# CS4830 Big Data Laboratory

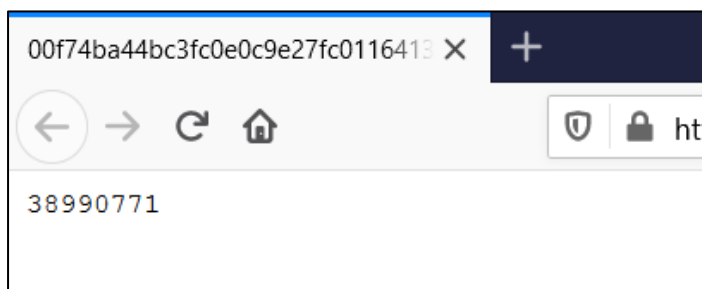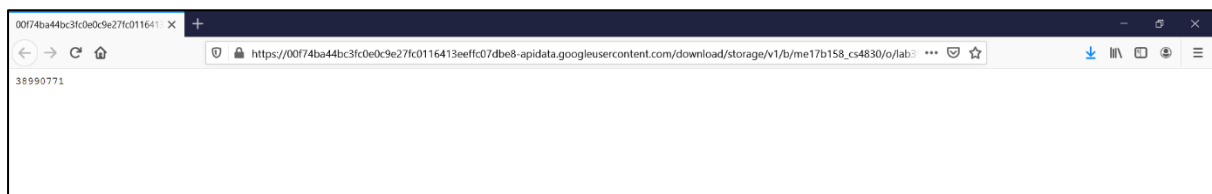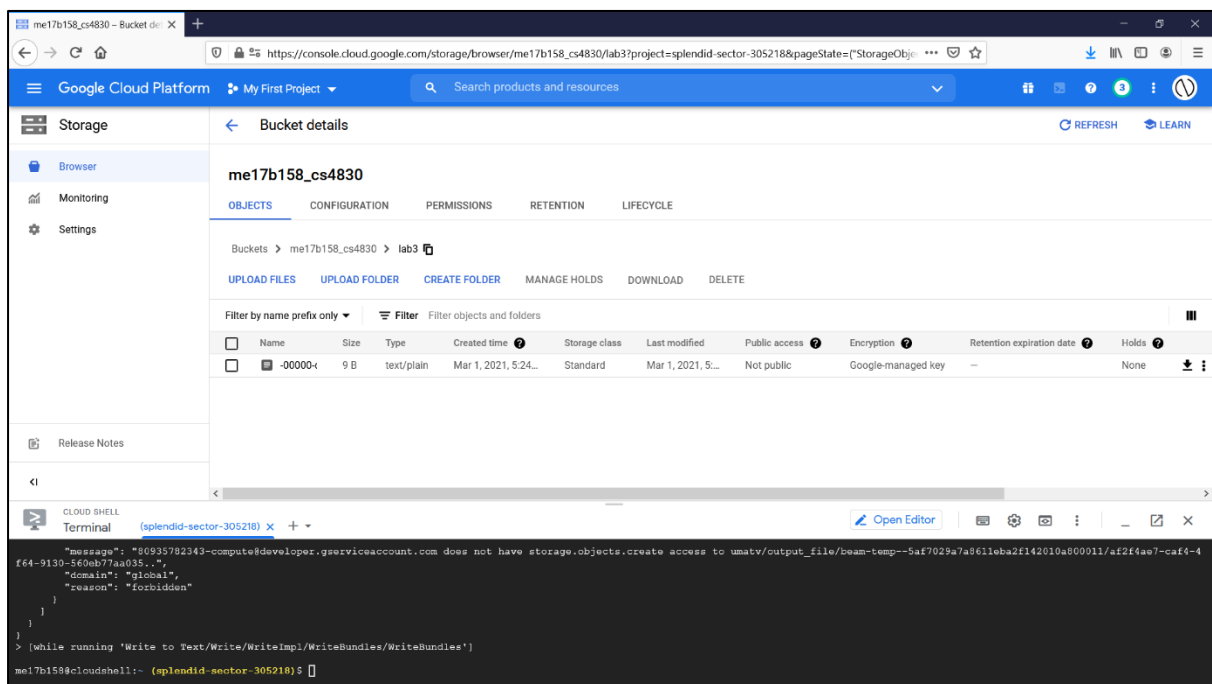## Jan – May 2021

# Lab 3 - Assignment



## ME17B158 - Omkar Nath

**1. Write a Python code to count lines of the file that is placed in the IITMBD bucket (gs://iitmbd/out.txt) using Dataflow and provide the screenshot of the file that is generated in your bucket. [2]**

Attached Code: **"line_count.py"**

Generated File: **"Line_Count.txt"**
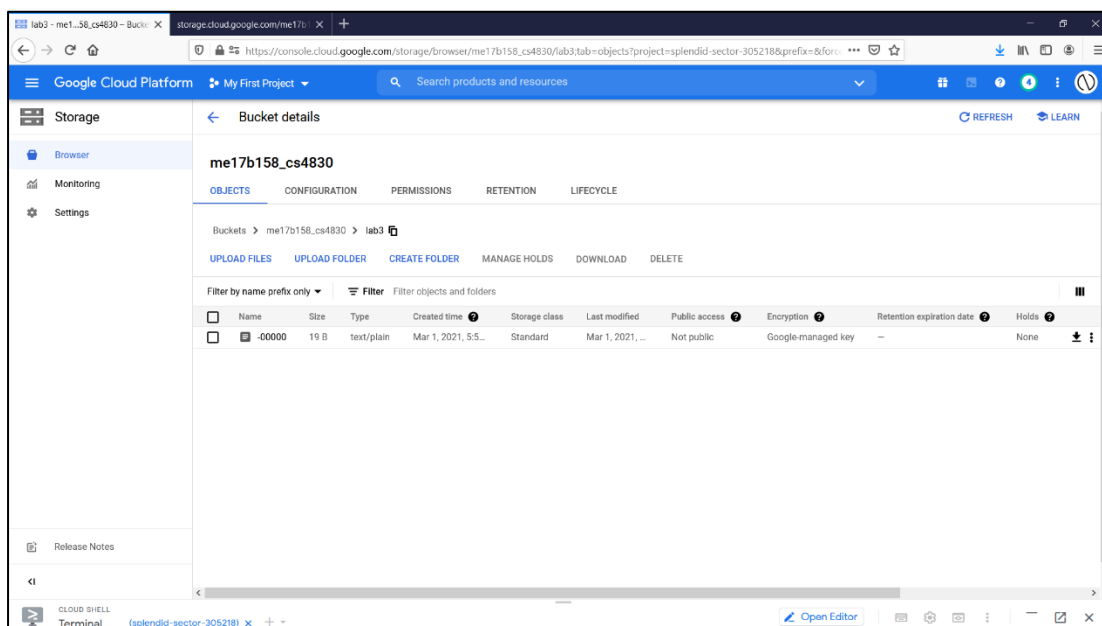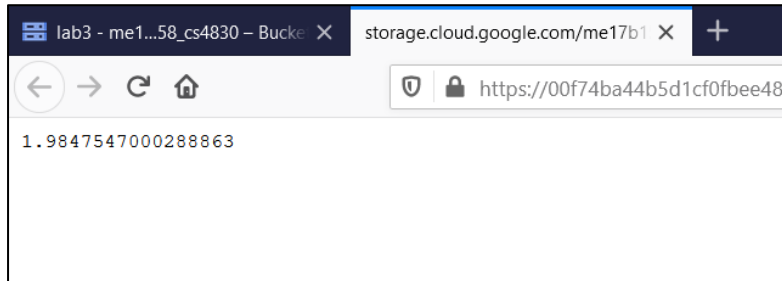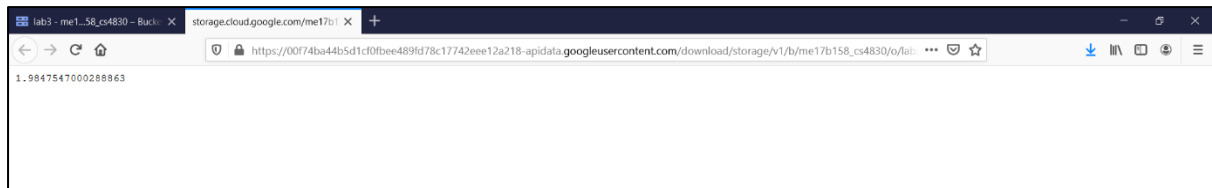
Screenshots of File:

Python Code:

```
import apache_beam as beam
from apache_beam.io import ReadFromText
from apache_beam.io import WriteToText
from apache_beam.options.pipeline_options import PipelineOptions
from apache_beam.options.pipeline_options import GoogleCloudOptions
from apache_beam.options.pipeline_options import StandardOptions
options = PipelineOptions()
google_cloud_options = options .view_as(GoogleCloudOptions)
google_cloud_options.project = 'splendid-sector-305218'
google_cloud_options.region = "us-central1"
google_cloud_options.job_name = 'lab3q1'
google_cloud_options.temp_location = "gs://me17b158_cs4830/tmp"
options.view_as(StandardOptions).runner = 'DataflowRunner'
p = beam.Pipeline( options = options )
lines = p | 'Read' >> beam.io.ReadFromText( 'gs://iitmbd/out.txt' ) |'counting lines' >>
beam.combiners.Count.Globally(sum) | 'Write' >>
beam.io.WriteToText('gs://me17b158_cs4830/lab3/')
result = p.run()
```

**2. Write a Python code to get the average number of words in a line of the file that is placed in the IITMBD bucket (gs://iitmbd/out.txt) using Dataflow and provide the screenshot of the file that is generated in your bucket. [4]**

Attached Code: **"avg_words.py"**

Generated File: **"Average_Words.txt"**

Screenshots of File:

Python Code:

```
import re
import apache_beam as beam
from apache_beam.io import ReadFromText
from apache_beam.io import WriteToText
from apache_beam.options.pipeline_options import PipelineOptions
from apache_beam.options.pipeline_options import GoogleCloudOptions
from apache_beam.options.pipeline_options import StandardOptions
options = PipelineOptions()
google_cloud_options = options .view_as(GoogleCloudOptions)
google_cloud_options.project = 'splendid-sector-305218'
google_cloud_options.region = "us-central1"
google_cloud_options.job_name = 'lab3q2'
google_cloud_options.temp_location = "gs://me17b158_cs4830/tmp"
options.view_as(StandardOptions).runner = "DataflowRunner"
with beam.Pipeline(options=options) as p:
    avgwords = p | 'Read' >> ReadFromText ('gs://iitmbd/out.txt') | 'Counting words per line' >>
beam.Map (lambda x: len(re.split('[\s,;!]+',x))) | 'Taking mean' >>
beam.CombineGlobally((beam.transforms.combiners.MeanCombineFn()) | 'Write to Text' >>
WriteToText('gs://me17b158_cs4830/lab3/')
```

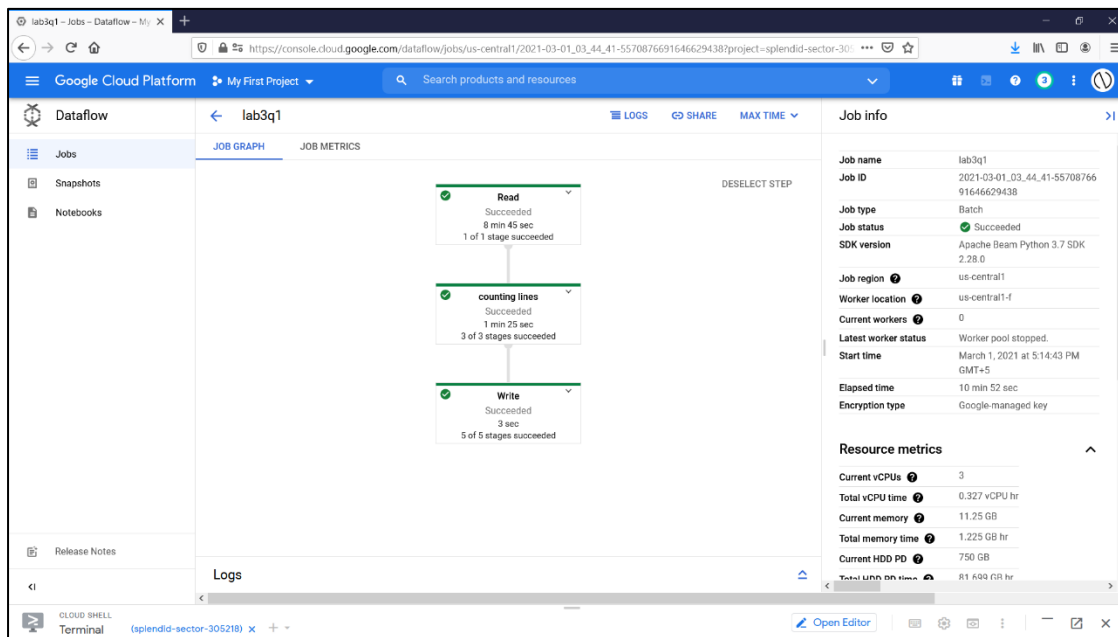

**3. Provide the screenshot for the execution graph created by Dataflow in the background for the pipeline object created for the questions 1 and 2. [2]**

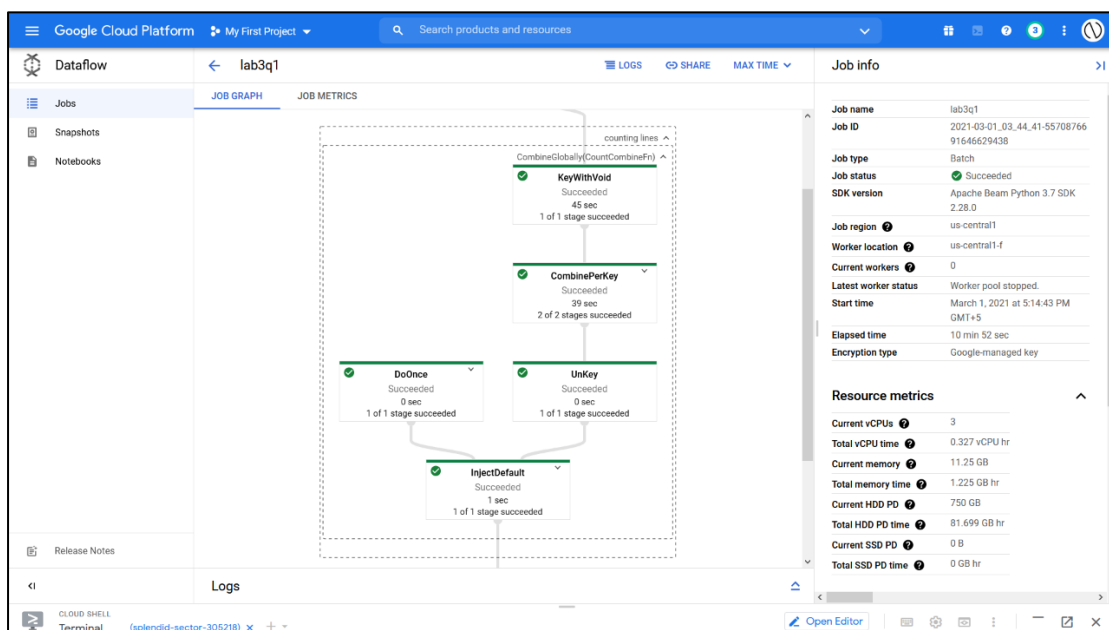

<u>Execution Graph for Question 1:</u>

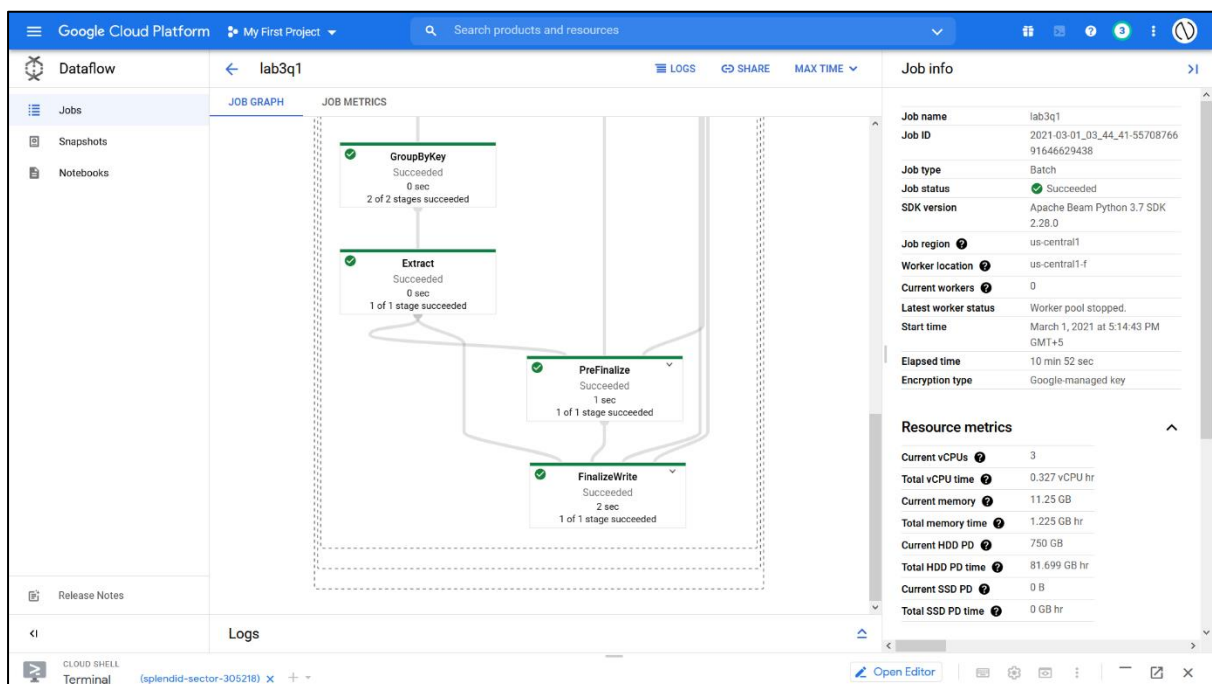

Dataflow:

| lab3q1 | Batch | Mar 1, 2021, 5:25:34 PM | 10 min 51 sec | Mar 1, 2021, 5:14:43 PM | Succeeded | 2.28.0 | 2021-03-01_03_44_41-5570876691646629438 | us-central1 |
|---|---|---|---|---|---|---|---|---|

Overview:


Counting Lines - Expanded:
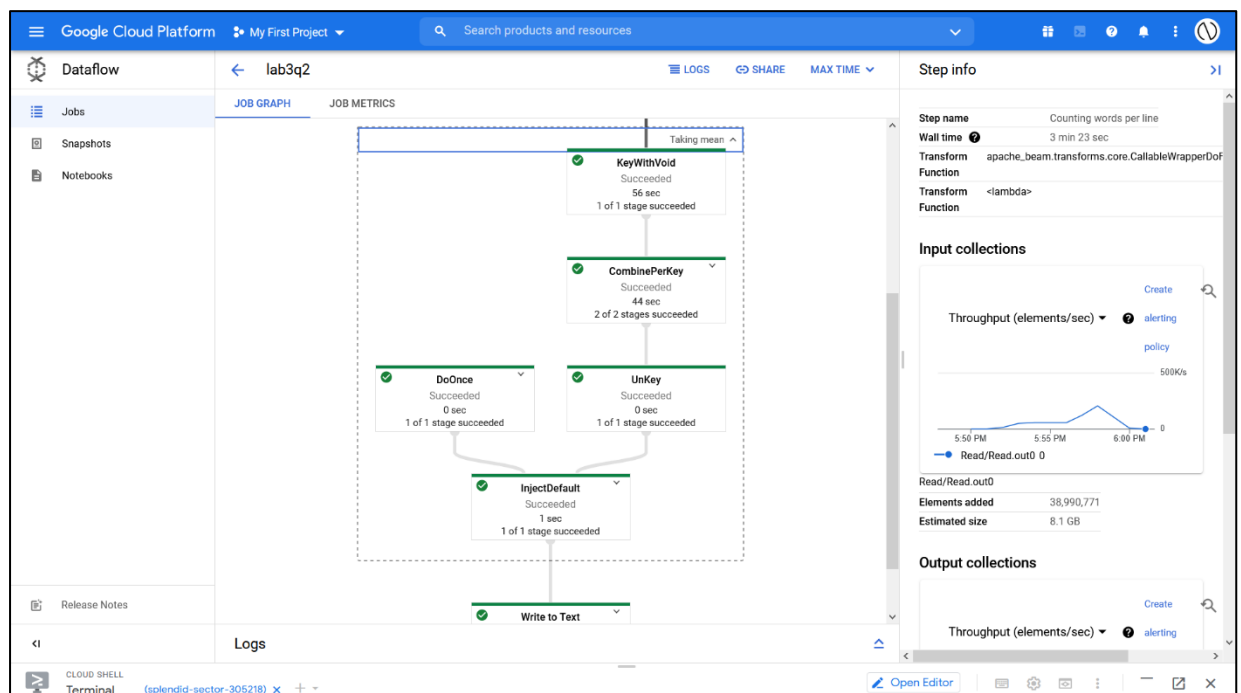
Write – Expanded:

## Execution Graph for Question 2:

Dataflow:



Overview:



Taking Mean:

Write to text:

**4. Explain the pipeline used in the first two questions. What issues did you face while trying to make the code work for the first two questions and how did you resolve them? [2]**

Usage of Google Cloud:
The apache beam pipeline is used in the Google cloud storage for running applications that involve the processing of large amounts of data. The broad steps for this are as follows:

1. Locate the data to be processed
2. Create a python code to perform the required function.
3. Use Google Cloud functions for more efficient processing.
4. Launch the process on Google Cloud platform.
5. Output is written to the specified location.

Pipeline for first question:

To count the number of lines:
1. Write the python script with the tasks to be performed.
    a. Read the text from out.txt given to us
    b. Count the global number of lines in the file
    c. Write the result to the output file.
2. Launch a virtual environment
3. Run the python script.

Pipeline for second question:

To count average number of words per line:
1. Write the python script with the tasks to be performed.
    a. Read the text from out.txt given to us
    b. Count the number of words per line
    c. Take the global average over all the lines
    d. Write the result to the output file.
2. Launch a virtual environment
3. Run the python script.

Challenges faced:

1. Understanding how to use the Google Cloud Platform, as well as the purpose of various functions.
   This was resolved by going through the various associated documentation

2. Various permissions had to be given, and additional add on packages had to be added.
   These were added as per the given prompts.

3. Specifically, for the addition of the Data Analytics Package, it took about ten minutes to add the package.

Had to wait that long for the change to be reflected.

4. Finding appropriate Cloud Functions, especially for use in the second questions. Going through the various relevant documentation was able to resolve this problem

5. The dataflow was not being triggered for Q5 by having the file in the same bin. A second bin had to be created in which the file had to be moved into to trigger the dataflow.

**5. [Bonus] Trigger a dataflow using GCF for any one of the first two questions. [2]**

Attached Code: **"main.py"**

Attached Requirements: **"requirements.txt"**

Generated File: **"Count_Lines.txt"**

Python Code:

```
def dataflow_count_lines(data, context):
    from uuid import uuid4
    import apache_beam as beam
    from apache_beam.io import ReadFromText
    from apache_beam.io import WriteToText
    from apache_beam.options.pipeline_options import PipelineOptions
    from apache_beam.options.pipeline_options import GoogleCloudOptions
    from apache_beam.options.pipeline_options import StandardOptions
    file_path = f"gs://{data['bucket']}/{data['name']}"
    unique_id = f"{data['name'].split('.')[0]}-{uuid4()}"
    output_path = f"gs://me17b158_cs4830/lab3/"
    options = PipelineOptions()
    google_cloud_options = options.view_as(GoogleCloudOptions)
    google_cloud_options.project = "splendid-sector-305218"
    google_cloud_options.job_name = f"{unique_id}"
    google_cloud_options.temp_location="gs://me17b158_cs4830/tmp/"
    options.view_as(StandardOptions).runner = "DataflowRunner"
    google_cloud_options.region="us-central1"
    with beam.Pipeline(options=google_cloud_options) as p:
        lines = p | 'Read' >> ReadFromText(file_path)
        counts = lines | 'Count elements' >> beam.combiners.Count.Globally()
        output = counts
        output | 'Write' >> WriteToText(output_path)
```

Requirements:

apache-beam[gcp]

Triggering the Google Cloud Dataflow:

*gcloud functions deploy dataflow_count_lines --runtime python37 --timeout 540 --trigger-resource gs://me17b158_cs4830_2 --trigger-event google.storage.object.finalize*

Copying the file (hence triggering dataflow):
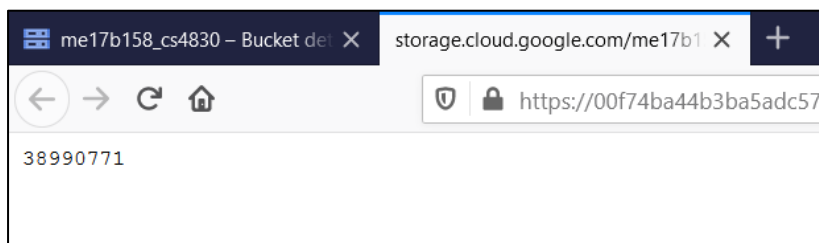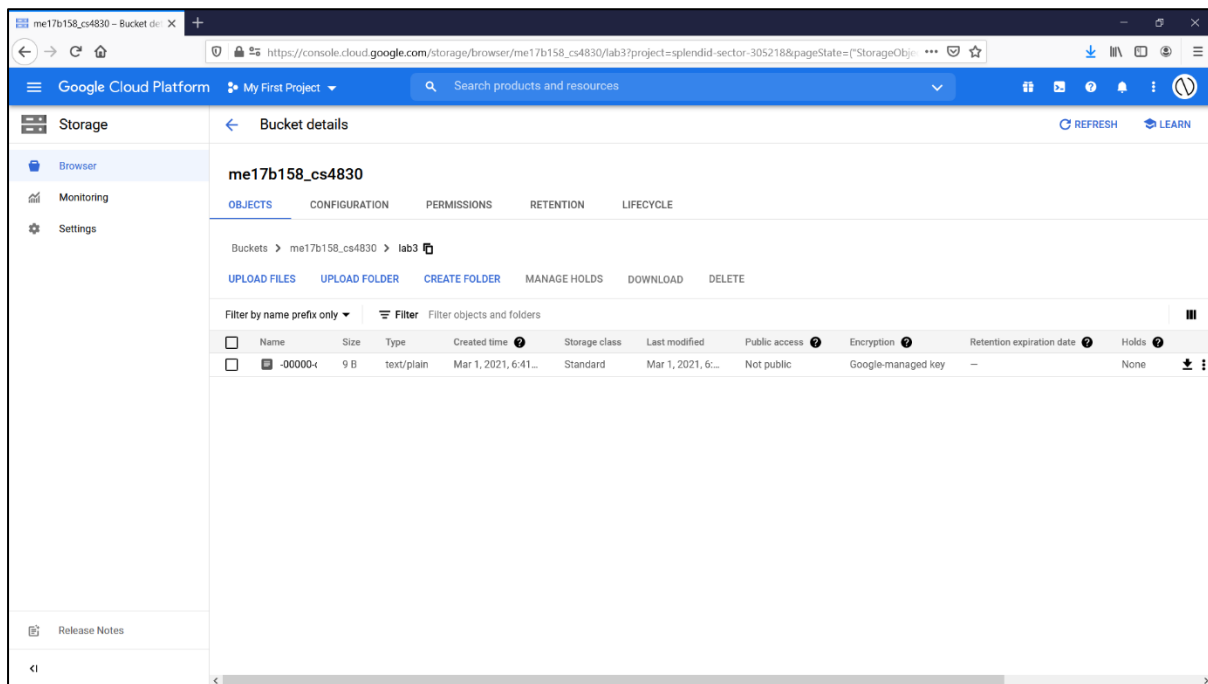
*gsutil cp gs://iitmbd/out.txt gs://me17b158_cs4830_2/*

Screenshots of created Process.

Screenshots of File:







The results obtained is the same as in Q1.

**Create a PDF file containing answers to the above questions. Zip it along with the output files (for the dataflow task) and your Python files. Then, submit this zip file on moodle.**