# Map Reducer Lab

## Docker Installation



## 1st Map Reducer Output

Activities  Firefox Web Browser ▾                                                                    Sun 20:54

Hue - Oozie Editor/Dash ×    Hue - 403 - CSRF error  ×    Hadoop Mapreduce on Clou ×   +

← → C ⌂          localhost:8889/oozie/list_oozie_workflow/0000000-201206150920391-oozie-oozi-W/      … �remote ☆         ↓ ⦙⦙\ ⏺ ⊕ ＡＢＰ  ≡

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

Oozie Dashboard   **Workflows**   Coordinators   Bundles   SLA   Oozie

**WORKFLOW**        Workflow paymentJob1

paymentJob1          Graph   Actions   Details   Configuration   Log   Definition

**SUBMITTER**

cloudera

**STATUS**                                                                         Error  Error
                    Logs  Id                  Name         Type    Status  External Id              Start Time  End Time   Code  Message  Transition  Data
SUCCEEDED
                    ▨    0000000-201206150920391-  paymentJob1  map-    OK     job_1607267270571_0002  Sun, 06    Sun, 06                        end
**PROGRESS**              oozie-oozi-W@paymentJob1                reduce                                Dec 2020   Dec 2020
                                                                                                        07:47:58   07:51:15
        100%

**ID**
                    Back
0000000-201206150920391-
oozie-oozi-W

**VARIABLES**

👁
oozie.wf.application....

**MANAGE**

---

Activities  Firefox Web Browser ▾                                                                    Sun 20:54

Hue - Oozie Editor/Dash ×    Hue - 403 - CSRF error  ×    Hadoop Mapreduce on Clou ×   +

← → C ⌂          localhost:8889/oozie/list_oozie_workflow/0000000-201206150920391-oozie-oozi-W/      … ☑ ☆         ↓ ⦙⦙\ ⏺ ⊕ ＡＢＰ  ≡
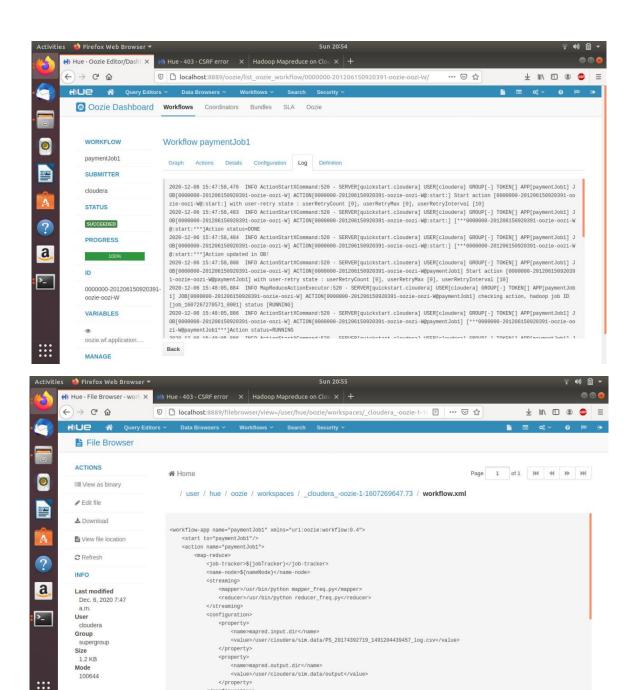
HUE  🏠  Query Editors ▾   Data Browsers ▾   Workflows ▾   Search   Security ▾

Oozie Dashboard   **Workflows**   Coordinators   Bundles   SLA   Oozie

**WORKFLOW**        Workflow paymentJob1

paymentJob1          Graph   Actions   Details   **Configuration**   Log   Definition

**SUBMITTER**

cloudera             Name                        Value

**STATUS**           hue-id-w                    1

SUCCEEDED            jobTracker                  localhost:8032

**PROGRESS**         mapreduce.job.user.name     cloudera

                     nameNode                    hdfs://quickstart.cloudera:8020
        100%
                     oozie.use.system.libpath    true
**ID**
                     oozie.wf.application.path   hdfs://quickstart.cloudera:8020/user/hue/oozie/workspaces/_cloudera_-oozie-1-1607269647.73
0000000-201206150920391-
oozie-oozi-W         user.name                   cloudera

**VARIABLES**        Back

👁
oozie.wf.application....

**MANAGE**

---

Activities  Firefox Web Browser ▾                                                                    Sun 20:54

Hue - Oozie Editor/Dash ×    Hue - 403 - CSRF error  ×    Hadoop Mapreduce on Clou ×   +

← → C ⌂          localhost:8889/oozie/list_oozie_workflow/0000000-201206150920391-oozie-oozi-W/      … ☑ ☆         ↓ ⦙⦙\ ⏺ ⊕ ＡＢＰ  ≡

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

Oozie Dashboard   **Workflows**   Coordinators   Bundles   SLA   Oozie

**WORKFLOW**        Workflow paymentJob1

paymentJob1          Graph   Actions   **Details**   Configuration   Log   Definition

**SUBMITTER**

cloudera             Group             -

**STATUS**           External Id       -

SUCCEEDED            Last Modified     Sun, 06 Dec 2020 07:47:58

**PROGRESS**         Start Time        Sun, 06 Dec 2020 07:47:58

                     Created Time      Sun, 06 Dec 2020 07:47:57
        100%
                     End Time          -
**ID**
                     Application Path  hdfs://quickstart.cloudera:8020/user/hue/oozie/workspaces/_cloudera_-oozie-1-1607269647.73
0000000-201206150920391-
oozie-oozi-W         Run               0

**VARIABLES**        Back

👁
oozie.wf.application....

**MANAGE**

Screenshot 1 — Hue Oozie Dashboard, Workflow paymentJob1, Log tab:

**Oozie Dashboard** — Workflows | Coordinators | Bundles | SLA | Oozie

**Workflow paymentJob1**

Graph | Actions | Details | Configuration | Log | Definition

WORKFLOW
paymentJob1

SUBMITTER
cloudera

STATUS
SUCCEEDED

PROGRESS
100%

ID
0000000-201206150920391-oozie-oozi-W

VARIABLES
oozie.wf.application....

MANAGE

```
2020-12-06 15:47:58,476  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB[0000000-201206150920391-oozie-oozi-W] ACTION[0000000-201206150920391-oozie-oozi-W@:start:] Start action [0000000-201206150920391-oozie-oozi-W@:start:] with user-retry state : userRetryCount [0], userRetryMax [0], userRetryInterval [10]
2020-12-06 15:47:58,483  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB[0000000-201206150920391-oozie-oozi-W] ACTION[0000000-201206150920391-oozie-oozi-W@:start:] [***0000000-201206150920391-oozie-oozi-W@:start:***]Action status=DONE
2020-12-06 15:47:58,484  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB[0000000-201206150920391-oozie-oozi-W] ACTION[0000000-201206150920391-oozie-oozi-W@:start:] [***0000000-201206150920391-oozie-oozi-W@:start:***]Action updated in DB!
2020-12-06 15:47:58,808  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB[0000000-201206150920391-oozie-oozi-W] ACTION[0000000-201206150920391-oozie-oozi-W@paymentJob1] Start action [0000000-201206150920391-oozie-oozi-W@paymentJob1] with user-retry state : userRetryCount [0], userRetryMax [0], userRetryInterval [10]
2020-12-06 15:48:05,884  INFO MapReduceActionExecutor:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB[0000000-201206150920391-oozie-oozi-W] ACTION[0000000-201206150920391-oozie-oozi-W@paymentJob1] checking action, hadoop job ID [job_1607267270571_0001] status [RUNNING]
2020-12-06 15:48:05,886  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] JOB[0000000-201206150920391-oozie-oozi-W] ACTION[0000000-201206150920391-oozie-oozi-W@paymentJob1] [***0000000-201206150920391-oozie-oozi-W@paymentJob1***]Action status=RUNNING
2020-12-06 15:48:05,886  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[paymentJob1] J
```

Back



Screenshot 2 — Hue File Browser, workflow.xml:

**File Browser**

ACTIONS
- View as binary
- Edit file
- Download
- View file location
- Refresh

INFO
Last modified
Dec. 6, 2020 7:47 a.m.
User
cloudera
Group
supergroup
Size
1.2 KB
Mode
100644

Home

/ user / hue / oozie / workspaces / _cloudera_-oozie-1-1607269647.73 / workflow.xml

```xml
<workflow-app name="paymentJob1" xmlns="uri:oozie:workflow:0.4">
    <start to="paymentJob1"/>
    <action name="paymentJob1">
        <map-reduce>
            <job-tracker>${jobTracker}</job-tracker>
            <name-node>${nameNode}</name-node>
            <streaming>
                <mapper>/usr/bin/python mapper_freq.py</mapper>
                <reducer>/usr/bin/python reducer_freq.py</reducer>
            </streaming>
            <configuration>
                <property>
                    <name>mapred.input.dir</name>
                    <value>/user/cloudera/sim.data/PS_20174392719_1491204439457_log.csv</value>
                </property>
                <property>
                    <name>mapred.output.dir</name>
                    <value>/user/cloudera/sim.data/output</value>
                </property>
```

# 2nd Map Reducer Output

Hue - Oozie Editor/Dashb ✕ | Hue - 403 - CSRF error ✕ | Hadoop Mapreduce on Clou ✕ | snip in ubuntu - Google S ✕ | +

localhost:8889/oozie/list_oozie_workflow/0000001-201206150920391-oozie-oozi-W/

**HUE**  🏠  Query Editors ▾   Data Browsers ▾   Workflows ▾   Search   Security ▾

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

maxAmount1

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000001-201206150920391-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.application....

**MANAGE**

Rerun

Graph  Actions  **Details**  Configuration  Log  Definition

| | |
|---|---|
| Group | - |
| External Id | - |
| Last Modified | Sun, 06 Dec 2020 08:08:46 |
| Start Time | Sun, 06 Dec 2020 08:08:46 |
| Created Time | Sun, 06 Dec 2020 08:08:45 |
| End Time | - |
| Application Path | hdfs://quickstart.cloudera:8020/user/hue/oozie/workspaces/_cloudera_-oozie-2-1607270916.24 |
| Run | 0 |

Back

---

Hue - Oozie Editor/Dashb ✕ | Hue - 403 - CSRF error ✕ | Hadoop Mapreduce on Clou ✕ | snip in ubuntu - Google S ✕ | +

localhost:8889/oozie/list_oozie_workflow/0000001-201206150920391-oozie-oozi-W/

**HUE**  🏠  Query Editors ▾   Data Browsers ▾   Workflows ▾   Search   Security ▾

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

maxAmount1

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000001-201206150920391-oozie-oozi-W

**VARIABLES**

👁
oozie.wf.application....

**MANAGE**

Rerun

Graph  Actions  Details  **Configuration**  Log  Definition

| Name | Value |
|---|---|
| hue-id-w | 2 |
| 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-2-1607270916.24 |
| user.name | cloudera |

Back

Hue - Oozie Editor/Dash✕   Hue - 403 - CSRF error ✕   Hadoop Mapreduce on Clou ✕   snip in ubuntu - Google S ✕   +

localhost:8889/oozie/list_oozie_workflow/0000001-201206150920391-oozie-oozi-W/

HUE   ⌂   Query Editors ▾   Data Browsers ▾   Workflows ▾   Search   Security ▾

Oozie Dashboard   **Workflows**   Coordinators   Bundles   SLA   Oozie

**maxAmount1**

**SUBMITTER**

cloudera

**STATUS**

SUCCEEDED

**PROGRESS**

100%

**ID**

0000001-201206150920391-oozie-oozi-W

**VARIABLES**

👁 oozie.wf.application....

**MANAGE**

Rerun

Graph   Actions   Details   Configuration   **Log**   Definition

```
2020-12-06 16:08:46,128  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount1] JO
B[0000001-201206150920391-oozie-oozi-W] ACTION[0000001-201206150920391-oozie-oozi-W@:start:] Start action [0000001-201206150920391-ooz
ie-oozi-W@:start:] with user-retry state : userRetryCount [0], userRetryMax [0], userRetryInterval [10]
2020-12-06 16:08:46,140  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount1] JO
B[0000001-201206150920391-oozie-oozi-W] ACTION[0000001-201206150920391-oozie-oozi-W@:start:] [***0000001-201206150920391-oozie-oozi-W
@:start:***]Action status=DONE
2020-12-06 16:08:46,141  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount1] JO
B[0000001-201206150920391-oozie-oozi-W] ACTION[0000001-201206150920391-oozie-oozi-W@:start:] [***0000001-201206150920391-oozie-oozi-W
@:start:***]Action updated in DB!
2020-12-06 16:08:46,617  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount1] JO
B[0000001-201206150920391-oozie-oozi-W] ACTION[0000001-201206150920391-oozie-oozi-W@maxAmount1] Start action [0000001-201206150920391-
oozie-oozi-W@maxAmount1] with user-retry state : userRetryCount [0], userRetryMax [0], userRetryInterval [10]
2020-12-06 16:08:48,219  INFO MapReduceActionExecutor:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount
1] JOB[0000001-201206150920391-oozie-oozi-W] ACTION[0000001-201206150920391-oozie-oozi-W@maxAmount1] checking action, hadoop job ID [j
ob_1607267270571_0003] status [RUNNING]
2020-12-06 16:08:48,225  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount1] JO
B[0000001-201206150920391-oozie-oozi-W] ACTION[0000001-201206150920391-oozie-oozi-W@maxAmount1] [***0000001-201206150920391-oozie-oozi
-W@maxAmount1***]Action status=RUNNING
2020-12-06 16:08:48,225  INFO ActionStartXCommand:520 - SERVER[quickstart.cloudera] USER[cloudera] GROUP[-] TOKEN[] APP[maxAmount1] JO
```

Back

---

Graph   Actions   Details   Configuration   Log   **Definition**

```xml
1  <workflow-app name="maxAmount1" xmlns="uri:oozie:workflow:0.4">
2      <start to="maxAmount1"/>
3      <action name="maxAmount1">
4          <map-reduce>
5              <job-tracker>${jobTracker}</job-tracker>
6              <name-node>${nameNode}</name-node>
7              <streaming>
8                  <mapper>/usr/bin/python mapper_max_amt.py</mapper>
9                  <reducer>/usr/bin/python reducer_max_amt.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_maxamount</value>
19                 </property>
20             </configuration>
21             <archive>/user/python/mapper_max_amt.py#mapper_max_amt.py</archive>
22             <archive>/user/python/reducer_max_amt.py#reducer_max_amt.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
```
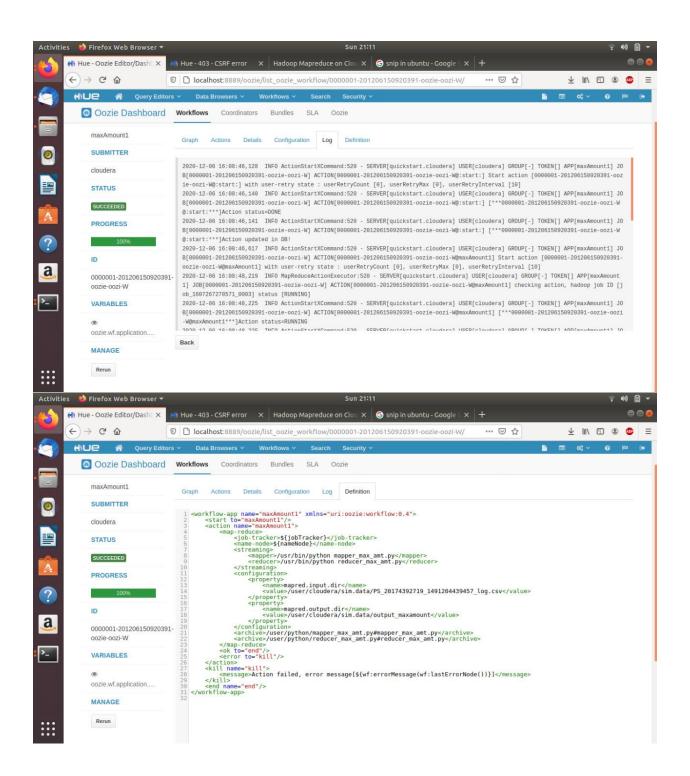
Q3: The process of transferring data from mappers to reducers is known as shuffling, the process the process by which the system performs the sort and transfers the map output to the reducer as input. This step is necessary for reducers otherwise they wont have any input to work on.

Q4: Sorting in Hadoop helps reducer to easily distinguish when a new reduce task should start. It does it by checking when the next key is different from the previous one in input sorted data.