# Assignment – Map Reduce Lab
## Hajra Abdul Hai 14893



```
hajra@hajra-Inspiron-15-3567:~$ sudo docker pull cloudera/quickstart:latest
[sudo] password for hajra:
latest: Pulling from cloudera/quickstart
Image docker.io/cloudera/quickstart:latest uses outdated schema1 manifest format
. Please upgrade to a schema2 image for better future compatibility. More inform
ation at https://docs.docker.com/registry/spec/deprecated-schema-v1/
1d00652ce734: Already exists
Digest: sha256:f91bee4cdfa2c92ea3652929a22f729d4d13fc838b00f120e630f91c941acb63
Status: Image is up to date for cloudera/quickstart:latest
docker.io/cloudera/quickstart:latest
hajra@hajra-Inspiron-15-3567:~$ sudo docker run --hostname=quickstart.cloudera -
-privileged=true -t -v/home/hajra/Desktop/dataset:/user/cloudera/shared -i -p 88
89:8888 -p7180:7181 cloudera/quickstart  /usr/bin/docker-quickstart
Starting mysqld:                                              [  OK  ]

if [ "$1" == "start" ] ; then
    if [ "${EC2}" == 'true' ]; then
        FIRST_BOOT_FLAG=/var/lib/cloudera-quickstart/.ec2-key-installed
        if [ ! -f "${FIRST_BOOT_FLAG}" ]; then
            METADATA_API=http://169.254.169.254/latest/meta-data
            KEY_URL=${METADATA_API}/public-keys/0/openssh-key
            SSH_DIR=/home/cloudera/.ssh
            mkdir -p ${SSH_DIR}
            chown cloudera:cloudera ${SSH_DIR}
```
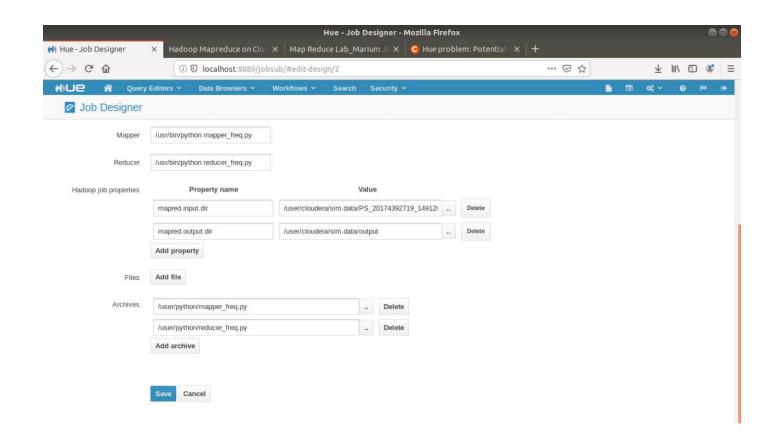


```
nohup: appending output to `nohup.out'
JMX enabled by default
Using config: /etc/zookeeper/conf/zoo.cfg
Starting zookeeper ... STARTED
starting datanode, logging to /var/log/hadoop-hdfs/hadoop-hdfs-datanode-quicksta
rt.cloudera.out
Started Hadoop datanode (hadoop-hdfs-datanode):               [  OK  ]
starting journalnode, logging to /var/log/hadoop-hdfs/hadoop-hdfs-journalnode-qu
ickstart.cloudera.out
Started Hadoop journalnode:                                   [  OK  ]
starting namenode, logging to /var/log/hadoop-hdfs/hadoop-hdfs-namenode-quicksta
rt.cloudera.out
Started Hadoop namenode:                                      [  OK  ]
starting secondarynamenode, logging to /var/log/hadoop-hdfs/hadoop-hdfs-secondar
ynamenode-quickstart.cloudera.out
Started Hadoop secondarynamenode:                             [  OK  ]

Setting HTTPFS_HOME:          /usr/lib/hadoop-httpfs
Using   HTTPFS_CONFIG:        /etc/hadoop-httpfs/conf
Sourcing:                     /etc/hadoop-httpfs/conf/httpfs-env.sh
Using   HTTPFS_LOG:           /var/log/hadoop-httpfs/
Using   HTTPFS_TEMP:          /var/run/hadoop-httpfs
Setting HTTPFS_HTTP_PORT:     14000
Setting HTTPFS_ADMIN_PORT:    14001
```

```
/etc/oozie/conf -Doozie.log.dir=/var/log/oozie -Doozie.data.dir=/var/lib/oozie -
Doozie.instance.id=quickstart.cloudera -Doozie.config.file=oozie-site.xml -Doozi
e.log4j.file=oozie-log4j.properties -Doozie.log4j.reload=10 -Doozie.http.hostnam
e=quickstart.cloudera -Doozie.admin.port=11001 -Doozie.http.port=11000 -Doozie.h
ttps.port=11443 -Doozie.base.url=http://quickstart.cloudera:11000/oozie -Doozie.
https.keystore.file=/var/lib/oozie/.keystore -Doozie.https.keystore.pass=passwor
d -Djava.library.path=:/usr/lib/hadoop/lib/native:/usr/lib/hadoop/lib/native

Using CATALINA_BASE:   /var/lib/oozie/tomcat-deployment
Using CATALINA_HOME:   /usr/lib/bigtop-tomcat
Using CATALINA_TMPDIR: /var/lib/oozie
Using JRE_HOME:        /usr/java/jdk1.7.0_67-cloudera
Using CLASSPATH:       /usr/lib/bigtop-tomcat/bin/bootstrap.jar
Using CATALINA_PID:    /var/run/oozie/oozie.pid
Starting Solr server daemon:                            [  OK  ]
Using CATALINA_BASE:   /var/lib/solr/tomcat-deployment
Using CATALINA_HOME:   /usr/lib/solr/../bigtop-tomcat
Using CATALINA_TMPDIR: /var/lib/solr/
Using JRE_HOME:        /usr/java/jdk1.7.0_67-cloudera
Using CLASSPATH:       /usr/lib/solr/../bigtop-tomcat/bin/bootstrap.jar
Using CATALINA_PID:    /var/run/solr/solr.pid
Started Impala Catalog Server (catalogd) :              [  OK  ]
Started Impala Server (impalad):                        [  OK  ]
[root@quickstart /]#
```

# Task 1

## Job Designer

### Designs

Search for design name    ▶ Submit    ✎ Edit    ⎘ Copy    ✖ Move to trash    ⌄          ⚑ View designs    ⊕ New action ⌄    🗑 View trash

#### Job Design (streaming type)

| | |
|---|---|
| Name | paymentJob1 |
| Description | Payment.Type.Frequency |

↪ advanced

You can parameterize the values, using `$myVar` or `${myVar}`. When the design is submitted, you will be prompted for the actual value of `myVar`.

| | |
|---|---|
| Mapper | /usr/bin/python mapper_freq.py |
| Reducer | /usr/bin/python reducer_freq.py |

| Hadoop job properties | Property name | Value |
|---|---|---|

ⵀⵓⵇ  🏠  Query Editors ⌄   Data Browsers ⌄   Workflows ⌄   Search   Security ⌄

📝 Job Designer

| | |
|---|---|
| Mapper | /usr/bin/python mapper_freq.py |
| Reducer | /usr/bin/python reducer_freq.py |

Hadoop job properties

| | Property name | Value | | |
|---|---|---|---|---|
| | mapred.input.dir | /user/cloudera/sim.data/PS_20174392719_14912( | .. | Delete |
| | mapred.output.dir | /user/cloudera/sim.data/output | .. | Delete |

Add property

Files    Add file

| | | | |
|---|---|---|---|
| Archives | /user/python/mapper_freq.py | .. | Delete |
| | /user/python/reducer_freq.py | .. | Delete |

Add archive

Save   Cancel

---

ⵀⵓⵇ  🏠  Query Editors ⌄   Data Browsers ⌄   Workflows ⌄   Search   Security ⌄

📝 Job Designer

## Designs

| Search for design name | ▶ Submit  ✏ Edit  ⧉ Copy  ✖ Move to trash ⌄ | ⊕ New action ⌄  🗑 View trash |
|---|---|---|

| ☐ Name | ⬍ Description | ⬍ Owner | ⬍ Type | ⬍ Status | ⬍ Last modified | ⬍ |
|---|---|---|---|---|---|---|
| ☐ paymentJob1 | Payment.Type.Frequency | cloudera | streaming | shared | December 07, 2020 01:29 PM | |

Showing 1 to 1 of 1 entries (filtered from 2 total entries)    ← Previous  1  Next →

50% - Hue - Oozie Editor/ ×    Hadoop Mapreduce on Clou ×    Map Reduce Lab_Marium Ja ×    (4) WhatsApp    ×    +

localhost:8889/oozie/list_oozie_workflow/0000004-201207103717106-oozie-oozi-W/

**HUE**    Query Editors ⌄    Data Browsers ⌄    Workflows ⌄    Search    Security ⌄

Oozie Dashboard    **Workflows**    Coordinators    Bundles    SLA    Oozie

**WORKFLOW**

paymentJob1

**SUBMITTER**

cloudera

**STATUS**

RUNNING

**PROGRESS**

50%

**ID**

0000004-201207103717106-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.application…

### Workflow paymentJob1

Graph    Actions    Details    Configuration    Log    Definition

Back

---

Hue - Oozie Editor/Dash ×    Hadoop Mapreduce on Clou ×    Map Reduce Lab_Marium Ja ×    Hue problem: Potential ×    +

localhost:8889/oozie/list_oozie_workflow/0000006-201207103717106-oozie-oozi-W/

**HUE**    Query Editors ⌄    Data Browsers ⌄    Workflows ⌄    Search    Security ⌄

Oozie Dashboard    **Workflows**    Coordinators    Bundles    SLA    Oozie

**WORKFLOW**

paymentJob1

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000006-201207103717106-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.application…

### Workflow paymentJob1

Graph    Actions    Details    Configuration    Log    Definition

Back

Hue - File Browser X | Hadoop Mapreduce on Clou X | Map Reduce Lab_Marium Ja X | Hue problem: Potential r X | +

localhost:8889/filebrowser/#/user/cloudera/sim.data/output

Hue   Query Editors ⌄   Data Browsers ⌄   Workflows ⌄   Search   Security ⌄

## File Browser

Search for file name     ⚙ Actions ⌄   ✖ Move to trash ⌄

⊕ Upload ⌄   ⊕ New ⌄

🏠 Home   / user / cloudera / sim.data / output  ✎

▼ History   🗑 Trash

| | Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|
| 📁 | ↟ | | cloudera | cloudera | drwxr-xr-x | December 07, 2020 05:30 AM |
| 📁 | . | | cloudera | cloudera | drwxr-xr-x | December 07, 2020 05:31 AM |
| 📄 | _SUCCESS | 0 bytes | cloudera | cloudera | -rw-r--r-- | December 07, 2020 05:31 AM |
| 📄 | part-00000 | 77 bytes | cloudera | cloudera | -rw-r--r-- | December 07, 2020 05:31 AM |

Show  45  ⌄  of 2 items

Page  1  of 1  |◀ ◀◀ ▶▶ ▶|

---

Hue - File Browser - part X | Hadoop Mapreduce on Clou X | Map Reduce Lab_Marium Ja X | Hue problem: Potential r X | +

localhost:8889/filebrowser/view=/user/cloudera/sim.data/output/part-00000

Hue   Query Editors ⌄   Data Browsers ⌄   Workflows ⌄   Search   Security ⌄

## File Browser

### ACTIONS

▥ View as binary

✎ Edit file

⬇ Download

🗎 View file location

🔄 Refresh

### INFO

**Last modified**
Dec. 7, 2020 5:31 a.m.
**User**
cloudera
**Group**
cloudera
**Size**
77 bytes
**Mode**
100644

🏠 Home   / user / cloudera / sim.data / output / **part-00000**

Page  1  of 1  |◀ ◀◀ ▶▶ ▶|

```
CASH_IN 1399283
CASH_OUT    2237500
DEBIT   41432
PAYMENT 2151493
TRANSFER    532909
```

| Logs | Id | Name | Type | Status | External Id | Start Time | End Time | Error Code | Error Message | Transition | Data |
|------|----|------|------|--------|-------------|------------|----------|-----------|--------------|------------|------|
| ☰ | 0000003-200511180426350-oozie-oozi-W@paymentJob1 | paymentJob1 | map-reduce | OK | job_1589220230939_0008 | Tue, 12 May 2020 05:04:09 | Tue, 12 May 2020 05:04:56 | | | end | |

## Workflow paymentJob1

Graph   Actions   **Details**   Configuration   Log   Definition

| | |
|---|---|
| Group | - |
| External Id | - |
| Last Modified | Mon, 07 Dec 2020 05:36:37 |
| Start Time | Mon, 07 Dec 2020 05:36:37 |
| Created Time | Mon, 07 Dec 2020 05:36:36 |
| End Time | - |
| Application Path | hdfs://quickstart.cloudera:8020/user/hue/oozie/workspaces/_cloudera_-oozie-2-1607347792.58 |
| Run | 0 |

03717106-   **Back**

| Name | Value |
|------|-------|
| hue-id-w | 4 |
| jobTracker | localhost:8032 |
| mapreduce.job.user.name | cloudera |
| nameNode | hdfs://quickstart.cloudera:8020 |
| oozie.use.system.libpath | true |
| oozie.wf.application.path | hdfs://quickstart.cloudera:8020/user/hue/oozie/workspaces/_cloudera_-oozie-4-1589284887.54 |
| user.name | cloudera |

Back

```
2020-12-07 13:36:37,173  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB
[0000007-201207103717106-oozie-oozi-W] ACTION[0000007-201207103717106-oozie-oozi-W@:start:] Start action [0000007-201207103717106-oozie-o
ozi-W@:start:] with user-retry state : userRetryCount [0], userRetryMax [0], userRetryInterval [10]
2020-12-07 13:36:37,185  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB
[0000007-201207103717106-oozie-oozi-W] ACTION[0000007-201207103717106-oozie-oozi-W@:start:] [***0000007-201207103717106-oozie-oozi-W@:sta
rt:***]Action status=DONE
2020-12-07 13:36:37,186  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB
[0000007-201207103717106-oozie-oozi-W] ACTION[0000007-201207103717106-oozie-oozi-W@:start:] [***0000007-201207103717106-oozie-oozi-W@:sta
rt:***]Action updated in DB!
2020-12-07 13:36:37,451  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB
[0000007-201207103717106-oozie-oozi-W] ACTION[0000007-201207103717106-oozie-oozi-W@paymentJob1] Start action [0000007-201207103717106-ooz
ie-oozi-W@paymentJob1] with user-retry state : userRetryCount [0], userRetryMax [0], userRetryInterval [10]
2020-12-07 13:36:39,199  INFO MapReduceActionExecutor:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1]
JOB[0000007-201207103717106-oozie-oozi-W] ACTION[0000007-201207103717106-oozie-oozi-W@paymentJob1] checking action, hadoop job ID [job_16
07337346855_0011] status [RUNNING]
2020-12-07 13:36:39,202  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB
[0000007-201207103717106-oozie-oozi-W] ACTION[0000007-201207103717106-oozie-oozi-W@paymentJob1] [***0000007-201207103717106-oozie-oozi-W@
```

Workflow paymentJob1

```xml
1  <workflow-app name="paymentJob1" xmlns="uri:oozie:workflow:0.4">
2      <start to="paymentJob1"/>
3      <action name="paymentJob1">
4          <map-reduce>
5              <job-tracker>${jobTracker}</job-tracker>
6              <name-node>${nameNode}</name-node>
7              <streaming>
8                  <mapper>/usr/bin/python mapper_freq.py</mapper>
9                  <reducer>/usr/bin/python reducer_freq.py</reducer>
10             </streaming>
11             <configuration>
12                 <property>
13                     <name>mapred.input.dir</name>
14                     <value>/user/cloudera/sim.data/PS_20174392719_1491204439457_log.csv</value>
15                 </property>
16                 <property>
17                     <name>mapred.output.dir</name>
18                     <value>/user/cloudera/sim.data/output</value>
19                 </property>
20             </configuration>
21             <archive>/user/python/mapper_freq.py#mapper_freq.py</archive>
22             <archive>/user/python/reducer_freq.py#reducer_freq.py</archive>
23         </map-reduce>
24         <ok to="end"/>
25         <error to="kill"/>
26     </action>
27     <kill name="kill">
28         <message>Action failed, error message[${wf:errorMessage(wf:lastErrorNode())}]</message>
29     </kill>
30     <end name="end"/>
31 </workflow-app>
32
```
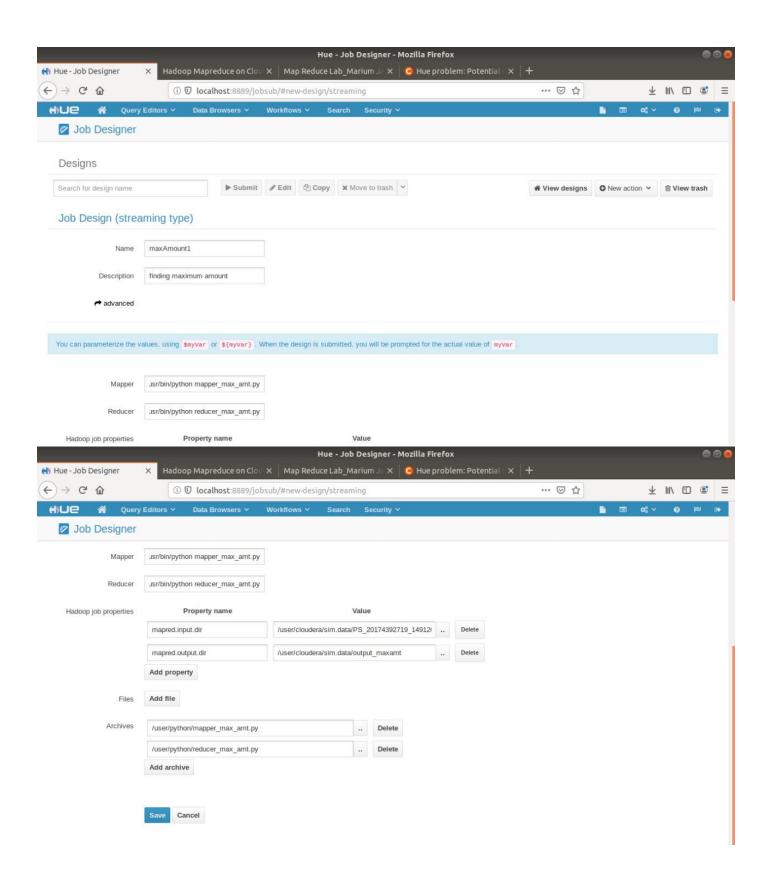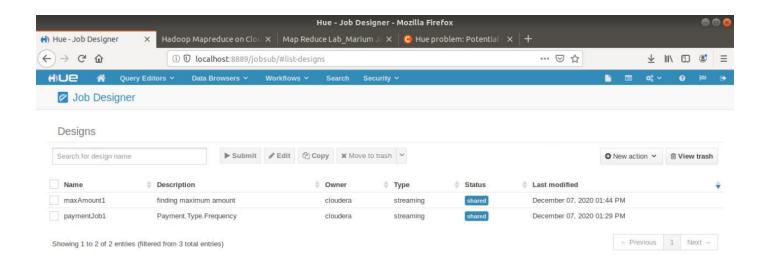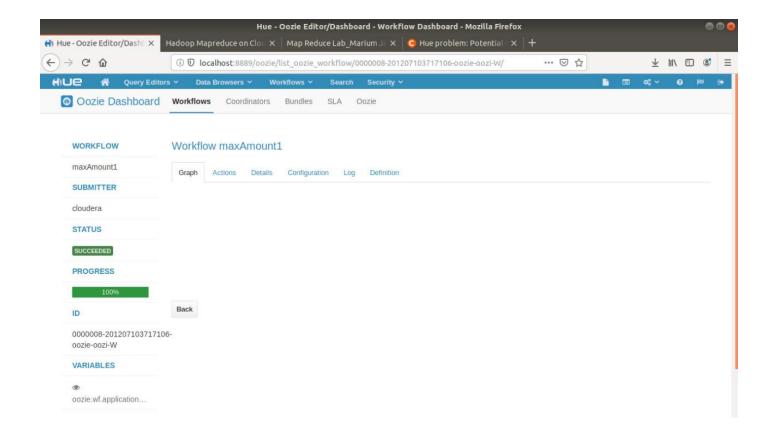
# Task 2

Hue - Job Designer ✕ | Hadoop Mapreduce on Clou ✕ | Map Reduce Lab_Marium Ja ✕ | Hue problem: Potential r ✕ | +

① ⓤ localhost:8889/jobsub/#list-designs

HUE 🏠 Query Editors ∨ Data Browsers ∨ Workflows ∨ Search Security ∨

## ✎ Job Designer

### Designs

| | Search for design name | | ▶ Submit | ✎ Edit | ⧉ Copy | ✖ Move to trash ∨ | | | | | ⊕ New action ∨ | 🗑 View trash |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| | Name | Description | Owner | Type | Status | Last modified | |
|---|---|---|---|---|---|---|---|
| ☐ | maxAmount1 | finding maximum amount | cloudera | streaming | shared | December 07, 2020 01:44 PM | |
| ☐ | paymentJob1 | Payment.Type.Frequency | cloudera | streaming | shared | December 07, 2020 01:29 PM | |

Showing 1 to 2 of 2 entries (filtered from 3 total entries)

← Previous  1  Next →

---

Hue - Oozie Editor/Dashb ✕ | Hadoop Mapreduce on Clou ✕ | Map Reduce Lab_Marium Ja ✕ | Hue problem: Potential r ✕ | +

① ⓤ localhost:8889/oozie/list_oozie_workflow/0000008-201207103717106-oozie-oozi-W/

HUE 🏠 Query Editors ∨ Data Browsers ∨ Workflows ∨ Search Security ∨

## ⊙ Oozie Dashboard  **Workflows**  Coordinators  Bundles  SLA  Oozie

**WORKFLOW**

maxAmount1

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000008-201207103717106-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.application…

### Workflow maxAmount1

Graph  Actions  Details  Configuration  Log  Definition

Back

Hue - File Browser - part- ✕ | Hadoop Mapreduce on Clou ✕ | Map Reduce Lab_Marium Ja ✕ | C Hue problem: Potential r ✕ | +

← → C ⟳ | ⓘ ⑦ localhost:8889/filebrowser/view=/user/cloudera/sim.data/output_maxamt/part-00000 | ⋯ ♡ ☆ | ↓ �III ⧉ ⑬ ≡

ℋUE 🏠 Query Editors ∨ Data Browsers ∨ Workflows ∨ Search Security ∨

📄 File Browser

**ACTIONS**

▥ View as binary

✏ Edit file

⬇ Download

📄 View file location

⟳ Refresh

**INFO**

**Last modified**
Dec. 7, 2020 5:46
a.m.
**User**
cloudera
**Group**
cloudera
**Size**
19 bytes
**Mode**
100644

🏠 Home / user / cloudera / sim.data / output_maxamt / **part-00000**     Page 1 of 1  |◄ ◄◄ ►► ►|

220696   92445516.64

localhost:8889/filebrowser/view=/user/cloudera/sim.data/output_maxamt/part-00000

# Task 3

**Q1:** Why is shuffling needed in MR? Give an example to justify.

The process of transferring data from the mappers to reducers is known as shuffling i.e. the process by which the system performs the sort, and transfers the map output to the reducer as input. So, MapReduce shuffle is necessary for the reducers, otherwise, they would not have any input (or input from every mapper).

**Q2:** Why is sorting needed in MR? Give an example to justify.

The keys generated by the mapper are automatically sorted by MapReduce Framework, i.e. before starting of reducer, all intermediate key-value pairs in MapReduce that are generated by mapper get sorted by key and not by value. Values passed to each reducer are not sorted; they can be in any order.
Sorting in Hadoop helps reducer to easily distinguish when a new reduce task should start. This saves time for the reducer. Reducer starts a new reduce task when the next key in the sorted input data is different than the previous.