

# Serverless Data Processing with Dataflow - Monitoring, Logging and Error Reporting for Dataflow Jobs

2 hours

No cost

## Overview

In this lab, you:

- Create an alerting policy.
- Perform simple troubleshooting for pipelines.
- See how the Diagnostic and BQ tabs work.
- Explore the Error Reporting page.

## Prerequisites

Basic familiarity with Python.

# Setup and requirements

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

1. Sign in to Qwiklabs using an **incognito window**.
2. Note the lab's access time (for example, **1:15:00**), and make sure you can finish within that time.  
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 (**Username** and **Password**). You will use them to sign in to the Google Cloud Console.
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 receive errors or **incur charges**.
7. Accept the terms and skip the recovery resource page.

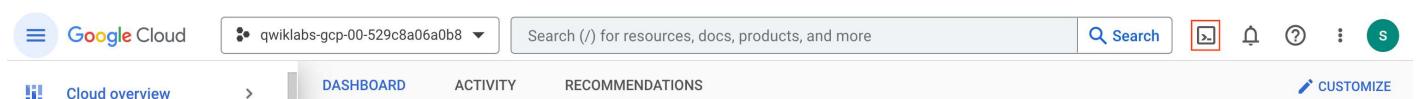
**Note:** Do not click **End Lab** unless you have finished the lab or want to restart it. This clears your work and removes the project.

## Activate Google Cloud Shell

Google Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud.

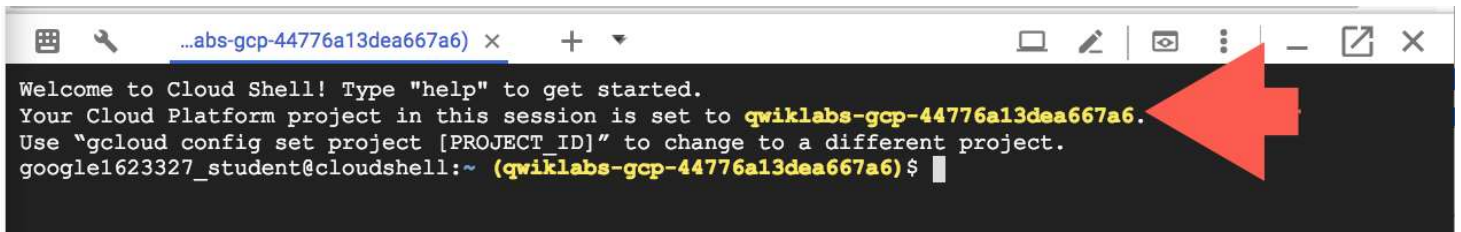
Google Cloud Shell provides command-line access to your Google Cloud resources.

1. In Cloud console, on the top right toolbar, click the Open Cloud Shell button.



## 2. Click **Continue**.

It takes a few moments to provision and connect to the environment. When you are connected, you are already authenticated, and the project is set to your *PROJECT\_ID*. For example:



```
...abs-gcp-44776a13dea667a6) x + ▾
Welcome to Cloud Shell! Type "help" to get started.
Your Cloud Platform project in this session is set to qwiklabs-gcp-44776a13dea667a6.
Use "gcloud config set project [PROJECT_ID]" to change to a different project.
google1623327_student@cloudshell:~ (qwiklabs-gcp-44776a13dea667a6) $
```

**gcloud** is the command-line tool for Google Cloud. It comes pre-installed on Cloud Shell and supports tab-completion.

- You can list the active account name with this command:

```
gcloud auth list
```

content\_co

### Output:

```
Credentialed accounts:
- @.com (active)
```

### Example output:

```
Credentialed accounts:
- google1623327_student@qwiklabs.net
```

- You can list the project ID with this command:

```
gcloud config list project
```

content\_co

### Output:

```
[core]  
project =
```


### Example output:

```
[core]  
project = qwiklabs-gcp-44776a13dea667a6
```

**Note:** Full documentation of **gcloud** is available in the [gcloud CLI overview guide](#) .

## 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** () , select **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 > Cloud Overview > Dashboard**.

## Permissions for project "qwiklabs-gcp-00-3f97701829bb"

These permissions affect this project and all of its resources. [Learn more](#)

☐ Include Google-provided role grants

VIEW BY PRINCIPALS

VIEW BY ROLES

+ GRANT ACCESS

- REMOVE ACCESS

Filter Enter property name or value

?

|||

<input type="checkbox"/> Type	Principal ↑	Name	Role	Security insights ?	Inheritance
<input type="checkbox"/>	96496971506-compute@developer.gserviceaccount.com	Compute Engine default service account	Editor Owner		
<input type="checkbox"/>	admiral@qwiklabs-services-prod.iam.gserviceaccount.com		Owner		
<input type="checkbox"/>	qwiklabs-gcp-00-3f97701829bb@qwiklabs-gcp-00-3f97701829bb.iam.gserviceaccount.com	Qwiklabs User Service Account	BigQuery Admin Owner Storage Admin		
<input type="checkbox"/>	student-03-93dbfa673ace@qwiklabs.net	student 7451284e	App Engine Admin BigQuery Admin Dataflow Admin Dataflow Developer Editor Owner Viewer		

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

1. In the Google Cloud console, on the **Navigation menu**, click **Cloud Overview > Dashboard**.
2. Copy the project number (e.g. 729328892908).
3. On the **Navigation menu**, select **IAM & Admin > IAM**.
4. At the top of the roles table, below **View by Principals**, click **Grant Access**.
5. For **New principals**, type:

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

content\_co

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

7. For **Role**, select **Project** (or Basic) > **Editor**.

8. Click **Save**.

## Launch Google Cloud Shell Code Editor

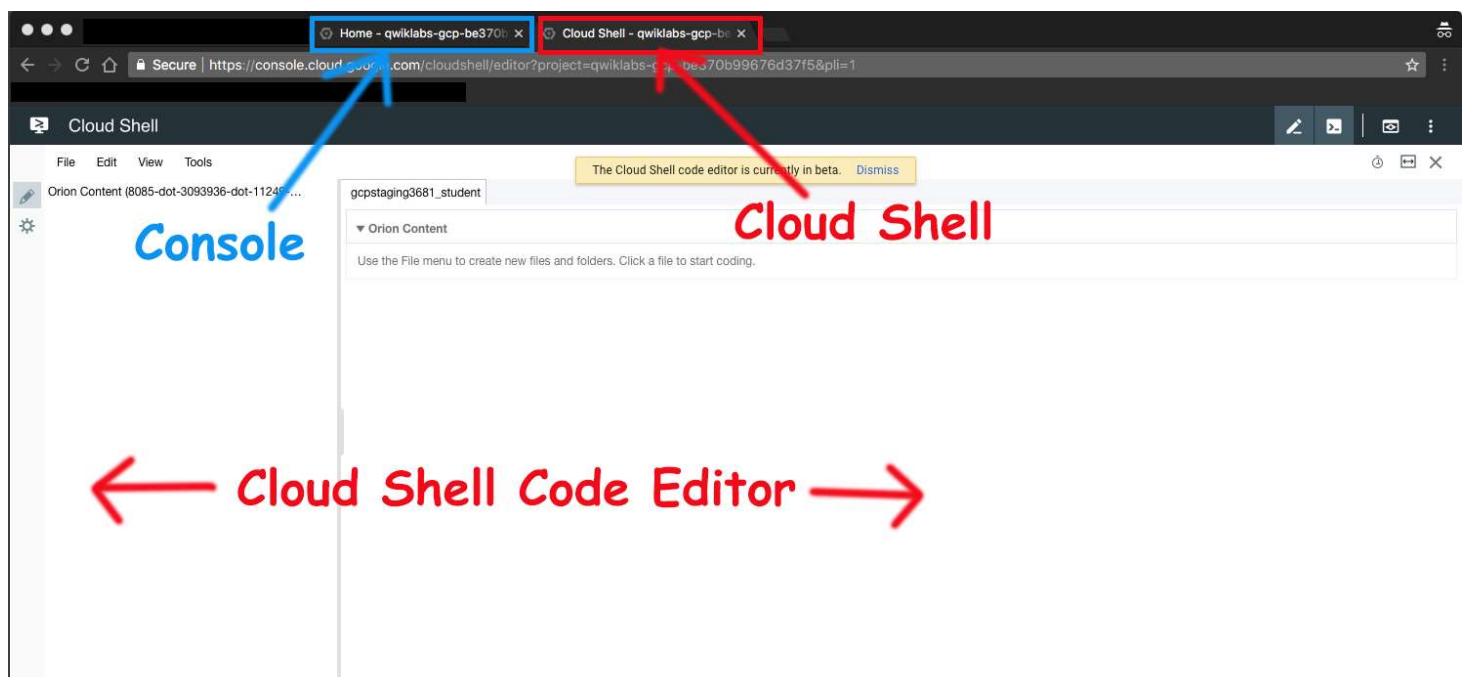
Use the Google Cloud Shell Code Editor to easily create and edit directories and files in the Cloud Shell instance.

- Once you activate the Google Cloud Shell, click **Open editor** to open the Cloud Shell Code Editor.



You now have three interfaces available:

- The Cloud Shell Code Editor
- Console (By clicking on the tab). You can switch back and forth between the Console and Cloud Shell by clicking on the tab.
- The Cloud Shell Command Line (By clicking on **Open Terminal** in the Console)



# Task 1. Create virtual environment and install dependencies

1. Create a virtual environment for your work in this lab:

```
sudo apt-get install -y python3-venv  
python3 -m venv df-env  
source df-env/bin/activate
```

content\_co

2. Next, install the packages you will need to execute your pipeline:

```
python3 -m pip install -q --upgrade pip setuptools wheel  
python3 -m pip install apache-beam[gcp]
```

content\_co

3. Finally, ensure that the Dataflow API is enabled:

```
gcloud services enable dataflow.googleapis.com
```

content\_co

## Task 2. Create alerting policy

In this section, you will create an alert policy that gets triggered when the condition is met.

1. In the Cloud Console, from the **Navigation menu**, select **Monitoring > Alerting**.
2. On the Alerting page, click on + **CREATE POLICY**.
3. Click on **Select a metric** dropdown and disable the **Show only active resources & metrics**.

a. Type **Dataflow Job** in filter by resource and metric name and click on **Dataflow Job > Job**. Select **Failed** and click **Apply**.

b. Set **Rolling windows function** to **Sum**.

c. Click **Next**. Set **0** as your **Threshold value**.

4. Click on the **NEXT** button to go to the **Who should be notified?** step.

5. On the form that loads, click on the **Notification Channels** drop-down menu, and click on **MANAGE NOTIFICATION CHANNELS**.

This opens a new window that lists the supported notification channel types.

6. Scroll down to the row that says **Email** and click on **ADD NEW** on the far right.

a. For **Email Address**, enter your personal email address.

**Note:** The Qwiklabs student account email will not work as it is a temporary account.

b. For **Display Name**, enter **Qwiklabs Student**.

c. Click on **Save**. You can now close this window and go back to the previous **Create alerting policy** window.

7. Click on the refresh icon to the left of **MANAGE NOTIFICATION CHANNELS**.

8. Next, click on the **Notification Channels** drop-down menu again. This time you should see the display name of the student account you just added.

9. Check the checkbox to the left of the name **Qwiklabs Student** and click **OK**.

10. Enter **Failed Dataflow Job** as the **Alert Name** in the textbox.

11. Click **Next** again.

12. Review the alert and click **Create Policy**.



Click *Check my progress* to verify the objective.



Create Alerting Policy

Check my progress

*Assessment Completed!*

## Task 3. Failure starting the pipeline on Dataflow

1. Review the code for your pipeline below:

```
import argparse
import logging
import argparse, logging, os
import apache_beam as beam
from apache_beam.io import WriteToText
from apache_beam.options.pipeline_options import PipelineOptions

class ReadGBK(beam.DoFn):

    def process(self, e):
        k, elems = e
        for v in elems:
            logging.info(f"the element is {v}")
            yield v

def run(argv=None):
    parser = argparse.ArgumentParser()
    parser.add_argument(
        '--output', dest='output', help='Output file to write results
to.')
    known_args, pipeline_args = parser.parse_known_args(argv)
    read_query = """(
```

content\_co

```

SELECT
    version,
    block_hash,
    block_number
FROM
    `bugquery-public-data.crypto_bitcoin.transactions`
WHERE
    version = 1
LIMIT
    1000000 )
UNION ALL (
SELECT
    version,
    block_hash,
    block_number
FROM
    `bigquery-public-data.crypto_bitcoin.transactions`
WHERE
    version = 2
LIMIT
    1000 ) ;"""
p = beam.Pipeline(options=PipelineOptions(pipeline_args))
(p
| 'Read from BigQuery' >>
beam.io.ReadFromBigQuery(query=read_query, use_standard_sql=True)
| "Add Hotkey" >> beam.Map(lambda elem: (elem["version"], elem))
| "Groupby" >> beam.GroupByKey()
| 'Print' >> beam.ParDo(ReadGBK())
| 'Sink' >> WriteToText(known_args.output))

result = p.run()

if __name__ == '__main__':
    logger = logging.getLogger().setLevel(logging.INFO)
    run()

```

2. In **Cloud Shell**, open a new file named `my_pipeline.py`. Copy and paste the code from above into this file and save it, using your preferred text editor (the example below refers to Vim):

```
vi my_pipeline.py
```

content\_copy

3. After you paste the code into the file, be sure to save it.
4. Run the command below to create a storage bucket:

```
export PROJECT_ID=$(gcloud config get-value project)
gsutil mb -l US gs://$PROJECT_ID
```

content\_copy

5. Next, you will attempt to launch the pipeline using the command below (note that the command will fail):

```
# Attempt to launch the pipeline
python3 my_pipeline.py \
  --project=${PROJECT_ID} \
  --region=us-central1 \
  --tempLocation=gs://${PROJECT_ID}/temp/ \
  --runner=DataflowRunner
```

content\_co

The above command to launch the pipeline will fail with a stacktrace similar to the one shown in the screenshot below:

```
Traceback (most recent call last):
  File "my_pipeline.py", line 56, in <module>
    run()
  File "my_pipeline.py", line 50, in run
    | 'Sink' >> WriteToText(known_args.output))
  File "/home/student_00_380e9clf411a/df-env/lib/python3.7/site-packages/apache_beam/io/textio.py", line 653, in __init__
    footer)
  File "/home/student_00_380e9clf411a/df-env/lib/python3.7/site-packages/apache_beam/io/textio.py", line 399, in __init__
    compression_type=compression_type)
  File "/home/student_00_380e9clf411a/df-env/lib/python3.7/site-packages/apache_beam/io/filebasedsink.py", line 90, in __init__
    'got %r instead' % file_path_prefix)
```

This failure is on the Beam end. In the code, we specify `WriteToText(known_args.output)`. Since we did not pass in the `--output` flag, the code did not pass Beam validation, and was not able to launch in Dataflow. As this did not reach Dataflow, no job ID is associated with this launch operation. This means you should not get any alert email in your inbox.

Click *Check my progress* to verify the objective.



Failure starting the pipeline on Dataflow

Check my progress

*Assessment Completed!*

## Task 4. Invalid BigQuery table

In this section, you try to launch the pipeline again with the required **--output** flag. While this does succeed in launching the Dataflow pipeline, the job will fail after a few minutes because we intentionally add an invalid BigQuery table. This should trigger the alerting policy and result in an alert email being sent.

1. In the terminal, execute the following command to launch the pipeline:

```
# Launch the pipeline
python3 my_pipeline.py \
  --project=${PROJECT_ID} \
  --region=us-central1 \
  --output gs://${PROJECT_ID}/results/prefix \
  --tempLocation=gs://${PROJECT_ID}/temp/ \
  --max_num_workers=5 \
  --runner=DataflowRunner
```

content\_co

2. In the Cloud Console, and from the **Navigation menu**, select **Dataflow > Jobs**.

3. Click on the Dataflow job listed on that page.

If you wait for about four minutes, you will see that the job fails. This triggers the alerting policy and you can look in your email inbox to view the alert email.

4. Below the Dataflow pipeline job graph, you will see the **Logs** panel. Expand the panel by clicking on the expansion icon to the far right of the word "Logs". This will expose the following tabs: **Jobs logs**, **Worker logs**, **Diagnostics**, and **BigQuery jobs**.

5. Click on the **Job logs** tab and explore the log entries.

6. Next, click on the **Worker logs** tab and explore those logs as well.

## Logs



### JOB LOGS

### WORKER LOGS

### ⚠️ DIAGNOSTICS

### BIGQUERY JOBS

Showing 19 messages

Severity

Info

Filter

Filter logs



**i** No older entries found matching current filter.

**i** 2021-05-19T00:09:52.727424267Z Worker configuration: n1-standard-1 in us-central1-c.

▶ **i** 2021-05-19T00:09:54.622737178Z Executing operation Read from BigQuery/FilesToRemoveImpulse/Read+Read f...

▶ **i** 2021-05-19T00:09:54.653035649Z Executing operation Sink/Write/WriteImpl/DoOnce/Read+Sink/Write/WriteIm...

**i** 2021-05-19T00:09:54.690022805Z Executing operation Groupby/Create

**i** 2021-05-19T00:09:54.700000000Z Starting 1 workers in us-central1-c.

**Note:** You will notice that the logs can be expanded when you click on them. Logs that get repeated can appear in the **Diagnostics** tab.

7. When you click on the **Diagnostics** tab, you will see a clickable link.

8. Click on it and this takes you to the Error Reporting page. This shows us an expanded view of the error.  
Looking at the stack trace, you will see the issue: a typo in the pipeline code!

Parsed Raw

```
apitools.base.py.exceptions.HttpForbiddenError: HttpError accessing <https://bigquery.googleapis.com/bigquery/v2/projects/qwiklabs-gcp-01-6d5e1292b78e/jobs?alt=json>: response: <{'vary': 'Origin, X-Origin, Referer', 'server': 'ESF', 'cache-control': 'private', 'x-xss-protection': '0', 'x-frame-options': 'SAMEORIGIN', 'x-content-type-options': 'nosniff', 'transfer-encoding': 'chunked', 'status': '403', 'error': {'code': 403, 'message': 'Access Denied: Table bugquery-public-data:crypto_bitcoin.transactions: User does not have permission to query table bugquery-public-data:crypto_bitcoin.transactions.', 'errors': data:crypto_bitcoin.transactions: User does not have permission to query table bugquery-public-data:crypto_bitcoin.transactions.', 'domain': 'global', 'reason': 'accessDenied' } }>, 'status': 'PERMISSION_DENIED' }>
  at __ProcessHttpResponse (/usr/local/lib/python3.7/site-packages/apitools/base/py/base_api.py:604)
  at ProcessHttpResponse (/usr/local/lib/python3.7/site-packages/apitools/base/py/base_api.py:737)
  at _RunMethod (/usr/local/lib/python3.7/site-packages/apitools/base/py/base_api.py:731)
  at Insert (/usr/local/lib/python3.7/site-packages/apache_beam/io/gcp/internal/clients/bigquery/bigquery_v2_client.py:345)
  at get_query_location (/usr/local/lib/python3.7/site-packages/apache_beam/io/gcp/bigquery_tools.py:416)
  at wrapper (/usr/local/lib/python3.7/site-packages/apache_beam/utils/retry.py:253)
  at _setup_temporary_dataset (/usr/local/lib/python3.7/site-packages/apache_beam/io/gcp/bigquery.py:846)
  at _f (/usr/local/lib/python3.7/site-packages/apache_beam/options/value_provider.py:193)
  at split (/usr/local/lib/python3.7/site-packages/apache_beam/io/gcp/bigquery.py:807)
  at _perform_source_split (/usr/local/lib/python3.7/site-packages/dataflow_worker/executor.py:305)
  at _perform_source_split_considering_api_limits (/usr/local/lib/python3.7/site-packages/dataflow_worker/executor.py:268)
  at execute (/usr/local/lib/python3.7/site-packages/dataflow_worker/executor.py:260)
  at do_work (/usr/local/lib/python3.7/site-packages/dataflow_worker/batchworker.py:651)
```

In the code, **bigquery** is intentionally misspelled as **bugquery**, so this fails the Dataflow job and triggers the alerting policy to send an email. In the next section, you will fix the code and relaunch the pipeline.

Click *Check my progress* to verify the objective.



Invalid BigQuery table

Check my progress

*Assessment Completed!*

## Task 5. Too much logging

1. Using a text editor (such as Vim or your preferred editor), open the `my_pipeline.py` file and fix the code by replacing `bugquery` with `bigquery`.

**Note:** In the code, you will find the misspelled **bugquery** in the **SELECT** statement as `bugquery-public-data.crypto_bitcoin.transactions`. Change it to `bigquery-public-data.crypto_bitcoin.transactions`.

2. In the terminal, execute the following command to launch the pipeline:

```
# Launch the pipeline
python3 my_pipeline.py \
  --project=${PROJECT_ID} \
  --region=us-central1 \
  --output gs://${PROJECT_ID}/results/prefix \
  --tempLocation=gs://${PROJECT_ID}/temp/ \
  --max_num_workers=5 \
  --runner=DataflowRunner
```

content\_co

The pipeline will take about seven minutes to complete.

3. In the Cloud Console, from the **Navigation menu**, select **Dataflow > Jobs**. When you click on this new job, in the job graph page on the far right you should see **Job status** as `succeeded`. If the job is not complete, please wait for it to complete and succeed.

4. Expand on the **Logs** at the bottom and review the logs in the **Worker logs** tab.

In the **Worker logs** tab, you will see log entries with the format `the element is {'version' : ..... }`, as shown in the screenshot below:

Logs

JOB LOGS

WORKER LOGS

DIAGNOSTICS

BIGQUERY JOBS

Showing 350 messages

Severity

Info

Filter

Filter logs

↺

↻

?

▶

i

2021-05-19T00:40:15.672689676Z

the element is {'version': 1, 'block\_hash': '0000000000000000047df4ffe50dfcf3b28d8c5...

▶

i

2021-05-19T00:40:15.672783851Z

the element is {'version': 1, 'block\_hash': '00000000000000000400f69fc0a650317f8c5cd...

▶

i

2021-05-19T00:40:15.672873973Z

the element is {'version': 1, 'block\_hash': '000000000000000005311386f8eeca5b549b72e...

▶

i

2021-05-19T00:40:15.672965049Z

the element is {'version': 1, 'block\_hash': '000000000000000000dabc3706fcd1b2bbd8fe...

▶

i

2021-05-19T00:40:15.673052787Z

the element is {'version': 1, 'block\_hash': '0000000000000000003ba3cd0299b7395ab6b77b...

▶

i

2021-05-19T00:40:15.673229932Z

the element is {'version': 1, 'block\_hash': '00000000000000000019de8e6eacfea8c46b108e...

▶

i

2021-05-19T00:40:15.673341989Z

the element is {'version': 1, 'block\_hash': '0000000000000000002c239827fe5f926a9c7838...

▶

i

2021-05-19T00:40:15.673447370Z

the element is {'version': 1, 'block\_hash': '000000000000000000376c170783fde4cf3c0b00...

▶

i

2021-05-19T00:40:15.673536539Z

the element is {'version': 1, 'block\_hash': '0000000000000000004c37d3e458ea02cfe7fd70...

These entries are being logged because of the following line in the pipeline code:

logging.info(f"the element is {v}")

content\_co

5. Click on the **Diagnostics** tab and you will see a message about throttling logs.

Logs

JOB LOGS

WORKER LOGS

DIAGNOSTICS

BIGQUERY JOBS


Occurrences	Count	Error	Users	First seen	Last seen
<div><div></div></div>	5	<div><div><b>Throttling logger worker. It used up its 30s quota for logs in only 1.747682095s</b></div><div>There was a high rate of log messages from the job and some of them were not sent to Cloud Logging.</div></div>	-	3 minutes ago	1 m ago


6. Click on it and navigate to the **Error Reporting** page. Excessive logging is raised as a job insight, with a public doc link showing what the issue is.


The simple fix to this issue would be to remove the logging line of code from your pipeline code and rerun the dataflow pipeline. Upon completion, you should no longer see the same message in the **Diagnostics** tab.




**Note:** You are not required to rerun this pipeline in order to complete this lab.

- Go back to the Dataflow Jobs page by clicking the back button on your browser.
- Under the expanded logs panel, click on the **BigQuery jobs** tab. From the **Data Location** drop-down menu, select us (multiple regions in United States) and click on **Load BigQuery jobs**. Then click **Confirm**.

**Logs** 

**JOB LOGS** **WORKER LOGS**  **DIAGNOSTICS** **BIGQUERY JOBS**

 This will incur a billing cost to query the underlying BigQuery table which stores this information. [Learn more](#)

**Data location**  
us (multiple regions in Uni...   **REFRESH** 

Last updated: Just now

ID	Type	End time	Elapsed time	Start time	State	More Info
beam_bq_job_EXP...	EXTRACT	Dec 24, 2021, 2:47:43 PM	2 sec	Dec 24, 2021, 2:47:41 PM	DONE	<b>COMMAND LINE</b>
beam_bq_job_QUE...	QUERY	Dec 24, 2021, 2:47:39 PM	3 sec	Dec 24, 2021, 2:47:35 PM	DONE	<b>COMMAND LINE</b> <b>VIEW QUERY</b>

You will see two jobs listed that were part of the pipeline.

- Click on **Command line** on the far right of the first job, and, on the pop-up, click on **Run in Cloud Shell**.

This will paste the command in your Cloud Shell. Make sure to run the command.

- The output of the command shows the file locations used to read/write to BigQuery, their size, and how long the job took. If a BigQuery job has failed, this would be the place to look for the cause of failure.



```
extract:
  compression: NONE
  destinationFormat: AVRO
  destinationUri: gs://dataflow-staging-us-central1-cla9299613c47de1a561cb7e8d52f29e/574ccccf-7/bigquery-table-dump-*.json
  destinationUris:
  - gs://dataflow-staging-us-central1-cla9299613c47de1a561cb7e8d52f29e/574ccccf-7/bigquery-table-dump-*.json
  printHeader: false
  sourceTable:
    datasetId: temp_dataset_40cd59091f064c75a753b9e4377c0e1b
    projectId: qwiklabs-gcp-02-f8c6a26b2f43
    tableId: temp_table_40cd59091f064c75a753b9e4377c0e1b
  useAvroLogicalTypes: true
  jobType: EXTRACT
  labels:
    beam_job_id: 2021-05-18_17_33_54-12699155996613753449
    step_name: readfrombigquery
etag: FpH0x+VAWS9qKinZMsg2HQ==
id: qwiklabs-gcp-02-f8c6a26b2f43:US.beam_bq_job_EXPORT_BQ_EXPORT_JOB_574ccccf-7_1621384641_564
```

Click *Check my progress* to verify the objective.



Too much logging

[Check my progress](#)

*Please fix "my\_pipeline.py" file code by replacing bugquery with bigquery, rerun your pipeline and wait until job status gets succeed.*

## End your lab