

## Hadoop Exercise to Create an Inverted List

For this project you will be creating an Inverted Index of words occurring in a set of English books. We'll be using a collection of 3,036 English books written by 142 authors acquired from [here](#). This collection is a small subset of the Project Gutenberg corpus that was further cleaned for the purpose of this assignment.

These books will be placed in a bucket on your Google cloud storage and the Hadoop job will be instructed to read the input from this bucket.


1. Uploading the input data into the bucket

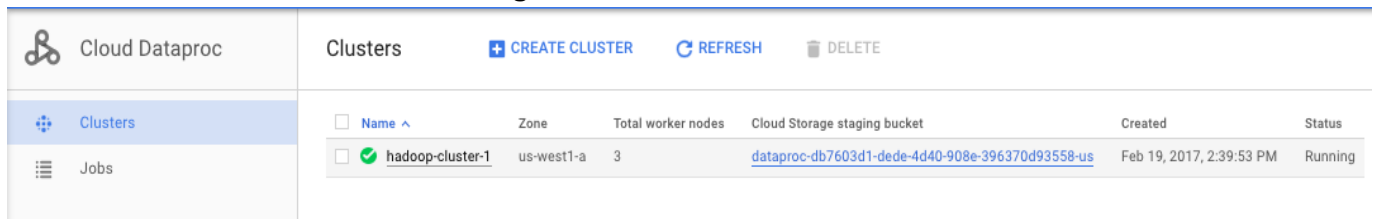
- a. Get the books from either of the links below

<http://www-scf.usc.edu/~csci572/2018Spring/hw3/DATA.zip>

<https://drive.google.com/open?id=0BxvEXzUaw-naNHNyNHBRYm5XRUE>

Use your USC account to get access to the data from the Google Drive link. The full data is around 385MB.

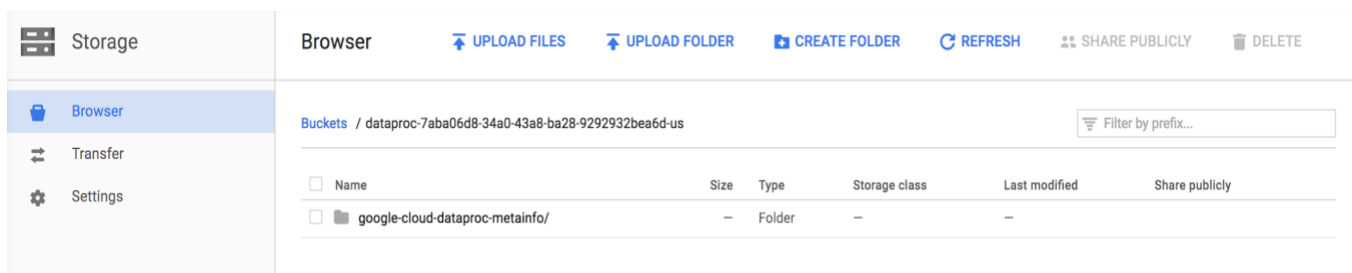
- b. Unzip the contents. You will find two folders inside named '**development**' and '**full data**'. Each of the folders contains the actual data (books) and a mapper file to map the docID to the file name. We suggest you use the development data initially while you are testing your code. Using the full data will take up to 2 hours for each run of the Map-Reduce (Hadoop) job and you may risk spending all your cloud credits while testing the code.
- c. Click on 'Dataproc' in the left navigation menu under . Next, locate the address of the default **Google cloud storage staging** bucket for your cluster. Underlined in blue in **Figure-1** below. If you've previously disabled billing, you need to re-enable it before you can upload the data. Refer to the "**Enable and Disable Billing account**" section to see how to do this.



Cloud Dataproc						
Clusters <span>+ CREATE CLUSTER</span> <span>REFRESH</span> <span>DELETE</span>						
<input type="checkbox"/>	Name ^	Zone	Total worker nodes	Cloud Storage staging bucket	Created	Status
<input checked="" type="checkbox"/>	hadoop-cluster-1	us-west1-a	3	<u>dataproc-db7603d1-dede-4d40-908e-396370d93558-us</u>	Feb 19, 2017, 2:39:53 PM	Running

**Figure 1:** The Default Cloud Storage Bucket.

- d. Go to the storage section in the left navigation bar and select your cluster's default bucket from the list of buckets (**Figure 2**). At the top you should see menu items UPLOAD FILES, UPLOAD FOLDER, CREATE FOLDER, etc. Click on the UPLOAD FOLDER button and upload the dev\_data folder and full\_data folder individually. This will take a while, but there will be a progress bar (**Figure 3**). You may not see this progress bar as soon as you start the upload but, it will show up eventually.



Storage						
Browser <span>UPLOAD FILES</span> <span>UPLOAD FOLDER</span> <span>CREATE FOLDER</span> <span>REFRESH</span> <span>SHARE PUBLICLY</span> <span>DELETE</span>						
Buckets / dataproc-7aba06d8-34a0-43a8-ba28-9292932bea6d-us <span>Filter by prefix...</span>						
<input type="checkbox"/>	Name	Size	Type	Storage class	Last modified	Share publicly
<input checked="" type="checkbox"/>	google-cloud-dataproc-metainfo/	—	Folder	—	—	

**Figure 2:** Cloud Storage Bucket.

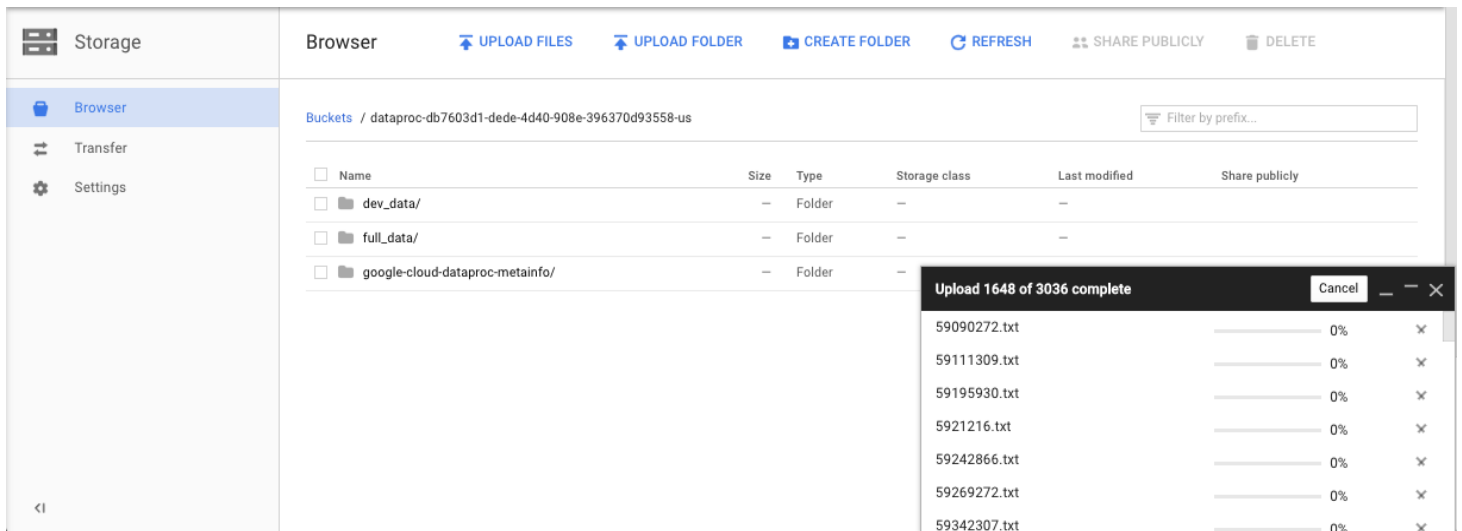


Figure 3: Progress of uploading

## Inverted Index Implementation using Map-Reduce

Now that you have the cluster and the books in place, you need to write the actual code for the job. As of now, Google Cloud allows us to submit jobs via the UI, only if they are packaged as a jar file. The following steps are focused on submitting a job written in Java via the Cloud console UI.

Refer to the below examples and write a Map-Reduce (Hadoop) job in Java that creates an Inverted Index given a collection of text files. You can very easily tweak a **word-count example** to create an inverted index instead (**Hint**: Change the mapper to output `word docID` instead of `word count` and in the reducer use a **HashMap**).

### Examples of Map-Reduce Jobs

1. <https://developer.yahoo.com/hadoop/tutorial/module4.html#wordcount>
2. <https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>

The example in the following pages explains a Hadoop word count implementation in detail. It takes one text file as input and returns the word count for every word in the file. Refer to the comments in the code for explanation.

## The Mapper Class:

```
/*
This is the Mapper class. It extends the Hadoop's Mapper class.
This maps input key/value pairs to a set of intermediate(output) key/value pairs.
Here our input key is a LongWritable and input value is a Text.
And the output key is a Text and value is an IntWritable.
*/
class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    /*
    Hadoop supported data types. This is a Hadoop specific datatype that is used to handle
    numbers and Strings in a hadoop environment. IntWritable and Text are used instead of
    Java's Integer and String datatypes.
    Here 'one' is the number of occurrences of the 'word' and is set to the value 1 during the
    Map process.
    */
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException
    {
        //Reading input one line at a time and tokenizing.
        String line = value.toString();
        StringTokenizer tokenizer = new StringTokenizer(line);

        //Iterating through all the words available in that line and forming the key value pair.
        while (tokenizer.hasMoreTokens())
        {
            word.set(tokenizer.nextToken());
            /*
            Sending to output collector(Context) which in-turn passes the output to Reducer.
            The output is as follows:
                'word1' 1
                'word1' 1
                'word2' 1
            */
            context.write(word, one);
        }
    }
}
```

## The Reducer Class:

```
/*
This is the Reducer class. It extends the Hadoop's Reducer class.
This maps the intermediate key/value pairs we get from the mapper to a set
of output key/value pairs, where the key is the word and the value is the word's count.
Here our input key is a Text and input value is a IntWritable.
And the output key is a Text and value is an IntWritable.
*/
class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable>
{
    /*
    Reduce method collects the output of the Mapper and adds the 1's to get the word's count.
    */
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException
    {
        int sum = 0;
        /*
        Iterates through all the values available with a key and add them together and give the
        final result as the key and sum of its values
        */
        for (IntWritable value : values)
        {
            sum += value.get();
        }
        context.write(key, new IntWritable(sum));
    }
}
```

## Main Class

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.*;
public class WordCount
{
    public static void main(String[] args)
        throws IOException, ClassNotFoundException, InterruptedException {
        if (args.length != 2) {
            System.err.println("Usage: Word Count <input path> <output path>");
            System.exit(-1);
        }
        //Creating a Hadoop job and assigning a job name for identification.
        Job job = new Job();
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        //The HDFS input and output directories to be fetched from the Dataproc job submission console.
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        //Providing the mapper and reducer class names.
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);
        //Setting the job object with the data types of output key(Text) and value(IntWritable).
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.waitForCompletion(true);
    }
}
```

We've already cleaned up the input data so you don't have to worry about any stray characters. Each input file consists of exactly one book that has been cleared of `\\n\\r`, `\\n` and all but one `\\t`. The only `\\t` separates the key(Document ID) from the value(Document). The input files are in a key value format as below:

DocumentID	document
------------	----------

Sample document:

```
1 51918182 four meetings by henry james 1885 i saw her only four times but i remember them
vividly she made an impression upon me i thought her very pretty and very interestinga charming
specimen of a type i am very sorry to hear of her death and yet when i think of it why should i
be sorry the last time i saw her she was certainly notbut i will describe all our meetings in
order i the first one took place in the country at a little teaparty one snowy night it must
have been some seventeen years ago my friend latouche going to spend christmas with his mother
```

The mapper's output is expected to be as follows:

```
james 51918182
people 51918182
people 51918182
of 51918182
of 51918182
```

The above example indicates that the word `james` occurred 1 time in the document with docID `51918182` and `people` 2 times.

The reducer takes this as input, aggregates the word counts using a Hashmap and creates the Inverted index. The format of the index is as follows.

word	docID:count	docID:count	docID:count...
------	-------------	-------------	----------------

```
1  ably      9931985:1
2  abnegate  85886314:1  80811098:1
3  abney     31694096:3  15109590:1  38583612:1  98115965:98
4  abnormal  47943267:1  94435826:1  80942074:1
5  abroad    73713297:1  11200532:1
```

The above sample shows the inverted index created by the reducer. The docID's can be mapped to their document names using the docID2name.csv file in the download package.

To write the Hadoop java code you can use the **VI** or **nano** editors that come pre-installed on the master node. You can test your code on the cluster itself. Be sure to use the development data while testing the code. You are expected to write a simple Hadoop job. You can just tweak [this](#) example if you'd like but, make sure you understand it first.

## Creating a jar for your code

Now that your code for the job is ready we'll need to run it. The Google Cloud console requires us to upload a Map-Reduce (Hadoop) job as a jar file. In the following example the Mapper and Reducer are in the same file called `InvertedIndexJob.java`. To create a jar for the Java class implemented please follow the instructions below. The following instructions were executed on the cluster's master node on the Google Cloud.

1. Say your Java Job file is called `InvertedIndex.java`. Create a JAR as follows:

- `hadoop com.sun.tools.javac.Main InvertedIndexJob.java`

If you get the following Note you can ignore them

Note: `InvertedIndexJob.java` uses or overrides a deprecated API.

Note: Recompile with `-Xlint:deprecation` for details.

- `jar cf invertedindex.jar InvertedIndex*.class`

Now you have a jar file for your job. You need to place this jar file in the default cloud bucket of your cluster. Just create a folder called JAR on your bucket and upload it to that folder. If you created your jar file on the cluster's master node itself use the following commands to copy it to the JAR folder.

- `hadoop fs -copyFromLocal ./invertedindex.jar`

- `hadoop fs -cp ./invertedindex.jar gs://dataproc-69070.../JAR`

The highlighted part is the default bucket of your cluster. It needs to be prepended by the `gs://` to tell the Hadoop environment that it is a bucket and not a regular location on the filesystem.

**Note:** This is not the only way to package your code into a jar file. You can follow any method that will create a single jar file that can be uploaded to the google cloud.

## Submitting the Hadoop job to your cluster

As mentioned before, a job can be submitted in two ways.

1. From the console's UI.
2. From the command line on the master node.

If you'd like to submit the job via the command line follow the instructions [here](#)

<https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>

Follow the instructions below to submit a job to the cluster via the console's UI.

1. Go to the “Jobs” section in the left navigation bar of the Dataproc page and click on “**Submit job**” (Figure 4).

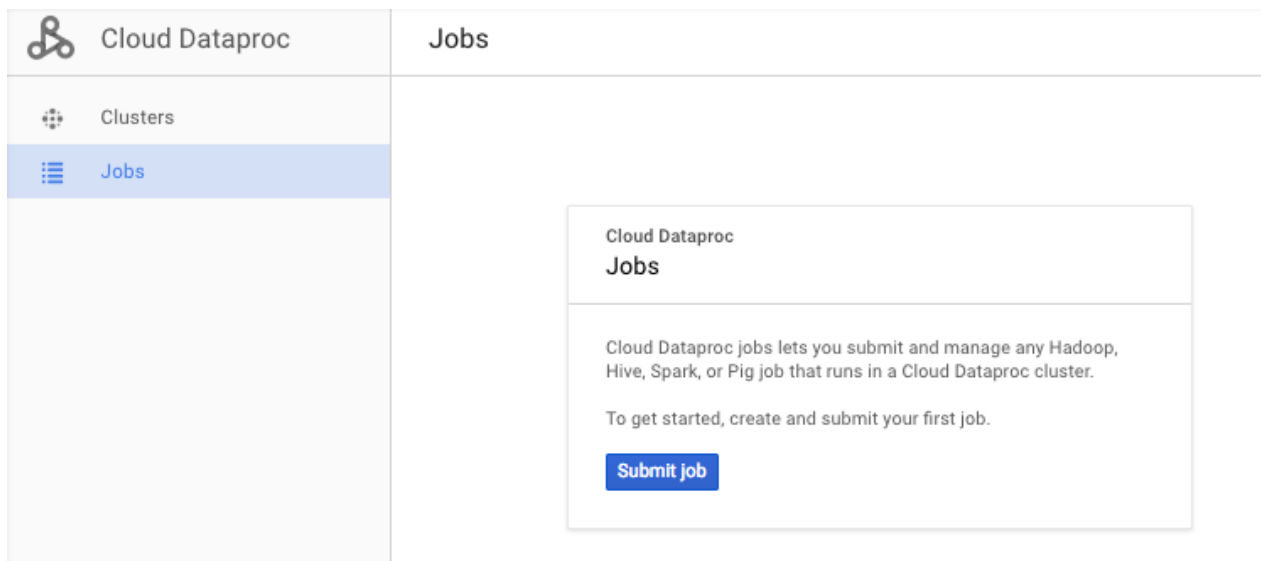


Figure 4: Dataproc jobs section

2. Fill the job parameters as follows (see **Figure 5** for reference):
  - **Cluster:** Select the cluster you created
  - **Job Type:** Hadoop
  - **Jar File:** Full path to the jar file you uploaded earlier to the Google storage bucket. Don't forget the `gs://`
  - **Main Class or jar:** The name of the java class you wrote the mapper and reducer in.
  - **Arguments:** This takes two arguments
    - i. **Input:** Path to the input data you uploaded
    - ii. **Output:** Path to the storage bucket followed by a **new** folder name. The folder is created during execution. You will get an error if you give the name of an existing folder.
  - Leave the rest at their default settings

Cloud Dataproc

← Submit a job

Clusters

Jobs

Cluster

hadoop-cluster-1

Job type

Hadoop

Jar files (Optional) ?

gs://dataproc-db7603d1-dede-4d40-908e-396370d93558-us/JAR/invertedindex.jar

Enter file path, for example, hdfs://example/example.jar

Main class or jar ?

InvertedIndexJob

Arguments (Optional) ?

gs://dataproc-db7603d1-dede-4d40-908e-396370d93558-us/dev\_data

gs://dataproc-db7603d1-dede-4d40-908e-396370d93558-us/output

Press <Return> to add more arguments

Properties (Optional) ?

+ Add item

Labels (Optional) ?

+ Add item

Submit Cancel

Equivalent [REST](#)

Figure 5: Job submission details

3. Submit Job. It will take quite a while. Please be patient. You can see the progress on the job's status section (Figure 6).

Job ID	Type	Cluster	Start time	Elapsed time	Status
<input type="checkbox"/> <a href="#">46ec16fa-5303-41ba-bb04-168fd5bbc57a</a>	Hadoop	hadoop-cluster-1	Feb 19, 2017, 5:14:20 PM	1 min 16 sec	Succeeded

Figure 6: Job ID generated. Click it to view the status of the job.

**NOTE:** If you encounter a **Java.lang.Interrupted exception** you can safely ignore it. Your submission will still execute.



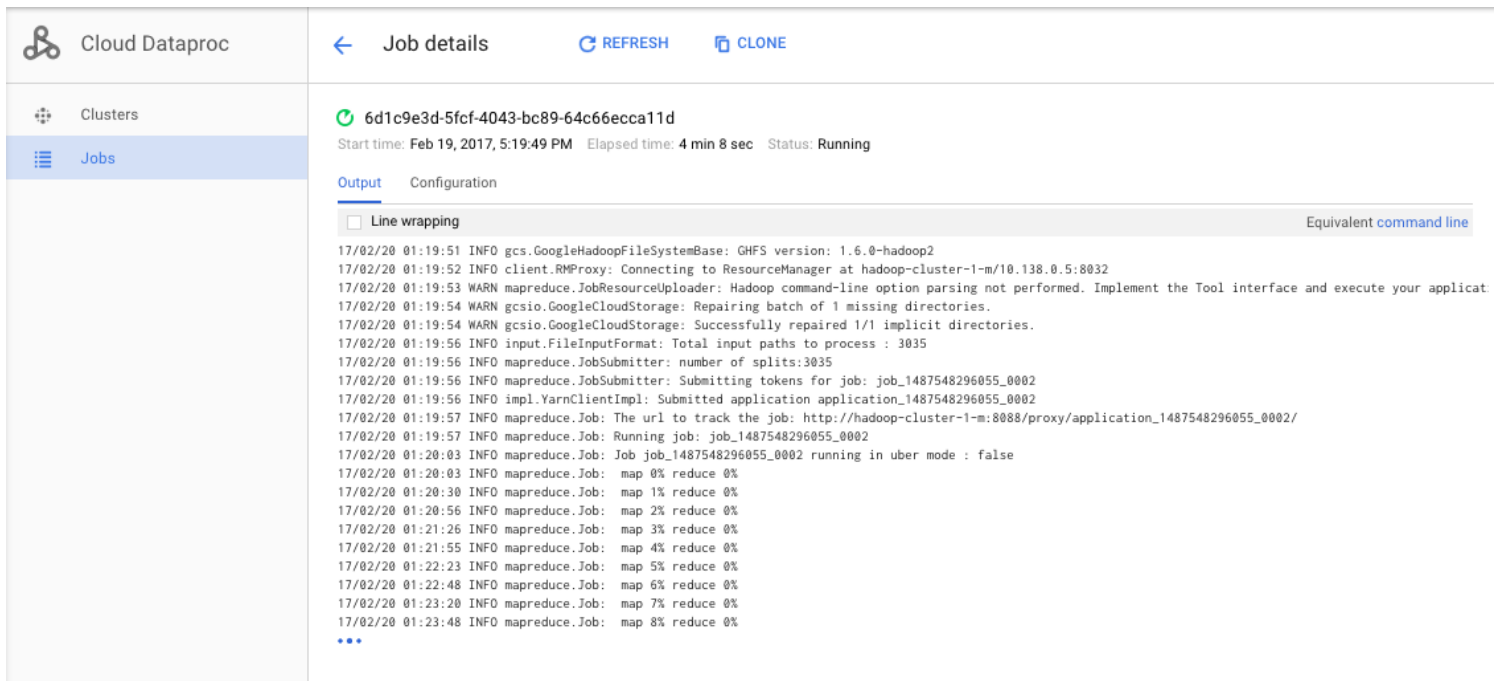


Figure 7: Job progress

4. Once the job executes copy all the log entries that were generated (**Figure 7**) to a text file called `log.txt`. You need to submit this log along with the java code. You need to do this only for the job you run on the full data. No need to submit the logs for the `dev_data`.
5. The output files will be stored in the `output` folder on the bucket. If you open this folder you'll notice that the inverted index is in several segments.(Delete the **\_SUCCESS** file in the folder before merging all the output files)

To merge the output files, run the following command in the master nodes command line(SSH)

- `hadoop fs -getmerge gs://dataproc-69070458-bbe2-.../output ./output.txt`
- `hadoop fs -copyFromLocal ./output.txt`
- `hadoop fs -cp ./output.txt gs://dataproc-69070458-bbe2-.../output.txt`

The `output.txt` file in the bucket contains the full Inverted Index for all the books.

Use `grep` to search for the words mentioned in the submissions section. Using `grep` is the fastest way to get the entries associated with the words.

For example to search for "string" use

```
grep -w '^string' fullindex.txt
```


**Note:->** Whitespace following the word (Eg;- 'little') is actually a tab rather than space

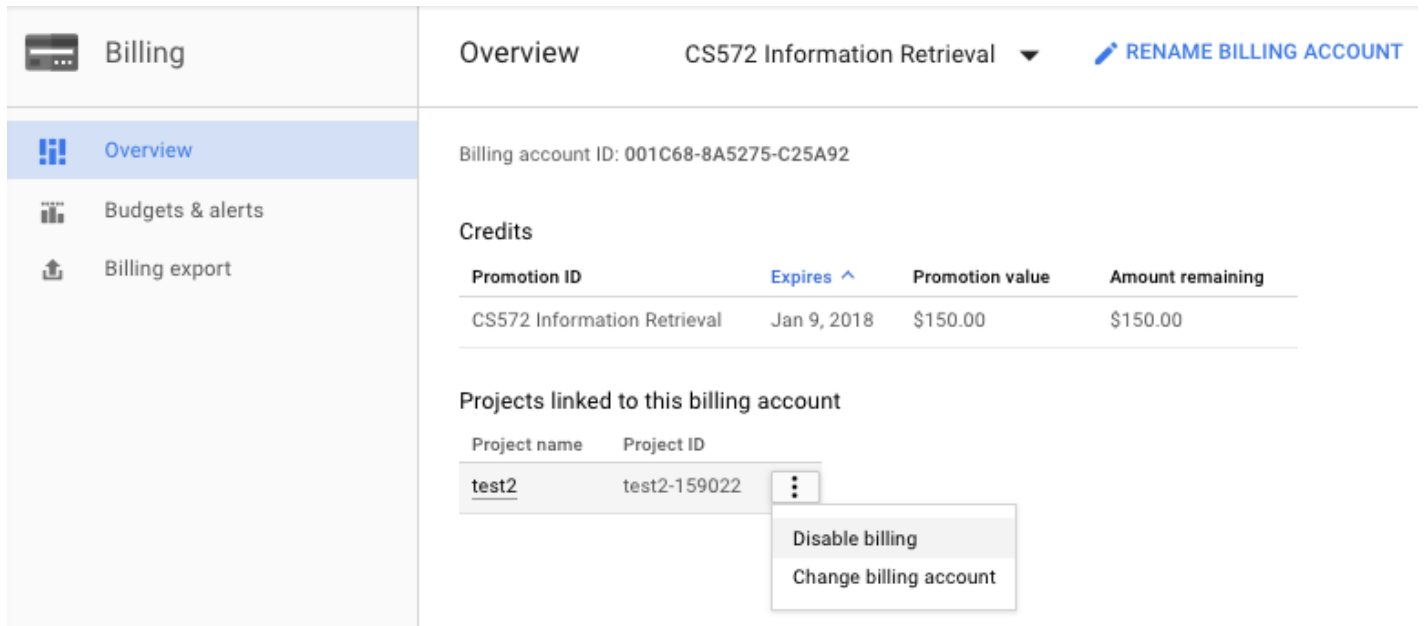


## Enabling and Disabling Billing accounts

We need to disable billing for the project (where the cluster was created) when we are not running the job to save some credits. Follow the steps below to disable and enable the billing for your project:

### Disable Billing:

1. Click the navigation button on the top left .
2. Navigate to the billing section.
3. Click on Disable billing for the project you created.(See screenshot below)

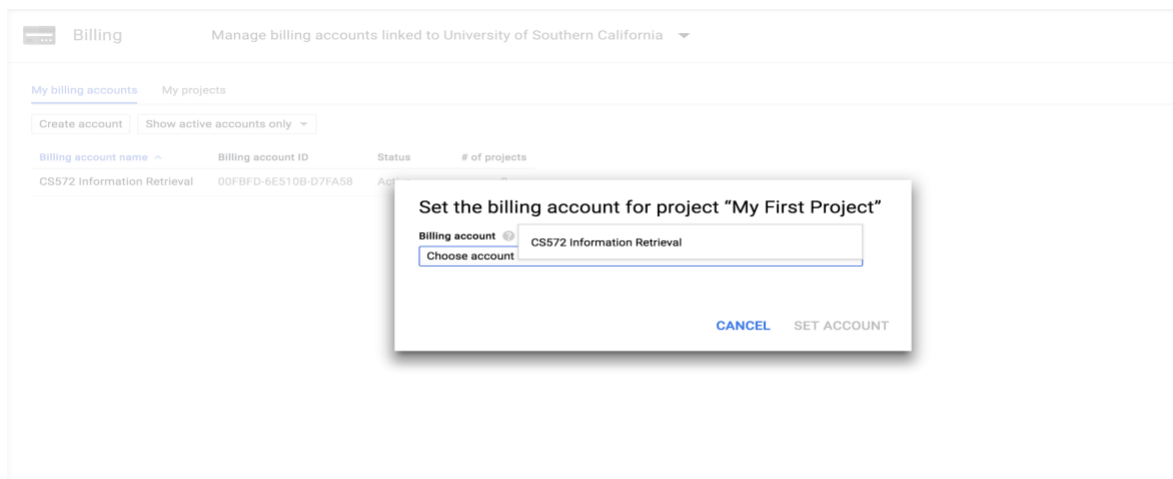


The screenshot shows the Google Cloud Billing console. On the left, the 'Billing' menu is open, showing 'Overview', 'Budgets & alerts', and 'Billing export'. The 'Overview' tab is selected. The main content area shows the billing account ID: 001C68-8A5275-C25A92. Below this, there is a 'Credits' section with a table showing a promotion for 'CS572 Information Retrieval' with a value of \$150.00. At the bottom, there is a section 'Projects linked to this billing account' with a table showing a project named 'test2' with ID 'test2-159022'. A dropdown menu is open for the 'test2' project, showing options to 'Disable billing' and 'Change billing account'.

**Figure 8:** Disabling the billing for the cluster.

### Enable Billing:

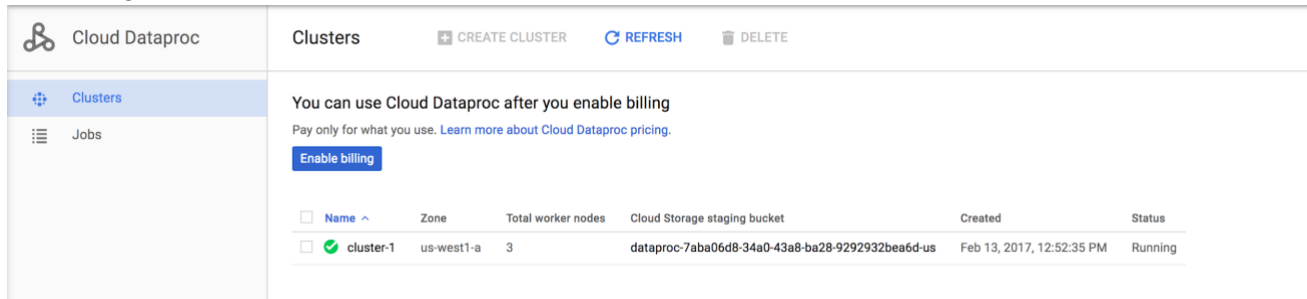
**Option 1:** When you navigate to the billing section you will be prompted to select the billing account. Select “CS572 Information Retrieval”. This billing account is created when you redeem the google credits.



The screenshot shows the Google Cloud Billing console. The top bar shows 'Billing' and 'Manage billing accounts linked to University of Southern California'. Below this, there is a table of billing accounts. A modal dialog is open, titled 'Set the billing account for project "My First Project"'. The dialog shows a list of billing accounts, with 'CS572 Information Retrieval' selected. The dialog also has a 'Choose account' button and 'CANCEL' and 'SET ACCOUNT' buttons.

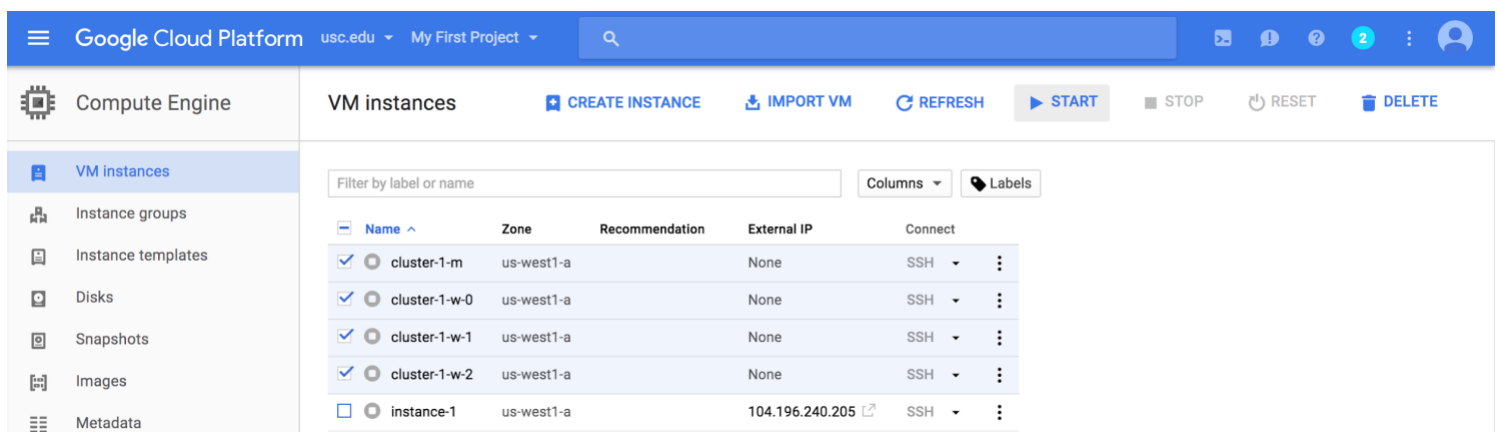
**Figure 9:** Select the account “CS572 Information Retrieval”

**Option 2:** Navigate to the Dataproc section. You will see a screen similar to the figure below. Click on Enable billing.



**Figure 10:** Enable billing

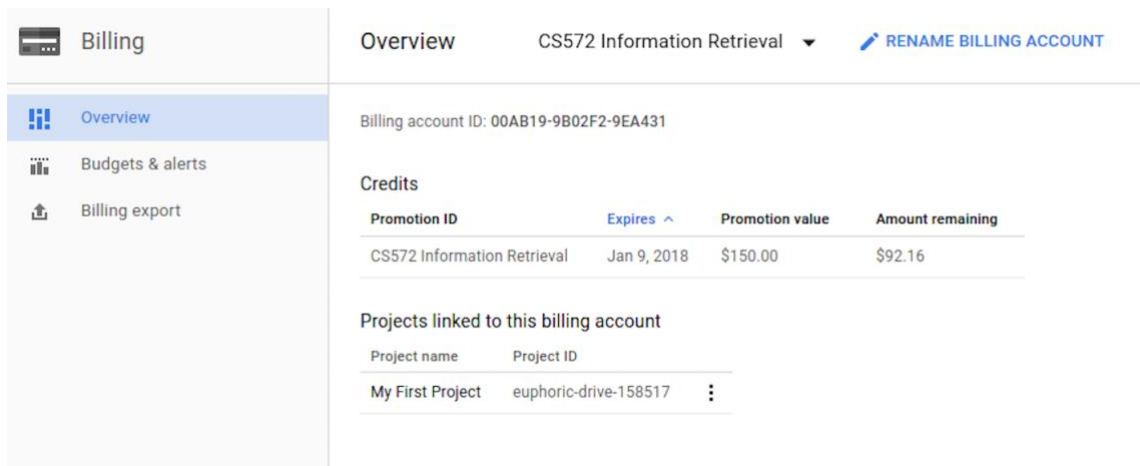
**NOTE :** Every time you disable and enable billing for a cluster, the Virtual Machines in the cluster don't start by themselves. We need to manually start the VMs. In the VM Instances section of the Cluster you might see all the VM's of the cluster disabled (See **Figure 11**). To enable the VM Instances, navigate to the Compute Engine section. Select all the instances corresponding to the cluster you created and click on the START button. Once activated navigate back to the Dataproc section to resume working on the cluster.



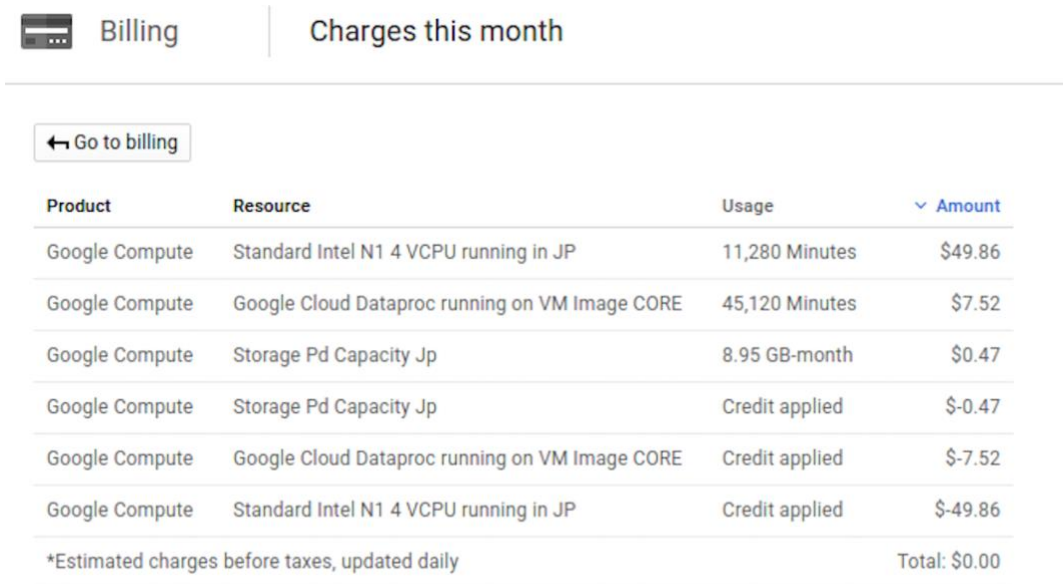
**Figure 11:** Select all virtual machines associated with the cluster.

### Credits Spent:

To check how much you've been charged for your cluster, navigate to the Billing section and click on the project name in the Overview section (see **Figure 12 & 13**). We suggest you check this section at least once every 24 hours.



**Figure 12:** Billing Overview section.



**Figure 13:** Cluster usage cost

## **Submission Instructions:**

1. Include all the code that you have written(java) and the log file created for the full data job submission.
2. Also include the inverted index file for the book “The American by Henry James ”
3. Create a text file named `index.txt` and include the index entries for the following words
  - a. little
  - b. jewel
  - c. believe
  - d. jovian
  - e. harriet
  - f. large
  - g. first
  - h. love

Add the full line from the index including the word itself.

4. Also submit a **screenshot of the output folder for the full data run in GCP.**
5. **Also submit Log file generated** from running the job on the full data.
6. **Do NOT submit your full index.**
7. Compress your code and the text file into a single zip archive and name it `index.zip`. Use a standard zip format and not zipx, rar, ace, etc.
8. To submit your file electronically to the csci572 account enter the following command from your UNIX prompt:

```
$ submit -user csci572 -tag hw3 index.zip
```

## **FAQ:**

**Q)** Can't seem to select a cluster for submitting a job?

**A)** Changing the region will do the trick

**Q)** How many files were there in `full_data` while uploading?

**A)** You need to upload .txt files only !!

The 3036th file that you see in the folder is the .DS\_Store. Number of .txt files is 3035.

It does not matter if you upload the entire folder, .DS\_Store is automatically ignored.

**Q)** Chrome suffers, in uploading almost 3000 files?

**A)** Consider opening the storage in another tab and checking the number of files. This way you will be able to know when the upload is complete.

**Q)** Do we have to use Java as the programming language?

**A)** Please go ahead and use any language binding of your choice.

**Note :** You may be on your own with language other than Java. TA's may not be able to help with other languages.

**Q)** How to Import and Export Java Projects as JAR Files in Eclipse?

**A)** <http://www.albany.edu/faculty/jmower/geog/gog692/ImportExportJARFiles.htm>

**Q)** Is it fine to submit only one .Java file, which has the all the (Mapper and Reducer Classes) inside it ?

**A)** One .java file containing your entire program should be good enough.

**Q)** Approximately how long does it take for a submitted job to finish in GCloud Dataproc?

**A)** It takes approximately 2 hours

**Q)** Should the postings list be in the sorted order of docIDs?

**A)** No need to sort the listings.

**Q)** Google cloud is not allowing to ssh?

**A)** You need to start VMs manually.

**Q)** Where can I find log files?

**A)** Cloud Dataproc -> Under Jobs

Click on one of the jobs you ran.

**Q)** Should identical words in different case be treated as same or different words?

**A)** Different words, you need not pre-process or case fold.

**Q)** How to check number of files in full\_data on storage bucket?

**A)** Go to your bucket, select the full\_data folder and click on delete. It'll list out the total files present. DO NOT PRESS DELETE in the dialog box that appears. Or run the following command from the hadoop cluster terminal:

**Hadoop fs -find gs://...//full\_data/\*.txt | wc -l**

**Q)** Different output file size?

**A)** **Scenario:-**My output size is coming out to be 186.62 MB, with each file being around 23MB. Can the output file size vary as compared to the 225MB that everyone is getting?

It shouldn't vary because data in the output files is the same.

You can perform certain sanity checks.

1) Check if your code run properly for the dev\_data?

2) Check if you used correct space / tab specifications as mentioned in the assignment description, sometimes it might be the problem with the storage space related to that.

3) You can debug with a single custom file to see, if everything is properly indexed or not.

**Q)** Different index order. Should we take the same index order (sorted) or can it be different (unsorted)?

**A)** Order does not matter. The accuracy of results is important.

**Q)** Code runs fine on development but strange file size with full data.

**A)** Check if the results produced by running on dev\_data produces huge file sizes as well. If so, that means you have to check your code. If not, check if your full\_data is uploaded correctly.

**Q)** I'm getting this error repeatedly, but I've already created the output directory and have set the argument path to that directory. Can someone help me with it?

**A)** You need to delete the output folder because the driver will attempt to create the output folder based on the argument provided.

**Q)** Am able to run the dev\_data and it is generating results. But if I ran the same code on the full data I am getting an error. The job is running for till map 25% and then it throws an error?

**A)** Please check that you have all the files uploaded just fine, and you should have 3035 files in full\_data.

**Q)** Starting VM instance failed

When I try to start the VM instances, for some of them it shows the message:

Error: Quota "CPUS" exceeded: Limit 8.0?

**A)** If you get an error saying that you've exceeded your quota, reduce the number of worker nodes or choose a Machine Type(for master and worker) with fewer vCPUs.

**Q)** Did anyone run into a situation where if you go under Dataproc > Clusters > (name of cluster instance) > VM instances > SSH, the only available option is to use another SSH client?

**A)** You probably didn't start the VM instances. Every time you disable billing and enable billing, you need to start VMs manually.

**Q)** Error enabling DataProc API

**A)** shut down project and create new one

**Q)** No space between DocID:count pairs in the output file after merge?

**A)** Happens due to copy-pasting the grep output from console to a text file. Pipe the grep output into a file and then download that file from gcloud

**Q)** "message" : "982699654446-compute@developer.gserviceaccount.com does not have storage.objects.get access to dataproc-60450493-bff5-4160-8156-fcb96702ebf0-

us/full\_data\_new/32229287.txt.",

"reason" : "forbidden"

**A)** If you're using a custom service account, you still have to give reader access to the Default service account <your-project-number>-compute@developer.gserviceaccount.com

### **Important Points:**

#P1) Output folder Number of parts generated - can be any number

#P2) Manually inspect output.txt and copy lines for the words from it and create a new txt file named index.txt. - for the 8 words

#P3) start worker nodes before submitting job to cluster

#P4) No Sysout - write in logs from reduce function

#P5) jar tvf jar\_file\_name - to list class files archived for a jar

#P6) space in your folder name which is treated as illegal character - throws error

#P7) Every time you disable billing and enable billing, you need to start VMs manually.