=6 FILE: Number of large read operations=0 FILE: Number of write operations=0 HDFS: Number of bytes read=1027541 HDFS: Number of bytes written=101959 HDFS: Number of read operations=6 HDFS: Number of large read operations=0 HDFS: Number of write operations=2 Job Counters Launched map tasks=1 </pre>	<pre> cli@hwi-MCS hdfs dfs -ls Found 2 items -rw-r--r- 1 cin supergroup 106273M 2018-09-27 02:09 wordContext.txt drwxr-xr-x 1 cin supergroup 0 2018-09-27 02:13 output cli@hwi-MCS hdfs dfs -ls output Found 2 items -rw-r--r- 1 cin supergroup 0 2018-09-27 02:13 output/SUCCESS -rw-r--r- 1 cin supergroup 19769M 2018-09-27 02:13 output/part-r-000000 cli@hwi-MCS hdfs dfs -cat output/* Hacking +-----+ +-----+ +-----+-----A 1 +-----+-----ashkin 1 +-----+-----spawning 1 +-----+-----ddt imp. 1 +-----+-----frank 1 +-----+-----ashkin 1 +-----+-----Piper 1 +-----+-----crowds 1 +-----+-----With 1 +-----+----- 2 +-----+----- 1 +-----+----- 1 </pre>

- (h) What are the differences between Hadoop master and slave nodes? Also name what functionalities are performed on each node.

In our experiment setting, our master node runs: **NameNode, SecondaryNameNode, ResourceManager**; our slave node runs: **DataNode, NodeManager**. In brief, master node mainly allocates the resource and assigns jobs, while slave node mainly execute the jobs.

From HDFS's point of view, the master node manages HDFS storage with NameNode, and the slave node have to handle the HDFS service to let NameNode store files. Moreover, the slave node will enable TaskTracker for managing the map/reduce tasks if a task is assigned to the node in Hadoop 1.X.

From YARN's point of view, the master node handles the scheduling of tasks and management of the resources in clusters with Resource Manager; on the other hand, the slave node manages its own resource and reports back to Resource Manager. Moreover, the slave nodes hold containers that can be Application Master or run tasks. When an application is assigned to YARN, the Application Master on the slave node will track the progress on the cluster, and the tasks will be executed on slave nodes.

- (i) Write a pseudo code to multiply large matrices using Hadoop. Also explain the function of your Mappers and Reducers.

Firstly, the problem does not provide the input of the large matrices, so we chose the one that is commonly used in the similar problem in Hadoop. The reason we need to specific the input format is that the mapping in Hadoop is divided with the new line in the file. Different kinds of input format will affect how the data is divided, and affect how mapping function should be designed according to the format. For example, in the following table, Case1 and Case2 has different data formats. In Case1, the mapper function receives "M,0,0,10", while in Case2, the mapper function receives "1 2 3 4". The

former one only receives one element of a matrix, and the latter one receives four element of the matrix. Therefore, in this answer, we take Case1 as our input format.

Case1. Input (Matrix M)	Case2. Input (Matrix M)
M,0,0,10 <Matrix ID, row, col, data>	1 2 3 4 <data data data data> (row 1)
M,0,2,9	5 6 7 8
M,0,3,9	9 10 11 12
M,0,5,9	13 14 15 16
M,0,6,9	17 18 19 20
...	...

Here is the pseudo code of our matrix multiplying function:

#### High level design:

$$MN = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} a & e \\ b & f \\ c & g \end{bmatrix} = \begin{bmatrix} 1a + 2b + 3c & 1e + 2f + 3g \\ 4a + 5b + 6c & 4e + 5f + 6g \end{bmatrix}$$

From the example, we find out that first row of matrix M repeats two times for multiplying the columns of matrix N. Therefore, we can divide the multiplications with MapReduce.

First, in map stage, we need to duplicate the rows of M and columns of N for the multiplication, and then shuffle the elements according to the position of answer matrix (i, k). After the stage, we can get key-value pairs with the (full or part) list of vectors that ready for inner product at position (i, k).

Then, in reduce stage, we perform the inner product for each key (i, k), and get the result.

#### Notations:

M, N: ID of first matrix (Our goal is MN)

M: i x j, N: j x k (where i, j, k are size of the matrices)

M<sub>ij</sub>, N<sub>jk</sub>: elements in M, N matrix

ANS<sub>ik</sub>: elements of answer matrix

#### Mapping Function

( input: [M, i, j, M<sub>ij</sub>] or [N, j, k, N<sub>jk</sub>])

1. **if** input[ID] equals to 'M':
2.     **for** count from 1 to k:
3.         Set (key, value pairs) as ((i, count), ('M', j, M<sub>ij</sub>))
4. **else if** input[ID] equals to 'N':
5.     **for** count from 1 to i:
6.         Set (key, value pairs) as ((count, k), ('N', j, N<sub>jk</sub>))
7. **return** Set of (key, value pairs)

#### Reduce Function

input: Set of shuffled (key, value pairs):

( (i,k), [ ('M',1,M<sub>i1</sub>), ('M',2,M<sub>i2</sub>),... ('M',j,M<sub>ij</sub>), ('N',1,N<sub>1k</sub>), ('N',2,N<sub>2k</sub>),... ('N',j,N<sub>jk</sub>) ] )  
value pair: (ID, j, data)

1. Sort the values with value pair[ID] equals to 'M' by value pair[j] and get arrayM //row i of M
2. Sort the values with value pair[ID] equals to 'N' by value pair[j] and get arrayN //col i of M

3.  $ANS_{ik} \leftarrow 0$
4. **for** *index* from 1 to *j*:
5.      $ANS_{ik} \leftarrow ANS_{ik} + arrayM[index] * arrayN[index]$
6. **return** (key, element): ( (i,k),  $ANS_{ik}$  )

(j) What is a combiner? Add a combiner to the last question and explain its function.

Combiner is a “sub-reducer” between mapper and reducer. It performs the pre-processing of the data and reduces the amount of data that needs to be sent to the reducer. The combiner runs on each node, so the input of the combiner is the output of the node instead of whole data. By doing the pre-processing, combiner can reduce the size of the output, and reduce the transition cost on the network.

To add the combiner, we need to first observe that in each node, what can be pre-processed. From the pseudo-code in (i), we can see that the multiplication is handled on reducer, so we can move some of the multiplication to combiner to avoid sending whole list of **value pair** to the reducer.

Before we start, we have to know that the list of might **value pair** not be fully gathered, so we need to first check if the data is gathered in the node. We defined a new ID: “S” to note that the multiplication of (i, k) is done in combiner. Thus on reducer, we can identify which are solved and which are not.

Mapping Function
( input: [M, i, j, Mij] or [N, j, k, Njk] ) <ol style="list-style-type: none"> <li>1. <b>if</b> input[ID] equals to ‘M’:</li> <li>2.     <b>for</b> <i>count</i> from 1 to <i>k</i>:</li> <li>3.         Set (key, value pairs) as ((i, <i>count</i>), (‘M’, j, Mij))</li> <li>4. <b>else if</b> input[ID] equals to ‘N’:</li> <li>5.     <b>for</b> <i>count</i> from 1 to <i>i</i>:</li> <li>6.         Set (key, value pairs) as ((<i>count</i>, k), (‘N’, j, Njk))</li> <li>7. <b>return</b> Set of (key, value pairs)</li> </ol>
Combiner Function
input: Set of shuffled (key, value pairs): ( (i,k), [ (‘M’,1,M <sub>i1</sub> ), (‘M’,2,M <sub>i2</sub> ),... (‘M’,j,M <sub>ij</sub> ), (‘N’,1,N <sub>1k</sub> ), (‘N’,2,N <sub>2k</sub> ),... (‘N’,j,N <sub>jk</sub> ) ] ) <b>value pair</b> : (ID, j, data) <ol style="list-style-type: none"> <li>1. <b>if</b> size of (value pairs) equals to 2*j:     <b>// All elements of this pair are gathered, then do “reduce”</b></li> <li>2.     Sort the values with <b>value pair</b>[ID] equals to ‘M’ by <b>value pair</b>[j] and get <i>arrayM</i>     //row i of M</li> <li>3.     Sort the values with <b>value pair</b>[ID] equals to ‘N’ by <b>value pair</b>[j] and get <i>arrayN</i>     //col i of M</li> <li>4.     <math>ANS_{ik} \leftarrow 0</math></li> <li>5.     <b>for</b> <i>index</i> from 1 to <i>j</i>:</li> <li>6.         <math>ANS_{ik} \leftarrow ANS_{ik} + arrayM[index] * arrayN[index]</math></li> <li>7.     Set (key, value pairs) as ((i, k), (‘S’, j, <math>ANS_{ik}</math>))</li> <li>8. <b>else:</b>     // If elements are not gathered, then keep the key-value pair</li> <li>9.     Set (key, value pairs) as ((i, k), (ID, j, data))</li> <li>10. <b>return</b> Set of (key, value pairs)</li> </ol>
Reduce Function
input: Set of shuffled (key, value pairs): ( (i,k), [ (‘M’,1,M <sub>i1</sub> ), (‘M’,2,M <sub>i2</sub> ),... (‘M’,j,M <sub>ij</sub> ), (‘N’,1,N <sub>1k</sub> ), (‘N’,2,N <sub>2k</sub> ),... (‘N’,j,N <sub>jk</sub> ) ] )

```

    value pair: (ID, j, data)
1.  if value pairs [first][ID] equals to 'S':      // Means that the multiplication is done in combiner
2.    Set (key, value pairs) as ((i, k), ANSik)
3.  else:
4.    Sort the values with value pair[ID] equals to 'M' by value pair[j] and get arrayM    //row i of M
5.    Sort the values with value pair[ID] equals to 'N' by value pair[j] and get arrayN    //col i of M
6.    ANSik ← 0
7.    for index from 1 to j:
8.      ANSik ← ANSik + arrayM[index] * arrayN[index]
9.    Set (key, value pairs) as ((i, k), ANSik)
10. return Set of (key, element): ( (i,k), ANSik)

```

Ref: Compile Java to Jar: <http://www.skylit.com/javamethods/faqs/createjar.html>

Install: <https://www.linode.com/docs/databases/hadoop/how-to-install-and-set-up-hadoop-cluster/>

**We have zero tolerance to forged or fabricated data!!** A single piece of forged/fabricated data would bring the total score down to zero.