# Map Reduce Lab
# Marium Jamal – 14881

## Before Tasks – downloading process:

## Task 1:

**Hue**    Query Editors ∨    Data Browsers ∨    Workflows ∨    Search    Security ∨
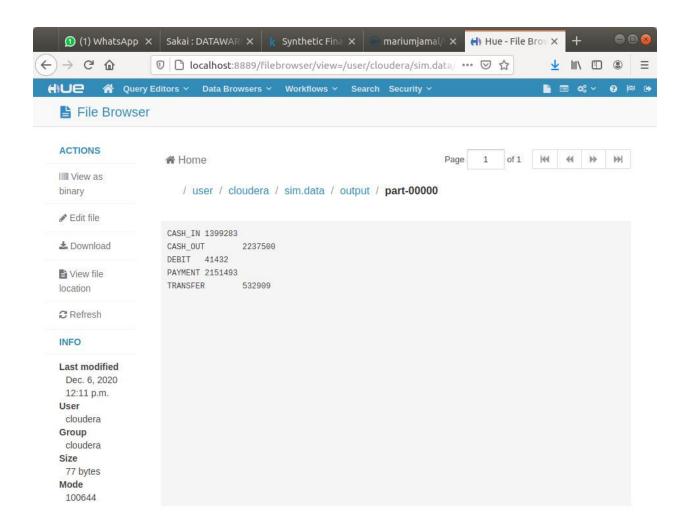
Oozie Dashboard    **Workflows**    Coordinators    Bundles    SLA    Oozie

Workflow paymentJob1

**WORKFLOW**

| **Graph** | Actions | Details | Configuration | Log | Definition |

paymentJob1

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

Back

**ID**

0000000-201206194654538-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.applicatio…

**MANAGE**

---

**Hue**    Query Editors ∨    Data Browsers ∨    Workflows ∨    Search    Security ∨

📄 File Browser

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

🏠 Home    / user / cloudera / sim.data / **output** ✏      ▾ History   🗑 Trash

| | Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|
| 📁 | ⬆ | | cloudera | cloudera | drwxr-xr-x | December 06, 2020 12:09 PM |
| 📁 | . | | cloudera | cloudera | drwxr-xr-x | December 06, 2020 12:11 PM |
| 📄 | _SUCCESS | 0 bytes | cloudera | cloudera | -rw-r--r-- | December 06, 2020 12:11 PM |
| 📄 | part-00000 | 77 bytes | cloudera | cloudera | -rw-r--r-- | December 06, 2020 12:11 PM |

Show 45 of 2 items      Page 1 of 1   ⏮ ⏪ ⏩ ⏭

**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. 6, 2020
12:11 p.m.
**User**
cloudera
**Group**
cloudera
**Size**
77 bytes
**Mode**
100644

🏠 Home

Page [ 1 ] of 1   |◀ ◀◀ ▶▶ ▶|

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

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

**Task 2:**



Job Designer

**Job Design (streaming type)**

| | |
|---|---|
| Name | maxAmount1 |
| Description | finding maximum amount |

↱ 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_max_amt.p |
| Reducer | /usr/bin/python reducer_max_amt.p |

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_maxamt | .. | Delete |

Add property

Files    Add file

Archives

| | | |
|---|---|---|
| /user/python/mapper_max_amt.py | .. | Delete |
| /user/python/reducer_max_amt.py | .. | Delete |

Add archive

Save    Cancel

---

**Oozie Dashboard**    **Workflows**    Coordinators    Bundles    SLA    Oozie

**WORKFLOW**

maxAmount1

**Workflow maxAmount1**

| Graph | Actions | Details | Configuration | Log | Definition |
|---|---|---|---|---|---|

**SUBMITTER**

cloudera

**STATUS**

RUNNING

**PROGRESS**

50%

Back

**ID**

0000001-201206194654538-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.applic…

**MANAGE**

**Oozie Dashboard**    **Workflows**    Coordinators    Bundles    SLA    Oozie

**WORKFLOW**

maxAmount1

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000001-201206194654538-
oozie-oozi-W

**VARIABLES**

👁
oozie.wf.applic…

**MANAGE**

## Workflow maxAmount1

Graph    Actions    Details    Configuration    Log    Definition

Back

---

## 📄 File Browser

**ACTIONS**

⬛ View as binary

✏ Edit file

⬇ Download

📄 View file location

🔄 Refresh

**INFO**

**Last modified**
Dec. 6, 2020
12:27 p.m.
**User**
cloudera
**Group**
cloudera
**Size**
19 bytes
**Mode**
100644

🏠 Home          Page 1 of 1

/ user / cloudera / sim.data / output_maxamt / **part-00000**
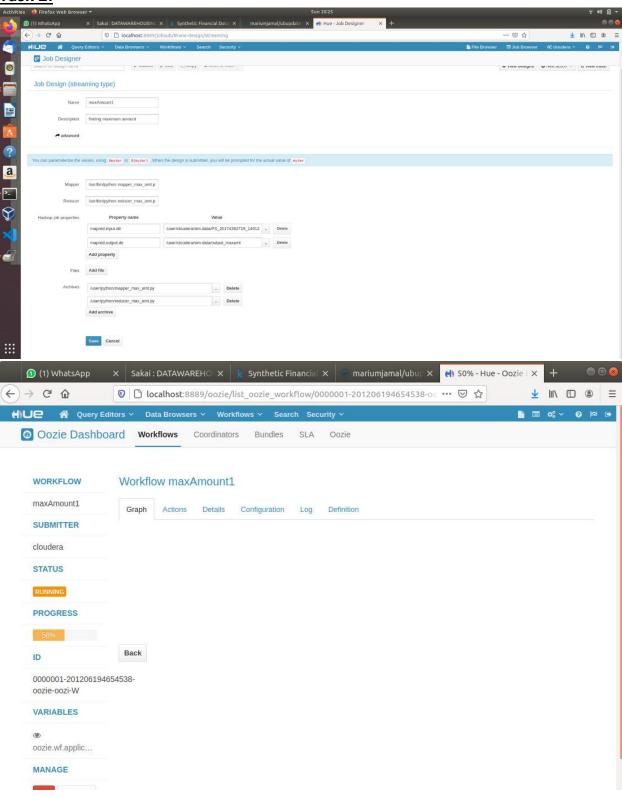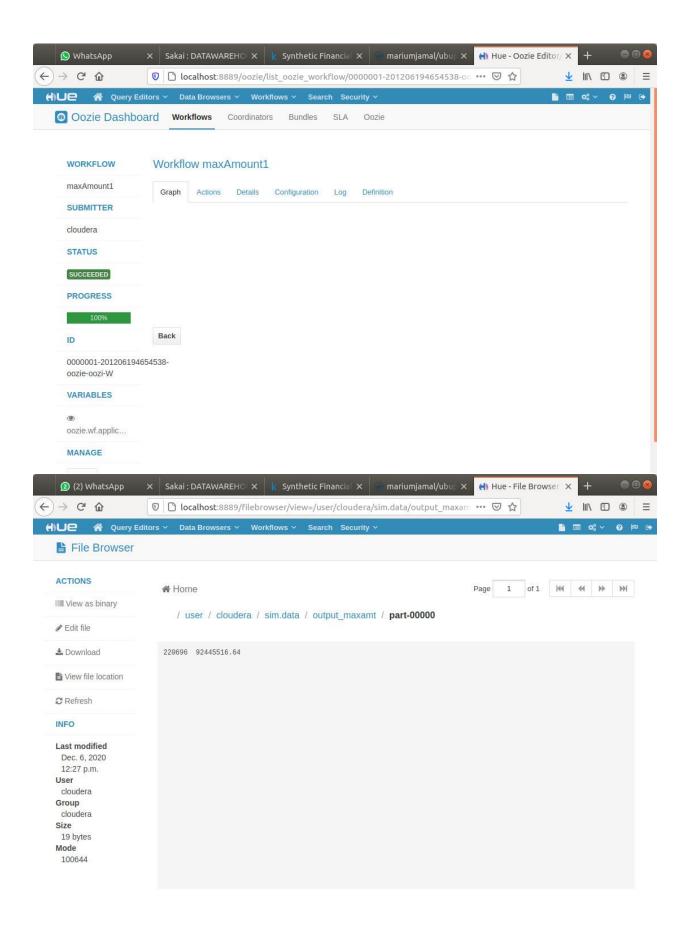
```
220696   92445516.64
```

## Task 3:

**a)** **Why is shuffling needed in MR? Give an example to justify.**

Shuffling is required in MR to transfer data from the mappers to the reducer because it is the process in which sorted intermediate output from mappers is transferred as the input to the reducer; otherwise, reducers would not have any input. An example of shuffling is given below.

**b)** **Why is sorting needed in MR? Give an example to justify.**

Sorting is needed in MR because it helps reducer to distiguish when a new reduce task should start i.e. it identifies when the next key in the reducer's input is different than the previous. An example for sorting is given below.

## SORT and SHUFFLE