National University of Sciences and Technology (NUST)
School of Electrical Engineering and Computer Science

# Department of Computing

## SE-315: Cloud Computing

**Lab 12: Serverless App Dev – Creating a Streaming Data Pipeline for a Real-Time Dashboard with Dataflow – Part B**

**CLO4: Display** skills to effectively use cloud centric solutions such as serverless application development.

**Date:  11.12.24**

## Lab 12: Serverless App Dev – Creating a Streaming Data Pipeline for a Real-Time Dashboard with Dataflow – Part B

### Introduction:

This lab task is an extension of previous lab task where you owned a fleet of New York City taxi cabs and you were live monitoring your business. In this lab, you have to build a similar streaming data pipeline but for the citywide payroll data to explore it and visualize the results in a management dashboard.

**Lab Objectives:** In this lab, the students will learn how to:

- Create a Dataflow job from a template

- Stream a Dataflow pipeline into BigQuery

- Monitor a Dataflow pipeline in BigQuery

- Analyze results with SQL

- Visualize key metrics in Looker Studio

### Lab Tasks

Recapture your knowledge of previous lab by going through the following link:

https://www.cloudskillsboost.google/focuses/19077?parent=catalog

Citywide Payroll Data download link:
https://data.cityofnewyork.us/City-Government/Citywide-Payroll-Data-Fiscal-Year-/k397-673e

The list of tasks is given below. Make sure to take screenshots of each task as you will need to add them in the solution section given below.

(Hint: You can take help from previous lab tasks as the list of tasks are same, however, the dataset is changed.)

**Setting up qwilabs account:**



## Task 1. Create a BigQuery dataset

create the **Citywide Payroll** dataset:



Uploading the data on a GCloud Bucket:

creating the payroll_data table:



**attributes within payroll_data table:**

## Task 2. Copy required lab artifacts

run the following commands to move files needed for the Dataflow job.

here, we copy the files from the google storage bucket using the **gcloud storage cp** command



**resulting bucket:**



## Task 3. Set up a Dataflow Pipeline

in this task, we set up a streaming data pipeline to read files from the Cloud Storage bucket and write data to BigQuery. (Dataflow is a serverless way to carry out data analysis.)

Restart the connection to the Dataflow API.

```
student_00_30944efe8ad1@cloudshell:~ (qwiklabs-gcp-04-8540b18fd6be)$ gcloud services disable dataf
.googleapis.com
gcloud services enable dataflow.googleapis.com
Operation "operations/acat.p17-71696123621-4c259990-eb36-427d-8f07-538c4118fcb2" finished successf
y.
Operation "operations/acf.p2-71696123621-f0bb7b27-dfb6-4970-8384-7842a4d363be" finished successful
```

**Create a new streaming pipeline:**

setting up job name, endpoint and dataflow template:



**setting up optional parameters:**

**National University of Sciences and Technology (NUST)**
**School of Electrical Engineering and Computer Science**

**setting up machine type:**

The machine type for Google Compute Engine instances used in your pipeline execution. e.g., n1-standard-1. Learn more ⧉

| ✓ General purpose | Compute optimized | Memory optimized | GPUs |

Machine types for common workloads, optimized for cost and flexibility

**Series**
E2 ▼

CPU platform selection based on availability

**Machine type**
e2-medium (2 vCPU, 1 core, 4 GB memory) ▼

| **vCPU** | **Memory** |
| 1-2 vCPU (1 shared core) | 4 GB |

Service account email ▼

The email address of the service account to run the job as

Additional experiments

Additional experiment flags for the job, e.g., experiment1, experiment2

**Worker IP Address Configuration**
Unspecified ▼

Dataflow workers can be configured to use public or internal IP addresses. Learn more ⧉

finally, our `payroll` pipeline job is created:



# Task 4. Analyze the taxi data using BigQuery

I ran the select query to display top 10 rows;



**top 10 rows in payroll_data table:**

## Task 5. Perform aggregations on the stream for reporting



**Saving the query:**

## Task 6. Stop the Dataflow Job

## Stop job

**Cancel**
Dataflow will immediately stop this job and abort all data ingestion and processing. Any buffered data may be lost.

**Drain**
Dataflow will cease all data ingestion, but will attempt to finish processing any remaining buffered data. Pipeline resources will be maintained until buffered data has finished processing and any pending output has finished writing.

**Force Cancel**
Dataflow will force cancel this job. This option terminates a job that has become stuck in the cancelation process.
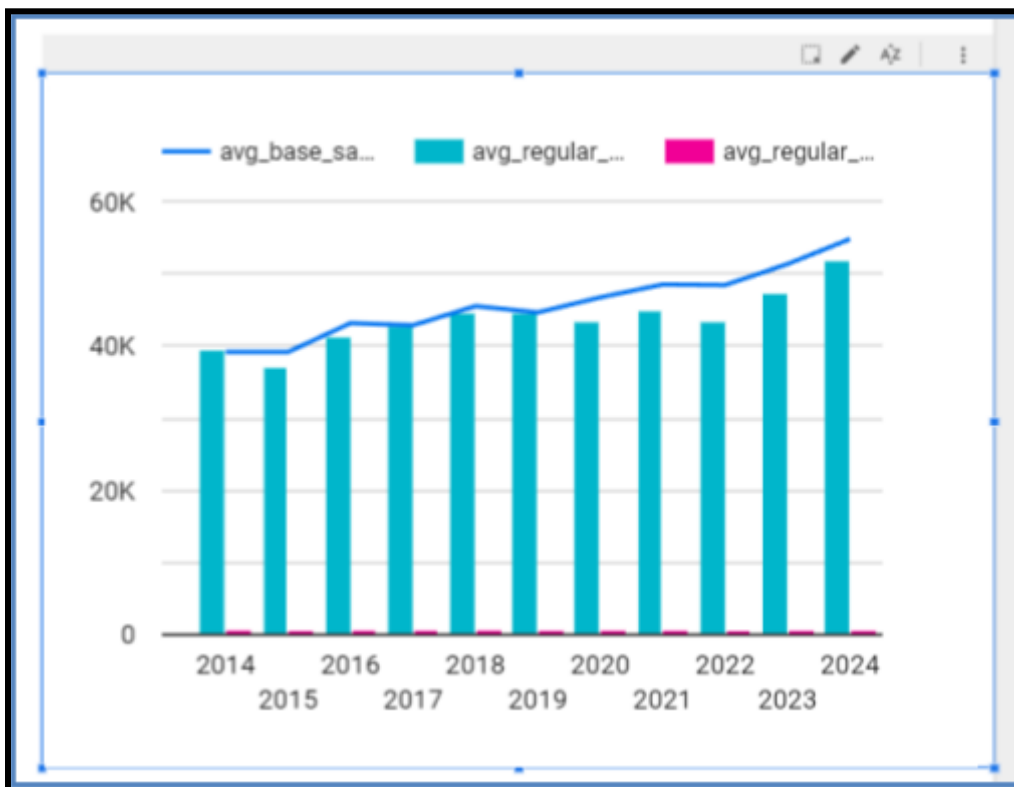
Read more about stopping Dataflow jobs ☑

DO NOTHING    **STOP JOB**

## Task 7. Create a real-time dashboard

**running the saved query and creating the chart  displaying salaries:**

## Task 8. Create a time series dashboard

**adding the custom query to display average salary from payroll_data:**



**Create a time series chart:**

National University of Sciences and Technology (NUST)
School of Electrical Engineering and Computer Science

Remaining credits:

## Credits

**ALL CREDITS**

View and download credit details here. Active committed use discounts are not included here and can be viewed on the Commitments page.

≡ Filter    Filter credits

| Credit name | Status ↑ | Percent remaining | | Remaining value | Original value | Type | Credit ID | |
|---|---|---|---|---|---|---|---|---|
| SE315: Cloud Computing | ✓ Available | ▬▬▬▬▬ | 96% | $48.16 | $50.00 | One-time | 4BPQ3J | ⌄ |