

Implementation of GCP Dataflow

NYC AIRBNB



Shriyanshi Shikha

Contact: 765-720-6222

Email: shriyanshikha@gmail.com

Problem Statement:

Create a pipeline in Dataflow that reads data from a csv file, applies transformations, and inserts resulting data into a BigQuery table.

Implementation:

I have implemented the Dataflow pipeline with the help of “Apache Beam SDK for Python” and “Google Cloud Tools”. The steps in the implementation of python application are as below:

- I. Created a GCP account and enabled necessary APIs
- II. Created a new project
- III. Created a Service account to access certain GCP services and also performed the setup of Google Cloud Shell
- IV. Created a service account key and downloaded the JSON file that contains the service account key
- V. Created a Cloud storage bucket to use the input data i.e. AB_NYC_2019
- VI. Created a BigQuery dataset manually under the same GCP project

After setting up the account I used the command prompt to set the virtual environment and once that's done I defined a pipeline with an Apache Beam program and chose the runner as Dataflow to run the pipeline.

I initialized an object class to leverage the Dataflow execution phases in the form of Directed Acyclic Graph. Find the stages of DAG execution as below:

- I. **Read** the input NYC-airbnb CSV file stored in GCS bucket
- II. **Extract** the field and rows from the input file and load into collections
- III. **Transform** the input rows to BigQuery compatible row format. Perform Group By aggregation on Neighbourhood field
- IV. **Load** the data into bigQuery table

```
with beam.Pipeline(options=pipeline_options) as p:
    lines=(p| "ReadFromFile" >> beam.Create([known_args.input])
    | "ParseCSV" >> beam.FlatMap(get_csv_reader)
    | "write to bigquery :" >>
        beam.io.WriteToBigQuery(table="springmltest",dataset="springmltest",project="springmltest",schema=table_schema,create_
        disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED))
```

Figure 1: Dataflow pipeline to read the csv file and write to a bigquery table

```

with beam.Pipeline(options=pipeline_options) as p:
    nb_count=(p|'QueryTableStdSQL' >> beam.io.ReadFromBigQuery(

        query='SELECT neighbourhood, sum(calculated_host_listings_count) as count FROM `
        `springmltest.springmltest.springmltest` group by neighbourhood',
        use_standard_sql=True)
    |"write count to bigquery :" >>
        beam.io.WriteToBigQuery(table="neighbourhoodcount",dataset="neighbourhoodcount",project="springmltest",schema=table_cou
        nt_schema,create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED))

```

Figure 2: Dataflow pipeline to read the original data from a bigquery table and run query to calculate listings count by neighbourhood

Output:

Below is a set of screenshots to leverage a successful ETL processing of NYC Airbnb dataset to BigQuery. Please find them below:

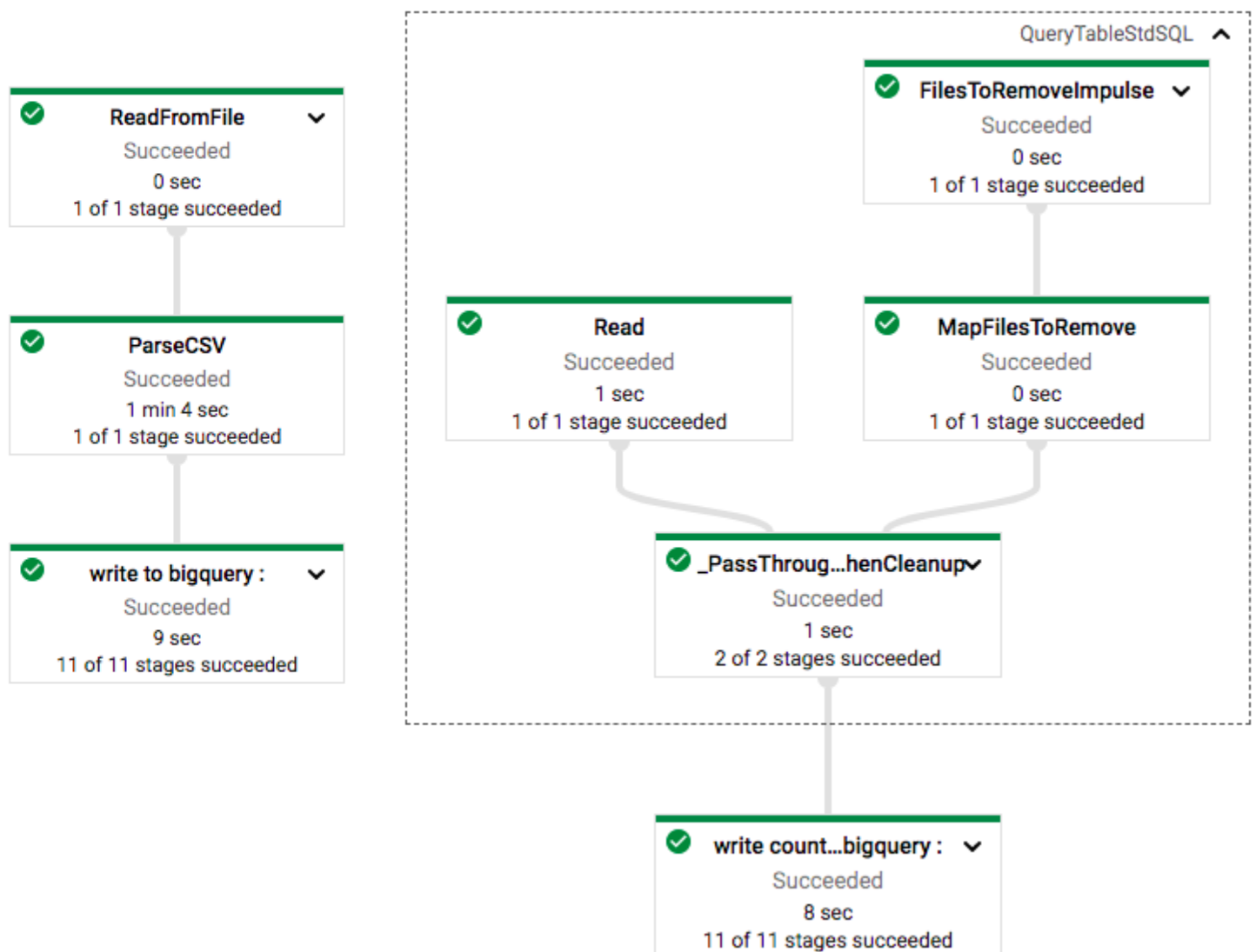


Figure 2: Job execution DAG in Dataflow UI

