# Lab 3 - Assignment

### MM16B023[1]

[a]*Indian Institute of Technology Madras*

**Keyword:** VM, Cloud Functions, DataFlow, pipeline, Pcollections
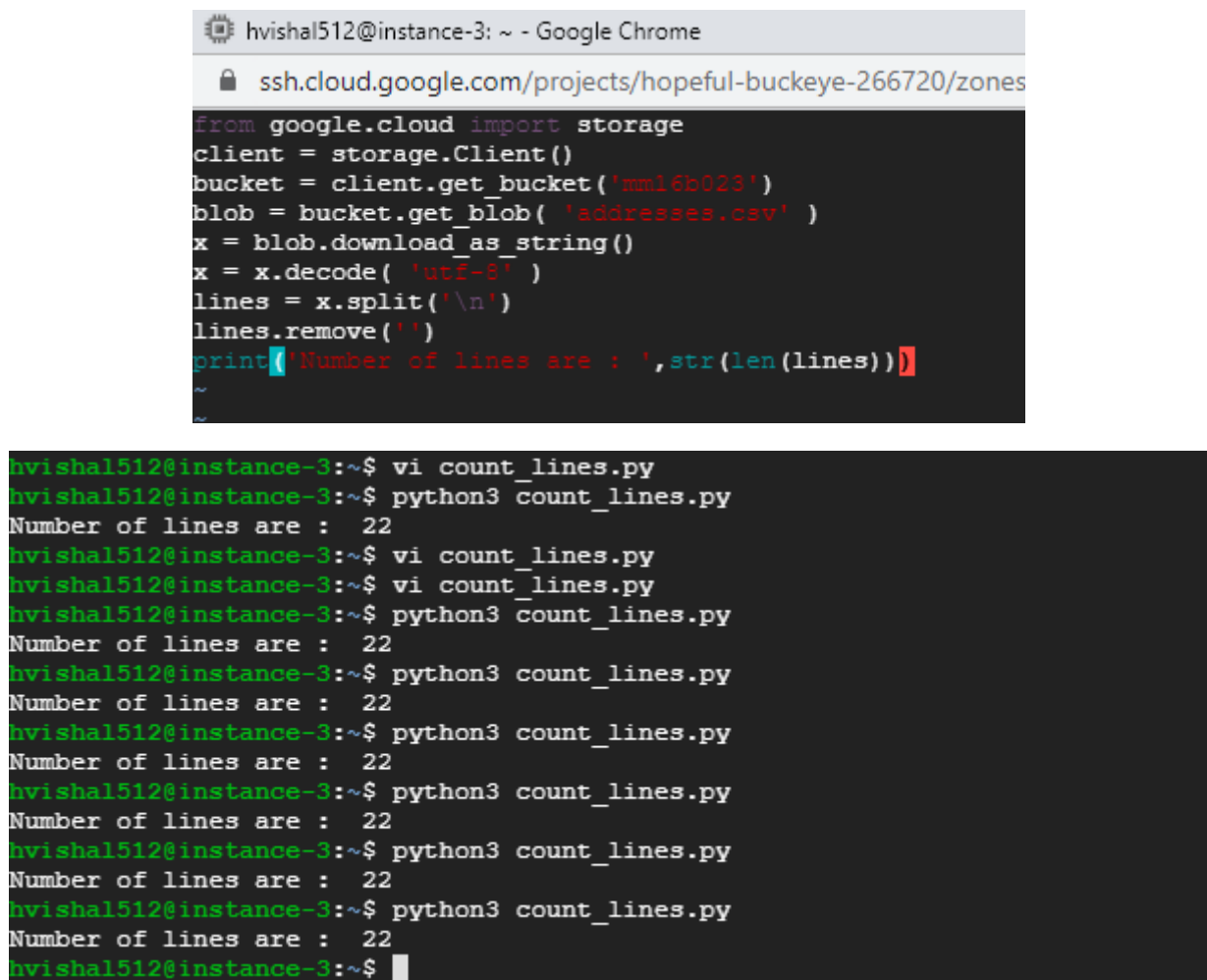
**Abstract:** This paper presents the solutions to first assignment of the Big Data Laboratory course (CS4830) at *IIT Madras.* All the notations used are as according with the textbook Mining of massive data sets by Anand Rajaraman

## Problem 1

Provide screenshot for logs of all the 3 tasks containing the required result - line count using VM, Cloud Functions and DataFlow, along with the python file for each task.
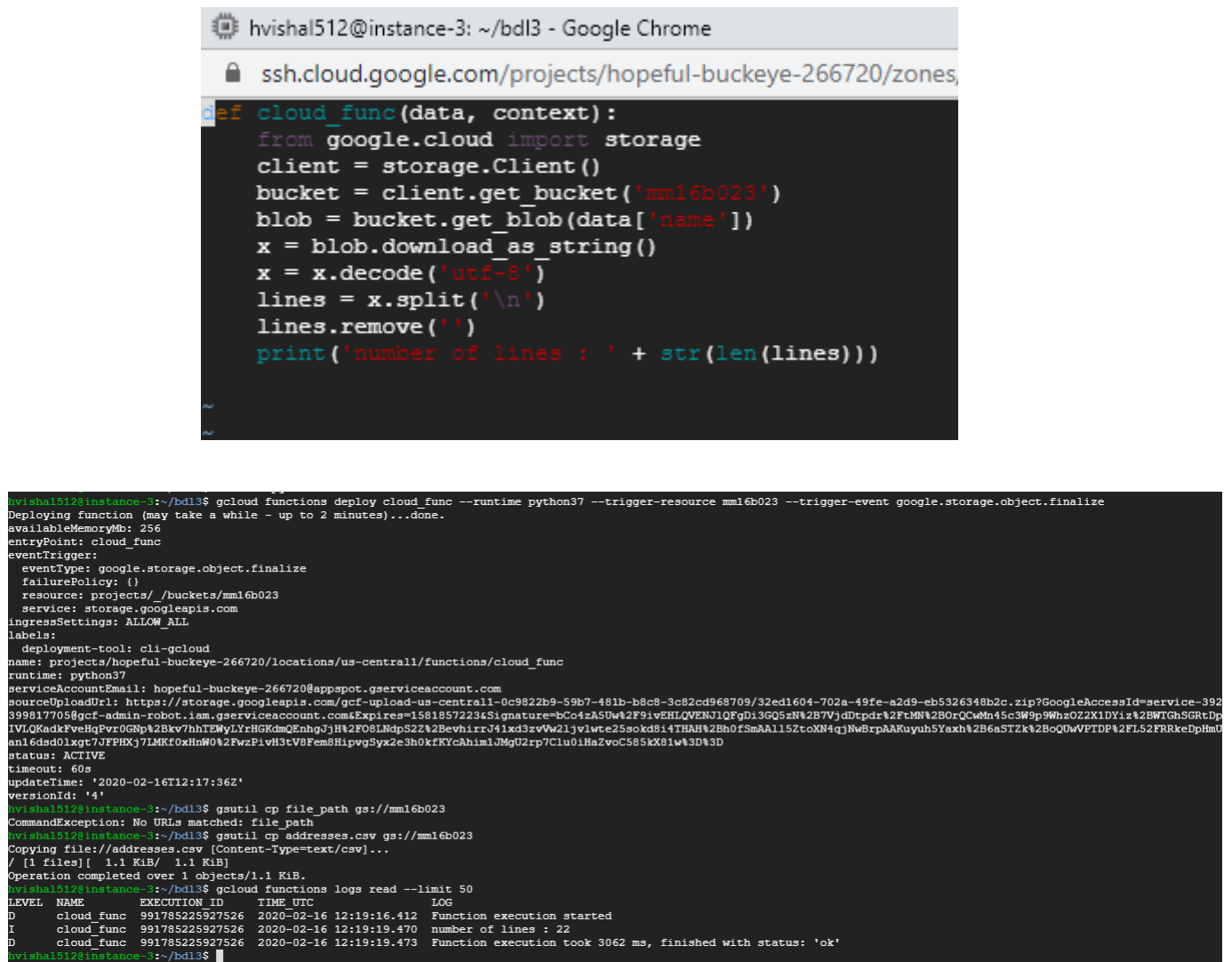
**Solution:**

**a) Virtual Machine**



**Fig1:** Counting number of lines using Virtual Machine
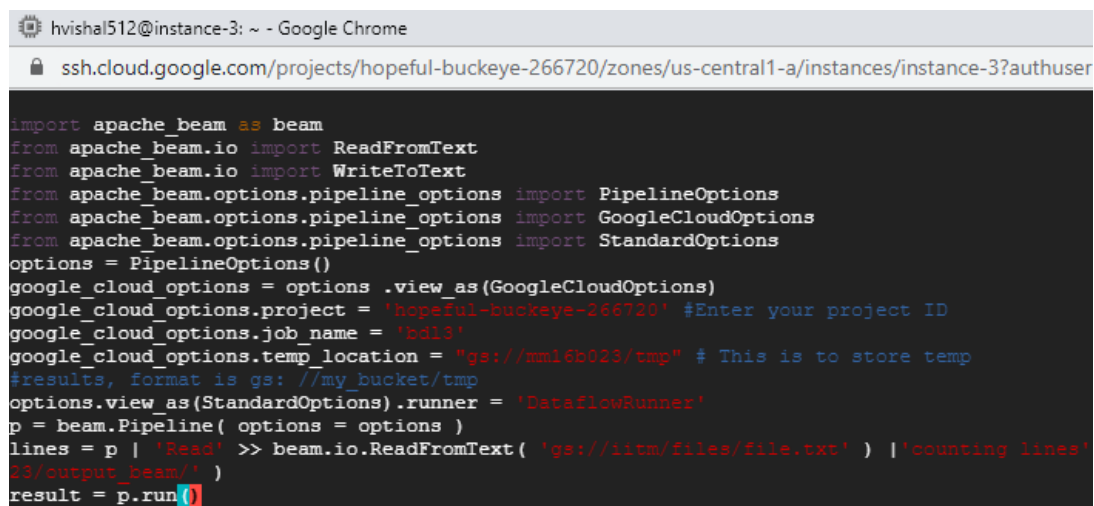
## b) Cloud Function



```python
def cloud_func(data, context):
    from google.cloud import storage
    client = storage.Client()
    bucket = client.get_bucket('mm16b023')
    blob = bucket.get_blob(data['name'])
    x = blob.download_as_string()
    x = x.decode('utf-8')
    lines = x.split('\n')
    lines.remove('')
    print('number of lines : ' + str(len(lines)))
```



**Fig 2:** Counting the number of lines using Cloud function

Using the above two methods, <u>number of lines in the file 'addresses.csv' = 22</u>

## c) Data flow



```python
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 = 'hopeful-buckeye-266720' #Enter your project ID
google_cloud_options.job_name = 'bdl3'
google_cloud_options.temp_location = "gs://mm16b023/tmp" # This is to store temp
#results, format is gs: //my_bucket/tmp
options.view_as(StandardOptions).runner = 'DataflowRunner'
p = beam.Pipeline( options = options )
lines = p | 'Read' >> beam.io.ReadFromText( 'gs://iitm/files/file.txt' ) |'counting lines'
23/output_beam/' )
result = p.run()
```

**Note:** Some portion of the image has been cropped out for aesthetics. Please refer to the code attached.



Buckets / mm16b023 / output_beam / -00000-of-00001

| | |
|---|---|
| Access | Not public |
| Type | text/plain |
| Size | 9 B |
| Created | February 16, 2020 at 6:14:24 PM UTC+5:30 |
| Last modified | February 16, 2020 at 6:14:24 PM UTC+5:30 |
| URI | gs://mm16b023/output_beam/-00000-of-00001 |
| Link URL | https://storage.cloud.google.com/mm16b023/output_beam/-0000 of-00001 |

00e9e64bacc921d6ecedb522dfb0fdf20c48b8f82a67063c94-apidata.googleusercontent.com/download/storage/v1/b/mm16b023/o/output_beam

25974026

## Logs



**Fig3 :** Counting the number of lines using DataFlow

*Number of lines in the data set = 25974026*

## Problem2

Provide the screenshot for the execution graph created by Dataflow in the background for the pipeline object created in task 3

**Solution:**

## Problem 3

Explain the pipeline used in task 3. What issues did you face while trying to make the code work for task 3 and how did you resolve them?
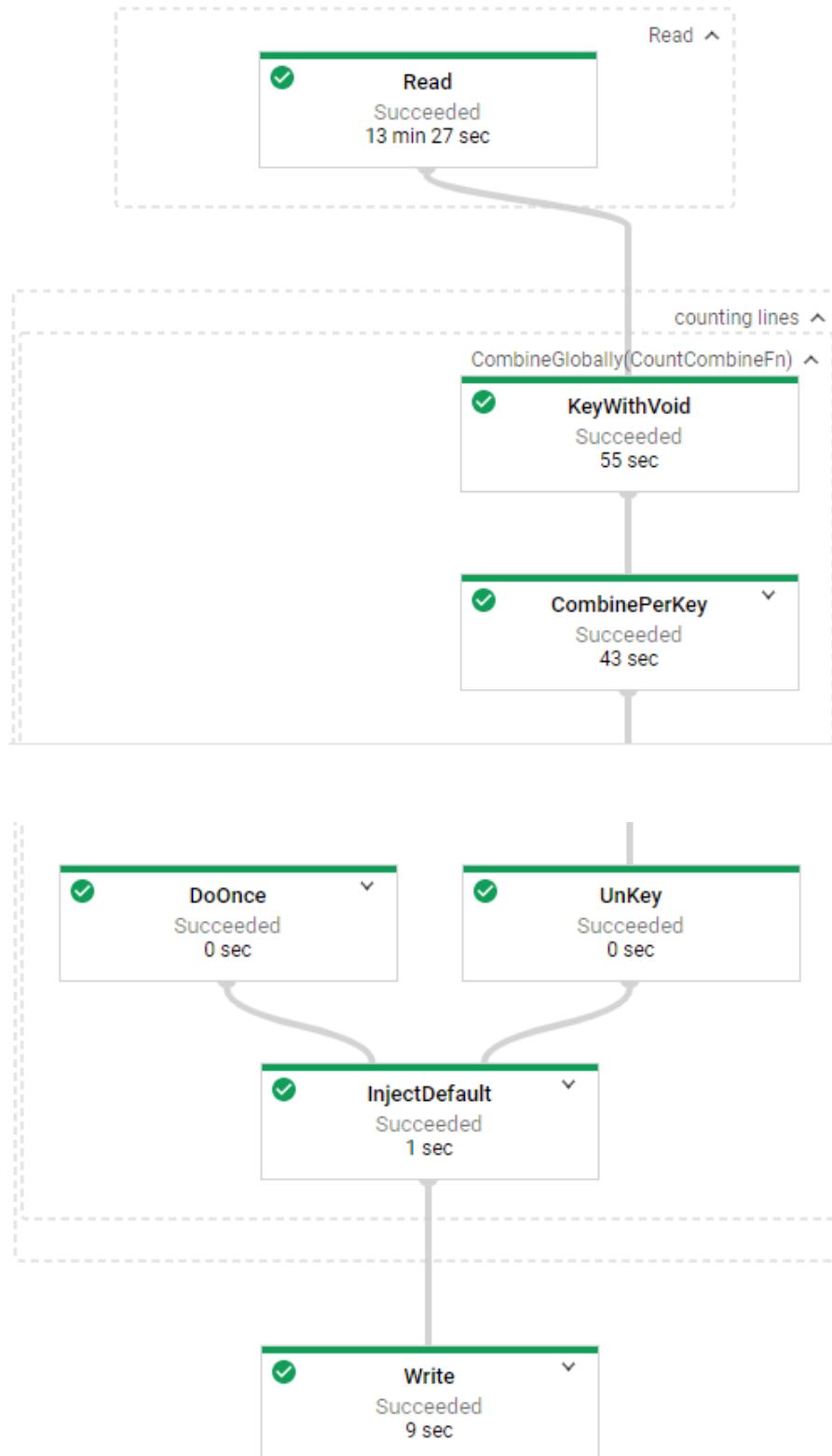
**Solution:** The pipeline used has the following transforms embedded in it

- **ReadFromText** - Navigates through the file reading each line of it

- **CombineGlobally** - Stores each line as a string into a bigger dataset called <u>PCollections</u>

- **WriteToText** - Writes the output (no. of lines in this case) to the specified directory

**Issues faced**

- There was an exception thrown in the worker code when multiple lines of the code were joined by a vertical line. Surprisingly, no error came when it was written as a big single line.

- Not an issue per se, but picking the appropriate transforms and designing the right pipeline structure took some time as there are many functions and hence combinations possible.

## Problem4

PCollections can handle unbounded data. What is meant by unbounded data and how do you think Pcollections can handle it? (Hint: Think on the line of triggers)

**Solution:**

- Data streamed from a constantly updating source is known as unbounded data. Since the data source continuously adds new elements, it has unlimited size in some sense and hence the name unbounded data.

- PCollections can handle it by modifying triggers in such a way that it maps to each collection in a streaming pipeline. This can be done using Apache Beam SDK which uses windowing (segregate unbounded data into windows) to operate the data on a combination of event time (indicated by time stamps), time of processing and the number of data elements. This modified trigger mechanism outputs whenever the watermark passes the edge of the window.