# Perform Foundational Data, ML, and AI Tasks in Google Cloud Challenge Lab
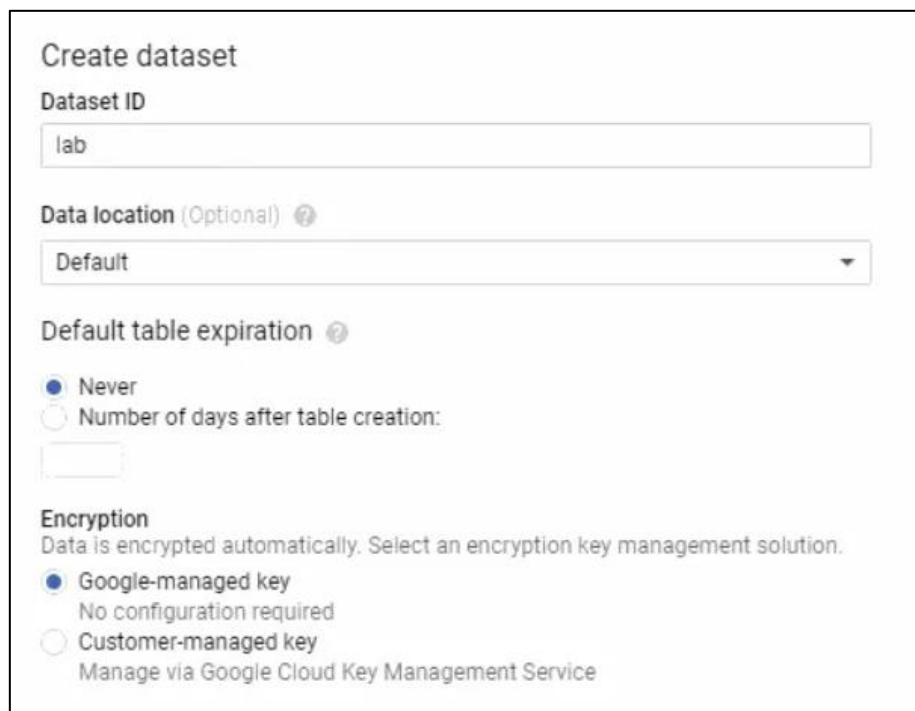
**Author:** Vedant Kakde | **GitHub Profile:** github.com/vedant-kakde | **LinkedIn Profile:** linkedin.com/in/vedant-kakde/

## Task 1: Run a simple Dataflow job

In this task, you have to transfer the data in a CSV file to BigQuery using Dataflow via Pub/Sub. First of all, you need to create a BigQuery dataset called `lab` and a Cloud Storage bucket called with your project ID.

### 1.1 Created a BigQuery dataset called *Lab*

1. In the Cloud Console, click on **Navigation Menu** > **BigQuery**.
2. Select your project in the left pane.
3. Click **CREATE DATASET**.
4. Enter `lab` in the Dataset ID, then click **Create dataset**.



5. Run `gsutil cp gs://cloud-training/gsp323/lab.schema .` in the Cloud Shell to download the schema file.
6. View the schema by running `cat lab.schema`.

```
[
  {"type":"STRING","name":"guid"},
  {"type":"BOOLEAN","name":"isActive"},
  {"type":"STRING","name":"firstname"},
  {"type":"STRING","name":"surname"},
  {"type":"STRING","name":"company"},
  {"type":"STRING","name":"email"},
  {"type":"STRING","name":"phone"},
  {"type":"STRING","name":"address"},
  {"type":"STRING","name":"about"},
  {"type":"TIMESTAMP","name":"registered"},
  {"type":"FLOAT","name":"latitude"},
  {"type":"FLOAT","name":"longitude"}
]
```

**Author:** Vedant Kakde | **GitHub Profile:** github.com/vedant-kakde | **LinkedIn Profile:** linkedin.com/in/vedant-kakde/

Copy this schema into notepad. We will use it in further step.

7. Go back to the Cloud Console, select the new dataset **lab** and click **Create Table**.
8. In the Create table dialog, select **Google Cloud Storage** from the dropdown in the Source section.
9. Copy `gs://cloud-training/gsp323/lab.csv` to **Select file from GCS bucket**.
10. Enter `customers` to "Table name" in the Destination section.
11. Enable **Edit as text** and copy the JSON data from the `lab.schema` file to the textarea in the Schema section.
12. Click **Create table**.

```
Create table

Source
Create table from:          Select file from GCS bucket: ⓘ                      File format:
[Google Cloud Storage ▼]    [✔] gs://cloud-training/gsp323/lab.csv  [Browse]    [CSV ▼]

[ ] Source Data Partitioning


Destination
● Search for a project      ○ Enter a project name

Project name                       Dataset name              Table type ⓘ
[qwiklabs-gcp-00-e66de107d4d4 ▼]   [lab              ▼]      [Native table ▼]

Table name
[customers|                                              ]


Schema

Auto detect
[ ] Schema and input parameters

   ⦾ Edit as text

    1  [
    2    {"type":"STRING","name":"guid"},
    3    {"type":"BOOLEAN","name":"isActive"},
    4    {"type":"STRING","name":"firstname"},
    5    {"type":"STRING","name":"surname"},
    6    {"type":"STRING","name":"company"},
    7    {"type":"STRING","name":"email"},
    8    {"type":"STRING","name":"phone"},
    9    {"type":"STRING","name":"address"},
   10    {"type":"STRING","name":"about"},
   11    {"type":"TIMESTAMP","name":"registered"},
   12    {"type":"FLOAT","name":"latitude"},
   13    {"type":"FLOAT","name":"longitude"}
   14  ]


Partition and cluster settings
Partitioning: ⓘ
[No partitioning        ▼]
```
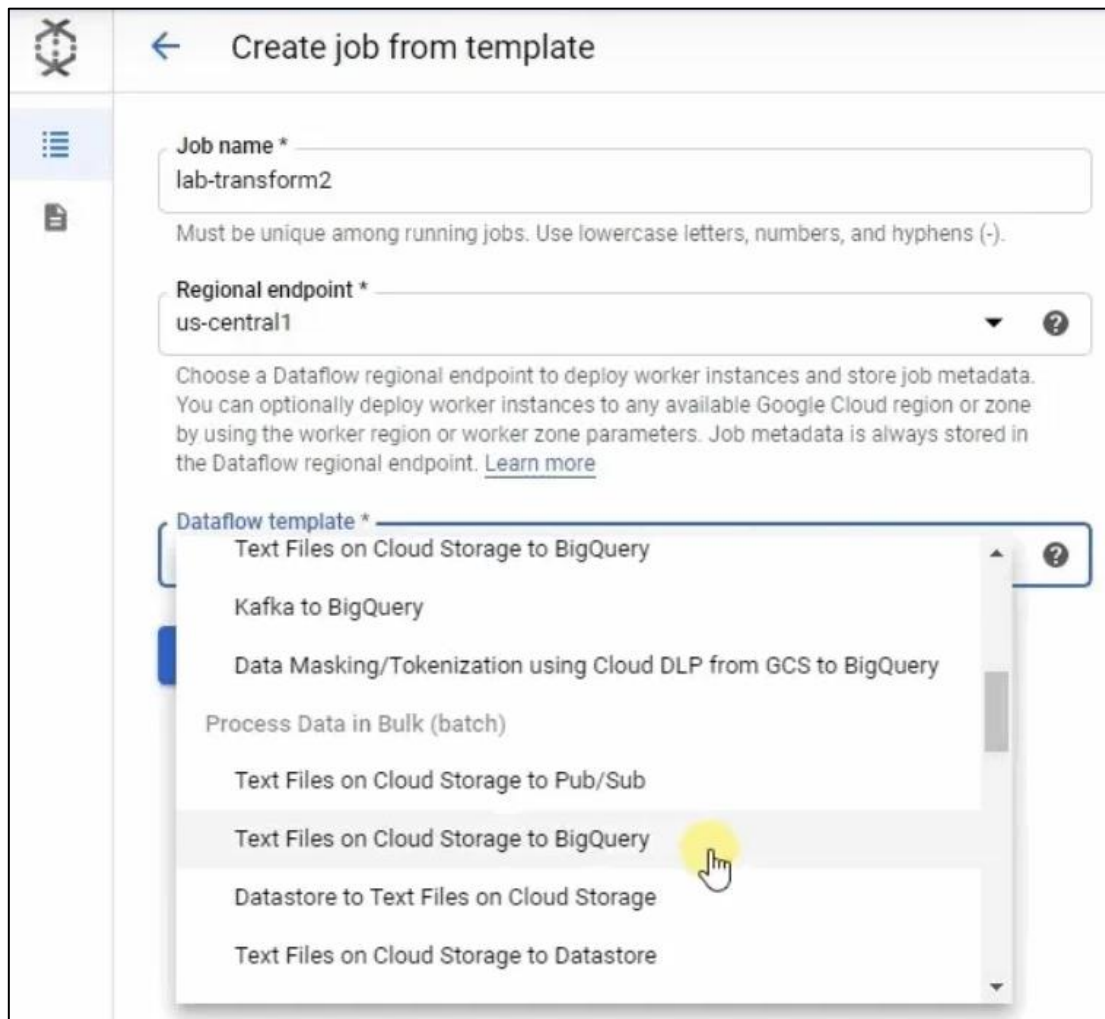
### 1.2 Create a Cloud Storage bucket

1. In the Cloud Console, click on **Navigation Menu** > **Storage**.
2. Click **CREATE BUCKET**.
3. Copy your GCP Project ID to Name your bucket.

**Author:** Vedant Kakde | **GitHub Profile:** github.com/vedant-kakde | **LinkedIn Profile:** linkedin.com/in/vedant-kakde/

4. Click **CREATE**.

### 1.3 Create a Dataflow job

1. In the Cloud Console, click on **Navigation Menu** > **Dataflow**.
2. Click **CREATE JOB FROM TEMPLATE**.
3. In Create job from template, give an arbitrary job name.
4. From the dropdown under Dataflow template, select **Text Files on Cloud Storage Pub/Sub** under "Process Data in Bulk (batch)". (**DO NOT** select the item under "Process Data Continuously (stream)").



5. Under the Required parameters, enter the following values:

| Field | Value |
|---|---|
| JavaScript UDF path in Cloud Storage | gs://cloud-training/gsp323/lab.js |
| JSON path | gs://cloud-training/gsp323/lab.schema |
| JavaScript UDF name | transform |
| BigQuery output table | YOUR_PROJECT:lab.customers |

| Field | Value |
|---|---|
| Cloud Storage input path | `gs://cloud-training/gsp323/lab.csv` |
| Temporary BigQuery directory | `gs://YOUR_PROJECT/bigquery_temp` |
| Temporary location | `gs://YOUR_PROJECT/temp` |

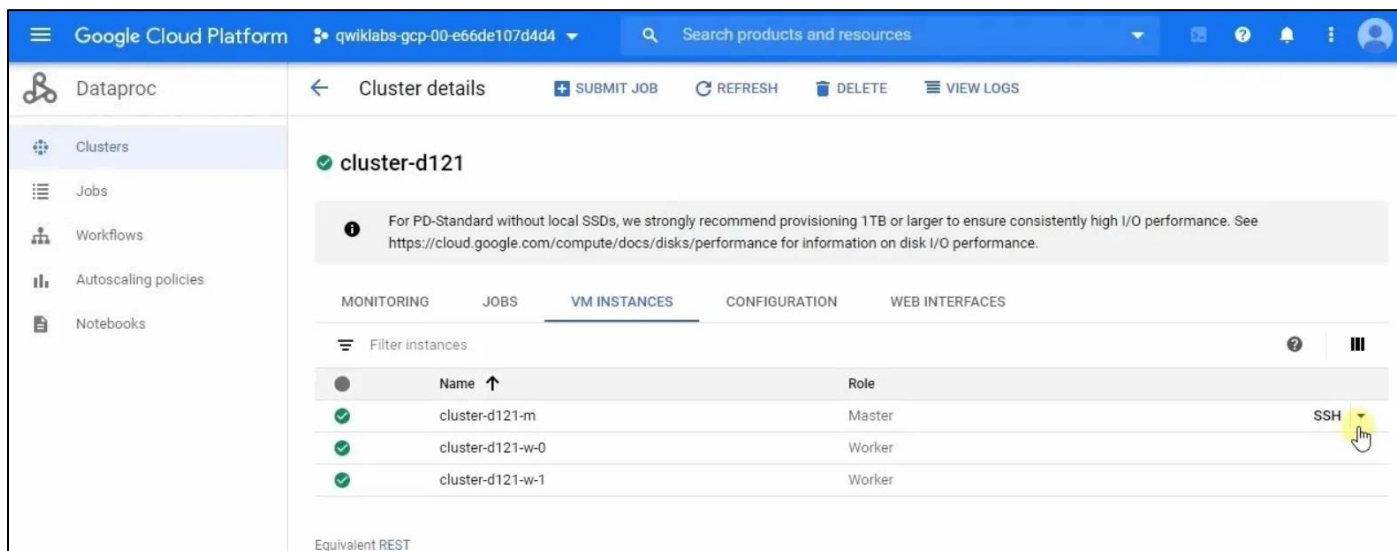**Replace** `YOUR_PROJECT` with your project ID.



6. Click **RUN JOB**.

# Task 2: Run a simple Dataproc job

### Create a Dataproc cluster

1. In the Cloud Console, click on **Navigation Menu** > **Dataproc** > **Clusters**.
2. Click **CREATE CLUSTER**.
3. Make sure the cluster is going to create in the region **us-central1**.
4. Click **Create**.

5. After the cluster has been created, click the **SSH** button in the row of the master instance.



6. In the SSH console, run the following command:

```
hdfs dfs -cp gs://cloud-training/gsp323/data.txt /data.txt
```

7. Close the SSH window and go back to the Cloud Console.
8. Click **SUBMIT JOB** on the cluster details page.
9. Select **Spark** from the dropdown of "Job type".
10. Copy `org.apache.spark.examples.SparkPageRank` to "Main class or jar".
11. Copy `file:///usr/lib/spark/examples/jars/spark-examples.jar` to "Jar files".
12. Enter `/data.txt` to "Arguments".



**Author:** Vedant Kakde | **GitHub Profile:** github.com/vedant-kakde | **LinkedIn Profile:** linkedin.com/in/vedant-kakde/

13.    Click **CREATE**.

### Import runs.csv to Dataprep

1. In the Cloud Console, click on **Navigation menu** > **Dataprep**.
2. After entering the home page of Cloud Dataprep, click the **Import Data** button.



3. In the Import Data page, select **GCS** in the left pane.
4. Click on the pencil icon under Choose a file or folder.
5. Copy `gs://cloud-training/gsp323/runs.csv` to the textbox, and click the **Go** button next to it.



6. After showing the preview of runs.csv in the right pane, click on the **Import & Wrangle** button.

### Transform data in Dataprep

1. After launching the Dataperop Transformer, scroll right to the end and select **column10**.
2. In the Details pane, click **FAILURE** under Unique Values to show the context menu.
3. Select **Delete rows with selected values** to Remove all rows with the state of "FAILURE".

4. Click the downward arrow next to **column9**, choose **Filter rows** > **On column value** > **Contains**.
5. In the Filter rows pane, enter the regex pattern `/(^0$|^0\.0$)/` to "Pattern to match".
6. Select **Delete matching rows** under the Action section, then click the **Add** button.



7. Rename the columns to be:
   - runid
   - userid
   - labid
   - lab_title

- start
- end
- time
- score
- state

8. Confirm the recipe. It should like the screenshot below.



9. Click **Run Job**.

# Task 4: AI

```
gcloud iam service-accounts create my-natlang-sa \
  --display-name "my natural language service account"

gcloud iam service-accounts keys create ~/key.json \
  --iam-account my-natlang-sa@${GOOGLE_CLOUD_PROJECT}.iam.gserviceaccount.com

export GOOGLE_APPLICATION_CREDENTIALS="/home/$USER/key.json"

gcloud auth activate-service-account my-natlang-
sa@${GOOGLE_CLOUD_PROJECT}.iam.gserviceaccount.com --key-
file=$GOOGLE_APPLICATION_CREDENTIALS

gcloud ml language analyze-entities --content="Old Norse texts portray Odin as
one-eyed and long-bearded, frequently wielding a spear named Gungnir and wearing a
cloak and a broad hat." > result.json

gcloud auth login
(Copy the token from the link provided)

gsutil cp result.json gs://YOUR_PROJECT-marking/task4-cnl.result
```

*Use Google Cloud Speech API to analyze the audio file*

1. In the Cloud Console, click on **Navigation menu** > **APIs & Services** > **Credentials**.
2. In the Credentials page, click on **+ CREATE CREDENTIALS** > **API key**.
3. Copy the API key to the clipboard, then click **RESTRICT KEY**.
4. Open the Cloud Shell, store the API key as an environment variable by running the following command:

```
export API_KEY=<YOUR-API-KEY>
```

5. **Replace** <YOUR-API-KEY> with the copied key value.

**Author:** Vedant Kakde | **GitHub Profile:** github.com/vedant-kakde | **LinkedIn Profile:** linkedin.com/in/vedant-kakde/

```
nano
request.json
```

```json
{
  "config": {
      "encoding":"FLAC",
      "languageCode": "en-US"
  },
  "audio": {
      "uri":"gs://cloud-training/gsp323/task4.flac"
  }
}
```

```
curl -s -X POST -H "Content-Type: application/json" --data-binary @request.json \
"https://speech.googleapis.com/v1/speech:recognize?key=${API_KEY}" > result.json

gsutil cp result.json gs://YOUR_PROJECT-marking/task4-gcs.result
```

# TASK 4 - PART 3 - VIDEO INTELLIGENCE:

```
gcloud iam service-accounts create quickstart

gcloud iam service-accounts keys create key.json --iam-account
quickstart@${GOOGLE_CLOUD_PROJECT}.iam.gserviceaccount.com

gcloud auth activate-service-account --key-file key.json

export ACCESS_TOKEN=$(gcloud auth print-access-token)
```

```
nano request.json
```

```json
{
    "inputUri":"gs://spls/gsp154/video/chicago.mp4",
    "features": [
        "TEXT_DETECTION"
    ]
}
```

```
curl -s -H 'Content-Type: application/json' \
    -H "Authorization: Bearer $ACCESS_TOKEN" \
    'https://videointelligence.googleapis.com/v1/videos:annotate' \
    -d @request.json
```

```
curl -s -H 'Content-Type: application/json' -H "Authorization:
Bearer $ACCESS_TOKEN"
'https://videointelligence.googleapis.com/v1/operations/OPERATION_FR
OM_PREVIOUS_REQUEST' > result1.json


gsutil cp result1.json gs://YOUR_PROJECT-marking/task4-gvi.result
```

**Congratulations! You completed this challenge lab.**