

[End Lab](#)

01:44:44

**Caution:** When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more](#).

[Open Google Console](#)

Username

student-04-0ab80e9a2f



Password

Ypj11mzr4yVg



GCP Project ID

qwiklabs-gcp-02-9f37c



QL Region

us-central1



# Streaming Data Processing: Streaming Data Pipelines into Bigtable

2 hours      Free      ★★★★1

[Overview](#)[Objectives](#)[Setup](#)[Task 1: Preparation](#)[Task 2: Simulate traffic sensor data into Pub/Sub](#)[Task 3: Launch Dataflow Pipeline](#)[Task 4: Explore the pipeline](#)[Task 5: Query Bigtable data](#)[Cleanup](#)[End your lab](#)

## Overview

In this lab, you will use Dataflow to collect traffic events from simulated traffic sensor data made available through Google Cloud PubSub, and write them into a Bigtable table.

At the time of this writing, streaming pipelines are not available in the DataFlow Python SDK. So the streaming labs are written in Java.

## Objectives

In this lab, you will perform the following tasks:

- Launch Dataflow pipeline to read from Pub/Sub and write into Bigtable.
- Open an HBase shell to query the Bigtable database.

## Setup

For each lab, you get a new GCP project and set of resources for a fixed time at no cost.

1. Make sure you signed into Qwiklabs using an **incognito window**.
2. Note the lab's access time (for example, **02:00:00**) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

3. When ready, click **START LAB**.

4. Note your lab credentials. You will use them to sign in to Cloud Platform Console.

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked.  
[Learn more.](#)

**Open Google Console**

**Username**  
student-01-23efd9347325@

**Password**  
gCXLv23N4fPN 

**GCP Project ID**  
qwiklabs-gcp-01-d7c92c04 

5. Click **Open Google Console**.

6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts.

If you use other credentials, you'll get errors or incur charges.

7. Accept the terms and skip the recovery resource page.

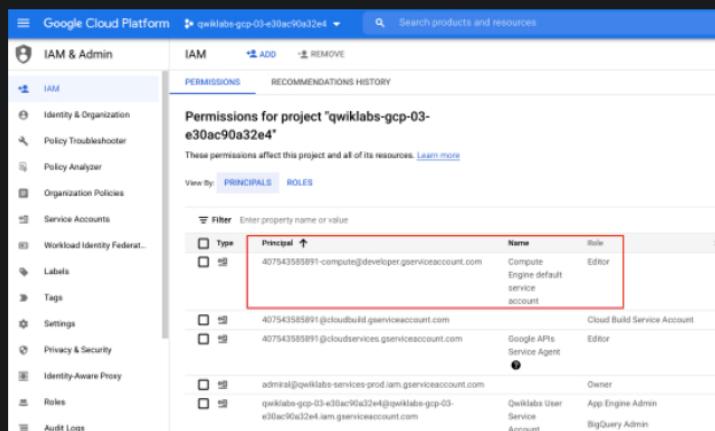
Do not click **End Lab** unless you are finished with the lab or want to restart it. This clears your work and removes the project.

## Check project permissions

Before you begin your work on Google Cloud, you need to ensure that your project has the correct permissions within Identity and Access Management (IAM).

1. In the Google Cloud console, on the **Navigation menu** (≡), click **IAM & Admin > IAM**.

2. Confirm that the default compute Service Account {project-number}-compute@developer.gserviceaccount.com is present and has the editor role assigned. The account prefix is the project number, which you can find on **Navigation menu > Home**.



The screenshot shows the Google Cloud Platform IAM & Admin interface. On the left, the navigation menu is open, showing options like IAM, Identity & Organization, Policy Troubleshooter, Policy Analyzer, Organization Policies, Service Accounts, Workload Identity Federation, Labels, Tags, Settings, Privacy & Security, Identity-Aware Proxy, Roles, and Audit Logs. The main area is titled 'Permissions for project "qwiklabs-gcp-03-e30ac90a32e4"'. It displays a table of permissions for the project. The table has columns for Type, Principal, Name, and Role. One row is highlighted with a red border:

Type	Principal	Name	Role
Principal	407543585891-compute@developer.gserviceaccount.com	Compute Engine default service account	Editor
Principal	407543585891@cloudbuild.gserviceaccount.com	Cloud Build Service Account	
Principal	407543585891@cloudservices.gserviceaccount.com	Google APIs Service Agent	Editor
Principal	admin@qwiklabs-services-prod.iam.gserviceaccount.com	Qwiklabs User	Owner
Principal	qwiklabs-gcp-03-e30ac90a32e4@qwiklabs-gcp-03-e30ac90a32e4.iam.gserviceaccount.com	App Engine Admin	
Principal	qwiklabs-gcp-03-e30ac90a32e4.iam.gserviceaccount.com	BigQuery Admin	

If the account is not present in IAM or does not have the editor role, follow the steps below to assign the required role.

- In the Google Cloud console, on the **Navigation menu**, click **Home**.
- Copy the project number (e.g. 729328892908).

- On the **Navigation menu**, click **IAM & Admin > IAM**.

- At the top of the **IAM** page, click **Add**.

- For **New principals**, type:

```
{project-number}-compute@developer.gserviceaccount.com
```



Replace `{project-number}` with your project number.

- For **Role**, select **Project (or Basic) > Editor**. Click **Save**.

## Task 1: Preparation

You will be running a sensor simulator from the training VM. There are several files and some setup of the environment required.

Open the SSH terminal and connect to the training VM

1. In the Console, on the **Navigation menu** (≡), click **Compute Engine > VM instances**.
2. Locate the line with the instance called **training-vm**.
3. On the far right, under **Connect** column, click on **SSH** to open a terminal window. Then click **Connect**.

In this lab, you will enter CLI commands on the **training-vm**.

Verify initialization is complete

4. The **training-vm** is installing some software in the background. Verify that setup is complete by checking the contents of the new directory.

```
ls /training
```



The setup is complete when the result of your list (`ls`) command output appears as in the image below. If the full listing does not appear, wait a few minutes and try again. **Note:** It may take 2 to 3 minutes for all background actions to complete.

```
student-04-2324ale71896@training-vm:~$ ls /training
bq_magic.sh project_env.sh sensor_magic.sh
student-04-2324ale71896@training-vm:~$
```

Download Code Repository

5. Next, you will download a code repository for use in this lab.

```
git clone https://github.com/GoogleCloudPlatform/training-data-analyst
```



Set environment variables

6. On the **training-vm** SSH terminal, enter the following:

```
source /training/project_env.sh
```



This script sets the `$DEVSHELL_PROJECT_ID` and `$BUCKET` environment variables.

### Prepare HBase quickstart files

7. In the **training-vm** SSH terminal, run the script to download and unzip the quickstart files (you will later use these to run the HBase shell.)

```
cd ~/training-data-analyst/courses/streaming/process/sandiego  
./install_quickstart.sh
```



Click Check my progress to verify the objective.



Copy sample files to the training\_vm home directory

[Check my progress](#)

*Assessment Complete!*

## Task 2: Simulate traffic sensor data into Pub/Sub

1. In the **training-vm** SSH terminal, start the sensor simulator. The script reads sample data from a csv file and publishes it to Pub/Sub.

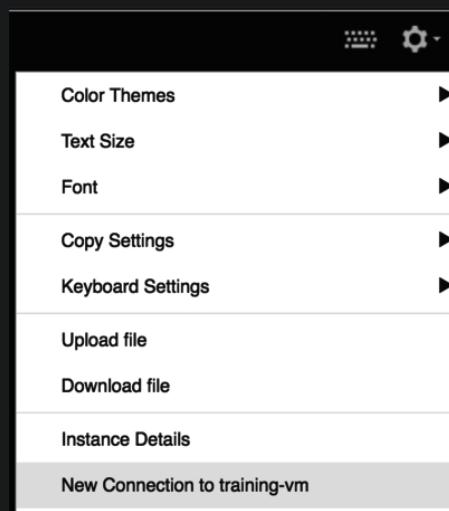
```
/training/sensor_magic.sh
```



This command will send 1 hour of data in 1 minute. Let the script continue to run in the current terminal.

### Open a second SSH terminal and connect to the training VM

2. In the upper right corner of the **training-vm** SSH terminal, click on the gear-shaped button (⚙️), and select **New Connection to training-vm** from the drop-down menu. A new terminal window will open.



[Change Linux Username](#)

[How to copy / paste](#)

[Send Feedback](#)

The new terminal session will not have the required environment variables. Complete the next step to set these variables.

3. In the new **training-vm** SSH terminal, enter the following:

```
source ~/training/project_env.sh
```



Click **Check my progress** to verify the objective.

Simulate traffic sensor data into Pub/Sub



[Check my progress](#)

*Assessment Complete!*

## Task 3: Launch Dataflow Pipeline

1. In the second **training-vm** SSH terminal, navigate to the directory for this lab. Examine the script in Cloud Shell or using nano. **Do not make any changes to the code.**

```
cd ~/training-data-analyst/courses/streaming/process/sandiego  
nano run_oncloud.sh
```



What does the script do?

The script takes 3 required arguments: project id, bucket name, classname and possibly a 4th argument: options. In this part of the lab, we will use the `--bigtable` option which will direct the pipeline to write into Cloud Bigtable.

2. Press **CTRL+X** to exit.

3. Run the following script to create the Bigtable instance.

```
cd ~/training-data-analyst/courses/streaming/process/sandiego  
../create_cbt.sh
```



4. Run the Dataflow pipeline to read from PubSub and write into Cloud Bigtable.

```
cd ~/training-data-analyst/courses/streaming/process/sandiego  
.run_oncloud.sh $DEVSHELL_PROJECT_ID $BUCKET CurrentConditions  
--bigtable
```



Example successful run:

```
[INFO] -----  
-----  
[INFO] BUILD SUCCESS  
[INFO] -----  
-----  
[INFO] Total time: 47.582 s  
[INFO] Finished at: 2018-06-08T21:25:32+00:00  
[INFO] Final Memory: 58M/213M  
[INFO] -----
```

Click Check my progress to verify the objective.

Launch Dataflow Pipeline



Check my progress

Assessment Complete!

## Task 4: Explore the pipeline

1. Return to the browser tab for Console. On the **Navigation menu** (≡), click

**Dataflow** and click on the new pipeline job. Confirm that the pipeline job is listed and verify that it is running without errors.

2. Find the **write:cbt** step in the pipeline graph, and click on the down arrow on the right to see the writer in action. Click on the given writer. Review the **Bigtable Options** in the **Step summary**.

## Task 5: Query Bigtable data

1. In the second **training-vm** SSH terminal, run the **quickstart.sh** script to launch the HBase shell.

```
cd ~/training-data-analyst/courses/streaming/process/sandiego/quickstart  
./quickstart.sh
```

2. When the script completes, you will be in an HBase shell prompt that looks like this:

```
hbase(main):001:0>
```

3. At the HBase shell prompt, type the following query to retrieve 2 rows from your Bigtable table that was populated by the pipeline. It may take a few minutes for results to return via the HBase query.

Repeat the 'scan' command until you see a list of rows returned.

```
scan 'current_conditions', {'LIMIT' => 2}
```

```
hbase(main):006:0> scan 'current_conditions', {'LIMIT' => 2}  
row=15481#9223370#134275#0#  
  column=lane:direction, timestamp=1225912000, value=S  
  column=lane:highway, timestamp=122512000, value=15  
  column=lane:lane, timestamp=1225572000, value=1.0  
  column=lane:latitude, timestamp=122512000, value=32.23248  
  column=lane:longitude, timestamp=122512000, value=-117.11563  
  column=lane:speed, timestamp=1225572000, value=74.8  
  column=lane:timestamp, timestamp=1225517000, value=2008-11-01 04:00:00  
  column=lane:vector, timestamp=1225517000, value=S  
  column=lane:highway, timestamp=1225117000, value=15  
  column=lane:lane, timestamp=1225517000, value=1.0  
  column=lane:latitude, timestamp=1225517000, value=32.705184  
  column=lane:longitude, timestamp=1225517000, value=-117.120565  
  column=lane:semanticid, timestamp=1225517000, value=32.705184,117.120565,15,8,1  
  column=lane:speed, timestamp=1225517000, value=74.8  
  column=lane:timestamp, timestamp=1225517000, value=2008-11-01 03:55:00  
2 row(s) in 0.2840 seconds
```

4. Review the output. Notice each row is broken into column, timestamp, value combinations.

5. Run another query. This time look only at the **lane: speed** column, limit to 10 rows, and specify **rowid patterns** for start and end rows to scan over.

Specify terms patterns for start and end rows to scan over.

```
scan 'current_conditions', {'LIMIT' => 10, STARTROW => '15#S#1',  
ENDROW => '15#S#999', COLUMN => 'lane:speed'}
```

6. Review the output. Notice that you see 10 of the column, timestamp, value combinations, all of which correspond to Highway 15. Also notice that column is restricted to **lane: speed**.

7. Feel free to run other queries if you are familiar with the syntax. Once you're satisfied, enter `quit` to exit the shell.

```
quit
```

## Cleanup

1. In the second **training-vm** SSH terminal, run the following script to delete your Bigtable instance.

```
cd ~/training-data-analyst/courses/streaming/process/sandiego  
../delete_cbt.sh
```

If prompted to confirm, enter `Y`.

2. On your Dataflow page in your Cloud Console, click on the pipeline job name.

3. Click **Stop** on the top menu bar. Select **Cancel**, and then click **Stop Job**.

4. Go back to the first SSH terminal with the publisher, and enter `Ctrl+C` to stop it.

5. In the BigQuery console, click on the three dots next to the **demos** dataset, and click **Delete**.

6. Type **delete** and then click **Delete**.

## End your lab

When you have completed your lab, click **End Lab**. Qwiklabs removes the resources you've used and cleans the account for you.

You will be given an opportunity to rate the lab experience. Select the applicable number of stars, type a comment, and then click **Submit**.

The number of stars indicates the following:

- 1 star = Very dissatisfied
- 2 stars = Dissatisfied
- 3 stars = Neutral
- 4 stars = Satisfied
- 5 stars = Very satisfied

You can close the dialog box if you don't want to provide feedback.

For feedback, suggestions, or corrections, please use the **Support** tab.

trademarks of Google LLC. All other company and product names may be trademarks of the respective companies with which they are associated.