# EEL 6761
## Assign1
## Word Count

**Github link: https://github.com/jeness/cloudComputingAssignment1**

**Haoran Yu**
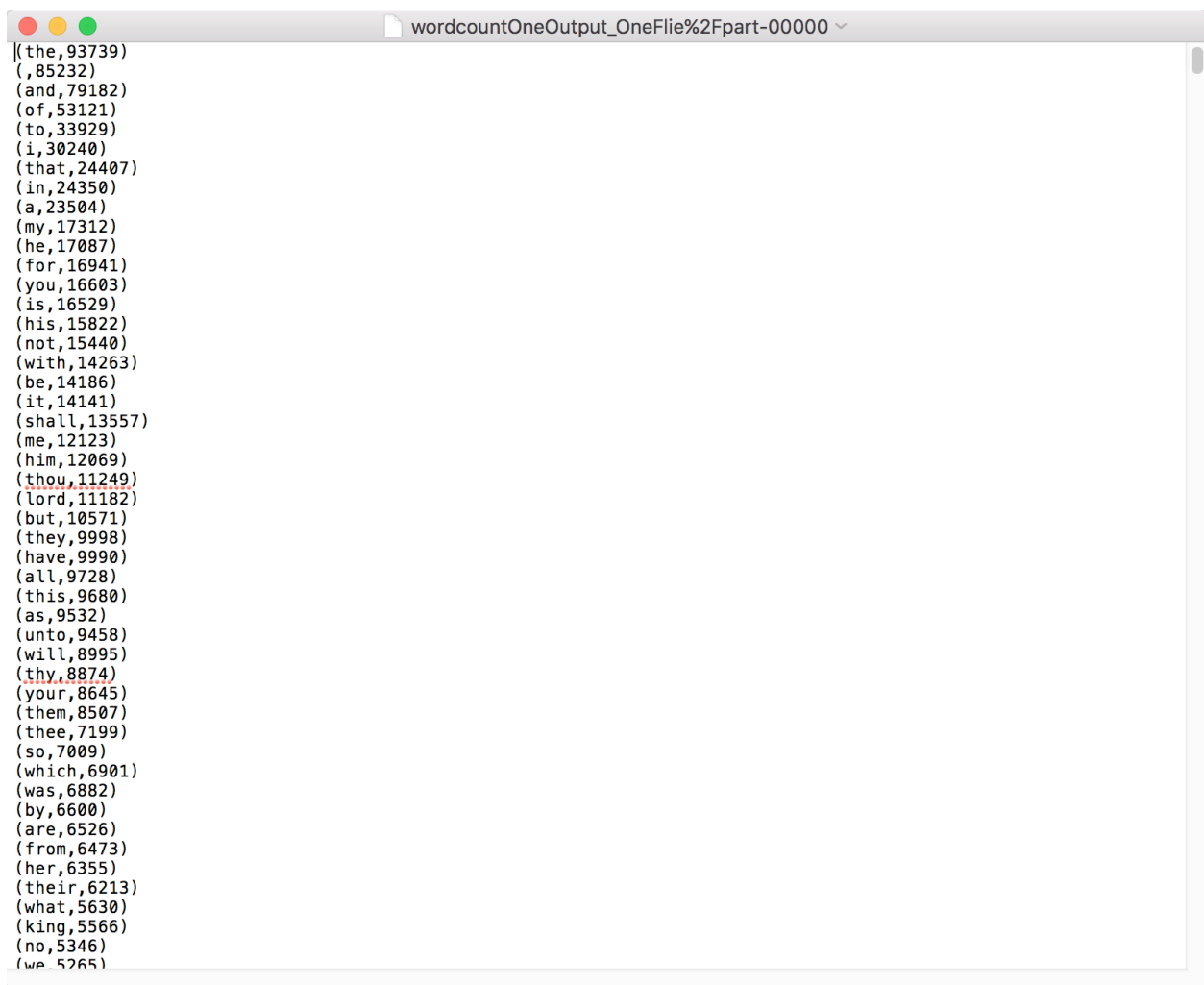**UFID:18106994**

Task1
1. Single word count
First, load bible data from google cloud storage bucket to RDD, split words with spaces, and then use filter to get non-space elements. Second, map elements to key-value pair like, <word, 1>. Third, reduce pair by the same key and get the result of counting words in descending order.

2. Result
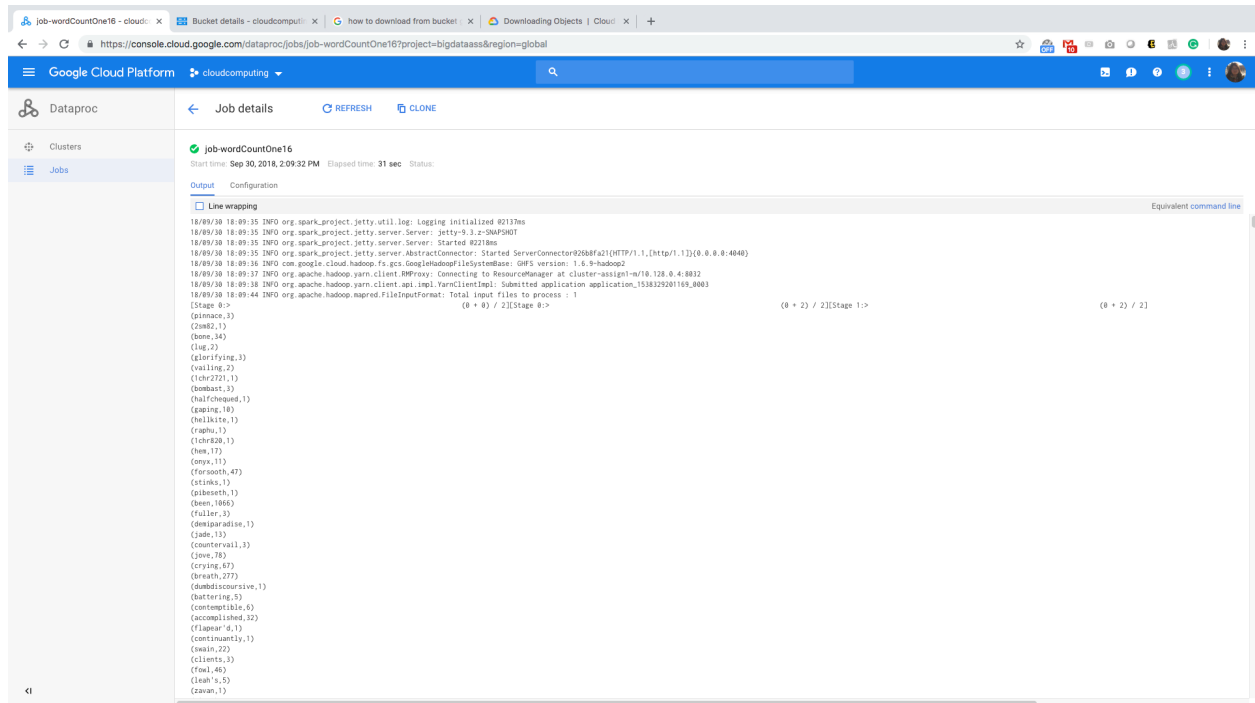Screenshots shows the result file of single word count which is download from google cloud storage.
File link:
https://storage.googleapis.com/ufcloudcomputing/wordcountOneOutput_OneFlie/part-00000



```
(the,93739)
(,85232)
(and,79182)
(of,53121)
(to,33929)
(i,30240)
(that,24407)
(in,24350)
(a,23504)
(my,17312)
(he,17087)
(for,16941)
(you,16603)
(is,16529)
(his,15822)
(not,15440)
(with,14263)
(be,14186)
(it,14141)
(shall,13557)
(me,12123)
(him,12069)
(thou,11249)
(lord,11182)
(but,10571)
(they,9998)
(have,9990)
(all,9728)
(this,9680)
(as,9532)
(unto,9458)
(will,8995)
(thy,8874)
(your,8645)
(them,8507)
(thee,7199)
(so,7009)
(which,6901)
(was,6882)
(by,6600)
(are,6526)
(from,6473)
(her,6355)
(their,6213)
(what,5630)
(king,5566)
(no,5346)
(we,5265)
```

3. Output console screenshot

Task 2

1. Double word count

First, load bible data from google cloud storage bucket to RDD, split words with spaces, then use function to get current word and the word after current word, and store every new double words combination to a new list. Second, map elements to key-value pair like, <<word 1, word 2>, 1>. Third, reduce pair by the same key and get the result of counting words with descending order.

2. Result

Screenshots shows the result file of double word count which is download from google cloud storage.

File Link:

https://storage.googleapis.com/ufcloudcomputing/wordcountTwoOutputFolder/part-00000

3. Console output screenshot



Task 3

1. Word count with distributed cache

First, load bible data and small list data from google cloud storage bucket to 2 different RDD, split words with spaces, and then use filter to get non-space elements. Second, map bible data
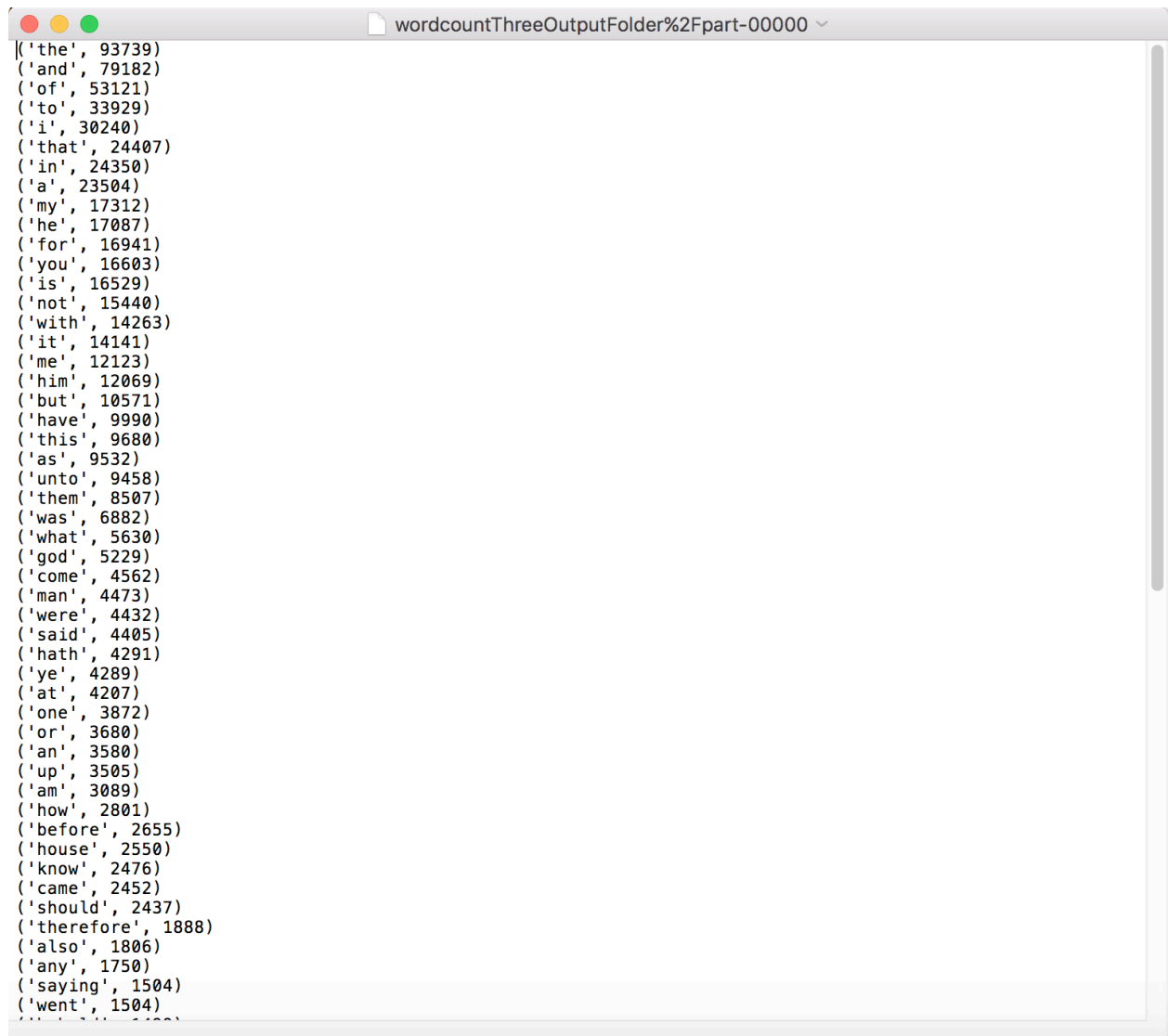
to key-value pair like, <word, 1> and reduce pair by the same key and get the result of counting words in bible. Third, map small list data to key-value pair like, <word, 0> and join two RDD and get the result of common words and the respect value from bible RDD with descending order.

2. Result

Screenshots shows the result file of distributed cache word count which is download from google cloud storage.

Link file:

```
●●●                          wordcountThreeOutputFolder%2Fpart-00000 ∨
('the', 93739)
('and', 79182)
('of', 53121)
('to', 33929)
('i', 30240)
('that', 24407)
('in', 24350)
('a', 23504)
('my', 17312)
('he', 17087)
('for', 16941)
('you', 16603)
('is', 16529)
('not', 15440)
('with', 14263)
('it', 14141)
('me', 12123)
('him', 12069)
('but', 10571)
('have', 9990)
('this', 9680)
('as', 9532)
('unto', 9458)
('them', 8507)
('was', 6882)
('what', 5630)
('god', 5229)
('come', 4562)
('man', 4473)
('were', 4432)
('said', 4405)
('hath', 4291)
('ye', 4289)
('at', 4207)
('one', 3872)
('or', 3680)
('an', 3580)
('up', 3505)
('am', 3089)
('how', 2801)
('before', 2655)
('house', 2550)
('know', 2476)
('came', 2452)
('should', 2437)
('therefore', 1888)
('also', 1806)
('any', 1750)
('saying', 1504)
('went', 1504)
```

3. Console output screenshot

```
ruhaoran99688cluster-assign1-m:~$ gcloud dataproc jobs submit pyspark    --cluster cluster-assign1 --region global    gs://ufcloudcomputing/wordCountThree.py
Job [dc0fe83fb9df40958b6f54cf7d0e2f88] submitted.
Waiting for job output...
18/09/30 22:14:27 INFO org.spark_project.jetty.util.log: Logging initialized @2602ms
18/09/30 22:14:27 INFO org.spark_project.jetty.server.Server: jetty-9.3.z-SNAPSHOT
18/09/30 22:14:27 INFO org.spark_project.jetty.server.Server: Started @2686ms
18/09/30 22:14:27 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@1c87579{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
18/09/30 22:14:27 INFO com.google.cloud.hadoop.fs.gcs.GoogleHadoopFileSystemBase: GHFS version: 1.6.9-hadoop2
18/09/30 22:14:28 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at cluster-assign1-m/10.128.0.4:8032
18/09/30 22:14:30 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1538331699096_0029
18/09/30 22:14:36 INFO org.apache.hadoop.mapred.FileInputFormat: Total input files to process : 1
18/09/30 22:14:36 INFO org.apache.hadoop.mapred.FileInputFormat: Total input files to process : 1
18/09/30 22:14:51 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@1c87579{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
Job [dc0fe83fb9df40958b6f54cf7d0e2f88] finished successfully.
driverControlFilesUri: gs://dataproc-23a43cfc-4681-427f-ae23-cf9d37bfdf1f-us/google-cloud-dataproc-metainfo/3fcaf3c6-8d28-4296-96d0-ac1ea95caf1a/jobs/dc0fe83fb9df40958b6f54cf7d0e2f88/
driverOutputResourceUri: gs://dataproc-23a43cfc-4681-427f-ae23-cf9d37bfdf1f-us/google-cloud-dataproc-metainfo/3fcaf3c6-8d28-4296-96d0-ac1ea95caf1a/jobs/dc0fe83fb9df40958b6f54cf7d0e2f88/driveroutput
placement:
  clusterName: cluster-assign1
  clusterUuid: 3fcaf3c6-8d28-4296-96d0-ac1ea95caf1a
pysparkJob:
  mainPythonFileUri: gs://ufcloudcomputing/wordCountThree.py
reference:
  jobId: dc0fe83fb9df40958b6f54cf7d0e2f88
  projectId: bigdataass
status:
  state: DONE
  stateStartTime: '2018-09-30T22:14:52.670Z'
statusHistory:
- state: PENDING
  stateStartTime: '2018-09-30T22:14:23.070Z'
- state: SETUP_DONE
  stateStartTime: '2018-09-30T22:14:23.129Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2018-09-30T22:14:23.355Z'
yarnApplications:
- name: PythonWordCount3
  progress: 1.0
  state: FINISHED
  trackingUrl: http://cluster-assign1-m:8088/proxy/application_1538331699096_0029/
```

## Configurations

1. VM instances

Three VM instances in total. One master node, and two worker nodes.



2. Submit jobs to run

## Task 1

Use web UI, see screenshot for configuration

## Task 2

Use command line tool to submit

```
gcloud dataproc jobs submit pyspark \
    --cluster cluster-assign1 --region global \
    gs://ufcloudcomputing/wordCountTwo.py
```

## Task 3

Use command line tool to submit

```
gcloud dataproc jobs submit pyspark \
    --cluster cluster-assign1 --region global \
    gs://ufcloudcomputing/wordCountThree.py
```