# Troubleshooting & Debugging
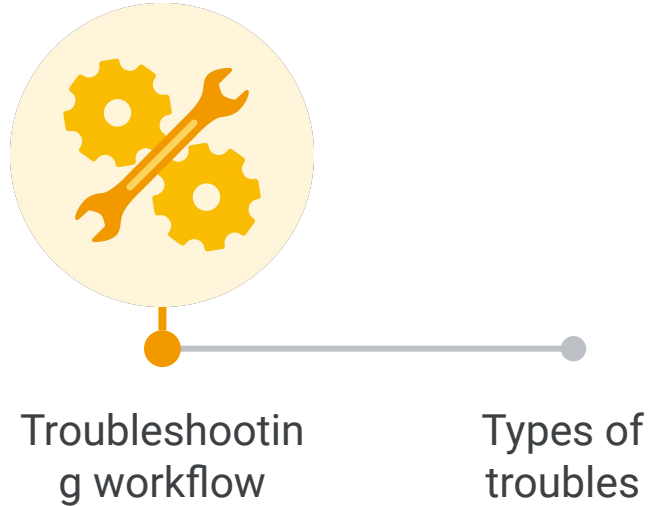
Chirag Shankar
Systems Engineer, Google Cloud
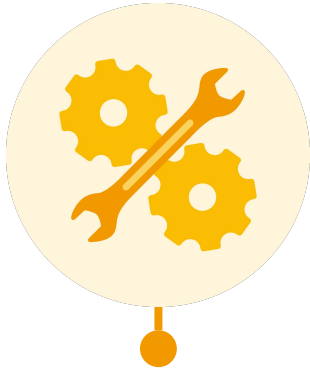
# Troubleshooting & Debug

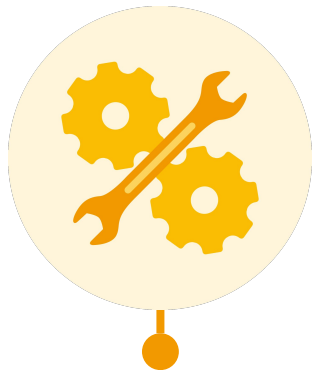Agenda

Troubleshooting workflow

Types of troubles

# Troubleshooting Workflow



Troubleshootin
g workflow

- Checking for errors
- Looking for anomalies in the Job metrics tab

# Troubleshooting Workflow



Troubleshooting workflow

- Checking for errors
  - Logs
  - Execution Graph
- Looking for anomalies in the Job metrics tab

# Checking for errors

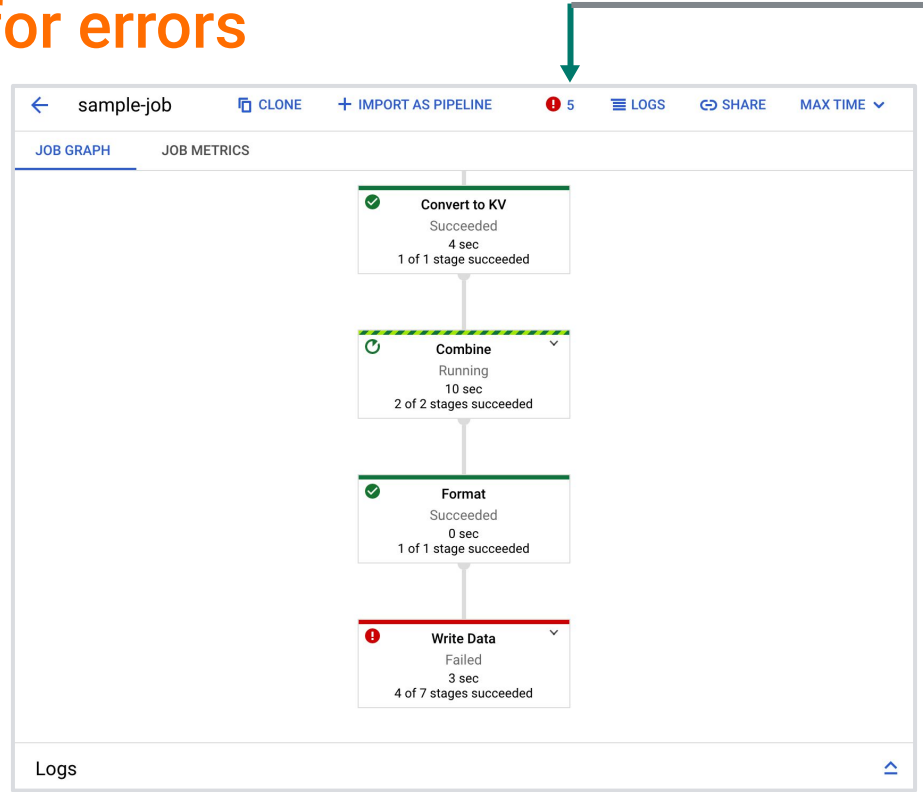| | Name | Type | End time | Elapsed time | Start time | Status | SDK version | ID | Region |
|---|---|---|---|---|---|---|---|---|---|
| ↻ | wordcount4 | Batch | — | 7 sec | Nov 21, 2018, 3:22:45 PM | Running | 2.7.0 | 2018-11-21_15_22_45- | us-central1 |
| ✓ | wordcount3 | Batch | Nov 21, 2018, 1:10:23 PM | 3 min 29 sec | Nov 21, 2018, 1:06:54 PM | Succeeded | 2.7.0 | 2018-11-21_13_06_53- | us-central1 |
| ! | wordcount2 | Batch | Nov 21, 2018, 12:57:05 PM | 3 min 56 sec | Nov 21, 2018, 12:53:09 PM | Failed | 2.7.0 | 2018-11-21_12_53_08- | us-central1 |
| ✓ | wordcount1 | Batch | Nov 21, 2018, 12:55:25 PM | 3 min 36 sec | Nov 21, 2018, 12:51:49 PM | Succeeded | 2.7.0 | 2018-11-21_12_51_48- | us-central1 |

Dataflow | Jobs + CREATE JOB FROM TEMPLATE

Filter jobs

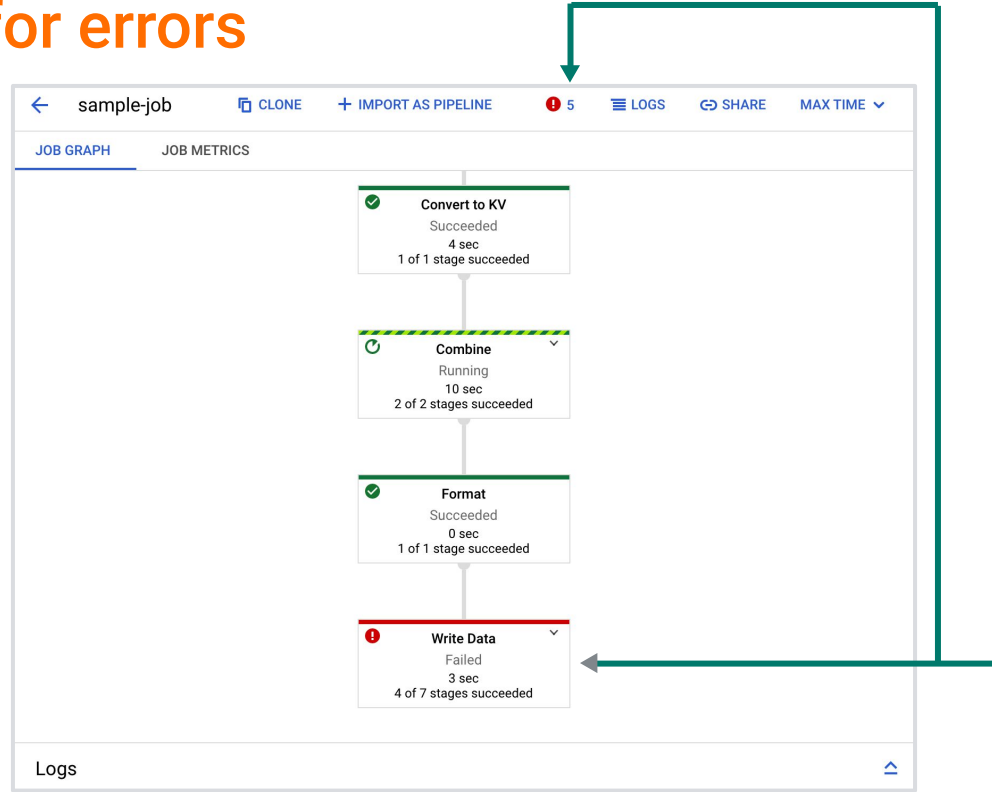# Checking for errors

# Checking for errors



Checking for errors

# Checking for errors



Checking
for errors

# Checking for errors

# Checking for errors

# Checking for errors

# Troubleshooting Workflow



Troubleshootin g workflow

- Checking for errors
- Looking for anomalies in the Job metrics tab

# Looking for anomalies in the Job metrics tab

Streaming pipelines

# Looking for anomalies in the Job metrics tab

Streaming pipelines

# Looking for anomalies in the Job metrics tab

All Dataflow Jobs

# Troubleshooting & Debug

Agenda

Troubleshooting workflow

Types of troubles

# Types of Troubles



Types of
troubles

- Failure building the pipeline
- Failure starting the pipeline on Dataflow
- Failure during pipeline execution
- Performance problems

# Types of Troubles

Types of troubles

- Failure building the pipeline
- Failure starting the pipeline on Dataflow
- Failure during pipeline execution
- Performance problems

# Failure building the pipeline

These errors occur when Apache Beam is building the pipeline:

- Validating "Beam Model" aspects of the pipeline

- Validating input/output specifications

- You can reproduce it with the direct runner, in a unit test

# Failure building the pipeline

```
java.lang.IllegalStateException: GroupByKey cannot be applied to non-bounded
PCollection in the GlobalWindow without a trigger. Use a Window.into or
Window.triggering transform prior to GroupByKey.
    at org.apache.beam.sdk.transforms.GroupByKey.applicableTo (GroupByKey.java:153)
    at org.apache.beam.sdk.transforms.GroupByKey.expand (GroupByKey.java:185)
    at org.apache.beam.sdk.transforms.GroupByKey.expand (GroupByKey.java:107)
    at org.apache.beam.sdk.Pipeline.applyInternal (Pipeline.java:537)
    at org.apache.beam.sdk.Pipeline.applyTransform (Pipeline.java:471)
    at org.apache.beam.sdk.values.PCollection.apply (PCollection.java:357)
    at ...
```

# Types of Troubles

Types of
troubles

- Failure building the pipeline
- **Failure starting the pipeline on Dataflow**
- Failure during pipeline execution
- Performance problems

# Failure starting the pipeline on Dataflow

Once the Dataflow service has received your pipeline's graph, the service will attempt to validate your job:

- Making sure the service can access your job's associated Cloud Storage buckets for file staging and temporary output.

# Failure starting the pipeline on Dataflow

Once the Dataflow service has received your pipeline's graph, the service will attempt to validate your job:

- Making sure the service can access your job's associated Cloud Storage buckets for file staging and temporary output.

- Checking for the required permissions in your Google Cloud project.

# Failure starting the pipeline on Dataflow

Once the Dataflow service has received your pipeline's graph, the service will attempt to validate your job:

- Making sure the service can access your job's associated Cloud Storage buckets for file staging and temporary output.

- Checking for the required permissions in your Google Cloud project.
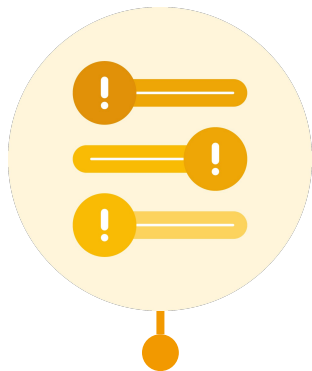
- Making sure the service can access input and output sources, such as files.

# Failure starting the pipeline on Dataflow

apitools.base.py.exceptions.HttpForbiddenError: HttpError accessing
<https://dataflow.googleapis.com/v1b3/projects/xyz/locations/us-central1/jobs?alt=json>:
response: <{'status': '403', 'content-length': '288', 'x-xss-protection': '0',
'x-content-type-options': 'nosniff', 'transfer-encoding': 'chunked', 'vary': 'Origin,
X-Origin, Referer', 'server': 'ESF', '-content-encoding': 'gzip', 'cache-control': 'private',
'date': 'Tue, 10 Dec 2019 17:56:00 GMT', 'x-frame-options': 'SAMEORIGIN', 'alt-svc':
'quic=":443"; ma=2592000; v="46,43",h3-Q050=":443"; ma=2592000,h3-Q049=":443";
ma=2592000,h3-Q048=":443"; ma=2592000,h3-Q046=":443"; ma=2592000,h3-Q043=":443"; ma=2592000',
'content-type': 'application/json; charset=UTF-8'}>, content <{
        "error": {
          "code": 403,
          "message": "(43abf11c2446750): Could not create workflow; user does not have write
      access to project: xyz Causes: (43abf11c244683b): Permission 'dataflow.jobs.create'
      denied on project: 'xyz'",
          "status": "PERMISSION_DENIED"
        }
      }

# Types of Troubles



Types of troubles

- Failure building the pipeline
- Failure starting the pipeline on Dataflow
- **Failure during pipeline execution**
- Performance problems

# Failure during pipeline execution

# Failure during pipeline execution

# Failure during pipeline execution

```python
class FilterMessagesFn(beam.DoFn):
    BAD_MESSAGE_TAG = 'bad_message'
    GOOD_MESSAGE_TAG = 'good_message'

    def process(self, element, window=beam.DoFn.WindowParam):
        try:
            data = element.decode()
            # tag the elements accordingly
            if 'bad' in data:
                yield pvalue.TaggedOutput(self.BAD_MESSAGE_TAG, element)
            else:
                yield pvalue.TaggedOutput(self.GOOD_MESSAGE_TAG, element)


        # handle any exceptions in the processing
        except Exception as exp:
            logging.getLogger.warning(exp)
            yield pvalue.TaggedOutput(self.BAD_MESSAGE_TAG, element)
```

# Failure during pipeline execution

To track failing elements:

- Log the failing elements and check the output using Cloud Logging.

# Failure during pipeline execution

To track failing elements:

- Log the failing elements and check the output using Cloud Logging.

- Check the Dataflow worker and worker startup logs for warnings or errors by following the instructions in Viewing logs.

# Failure during pipeline execution

To track failing elements:

- Log the failing elements and check the output using Cloud Logging.

- Check the Dataflow worker and worker startup logs for warnings or errors by following the instructions in Viewing logs.

- Set your ParDo to write the failing elements to an additional output for later inspection.

# Failure during pipeline execution

Batch vs Streaming

**Batch pipelines** ⟶ **Task** ⟶ Retried upto 4 times

Job fails.

**Streaming pipelines** ⟶ **Task** ⟶ Stall indefinitely

- High Data Freshness
- Job logs
- Cloud Monitoring
- Metrics for pipeline progress and error count

# Types of Troubles



Types of
troubles

- Failure building the pipeline
- Failure starting the pipeline on Dataflow
- Failure during pipeline execution
- Performance problems

# Performance problems

Factors that can influence the performance of a pipeline:

- Pipeline design

- Data shape

- Interactions with sources, sinks and external systems

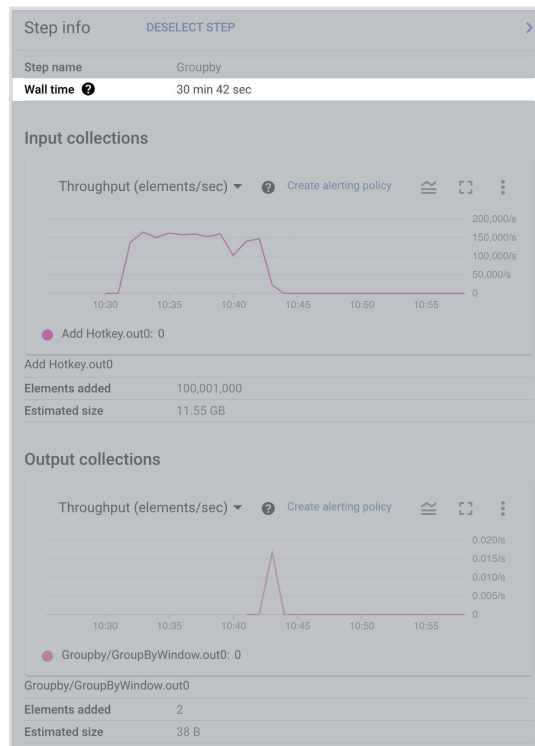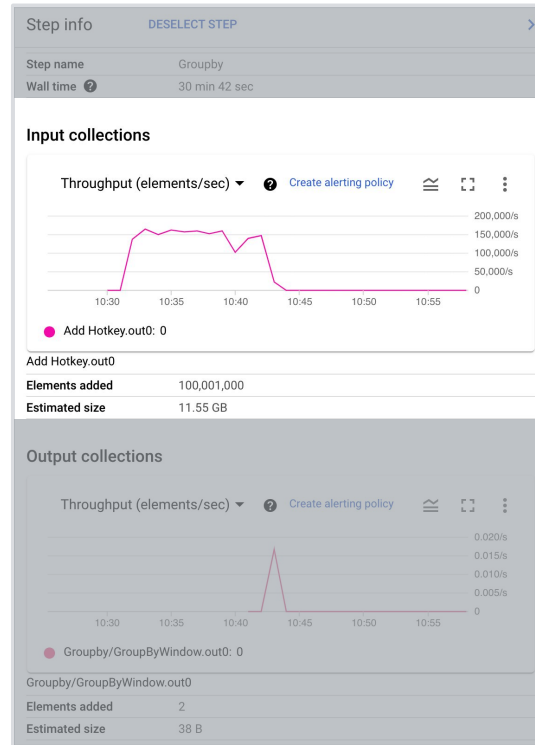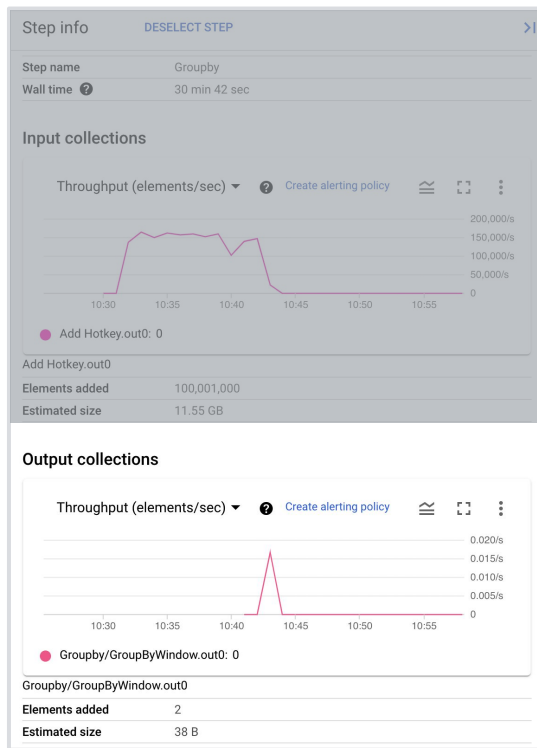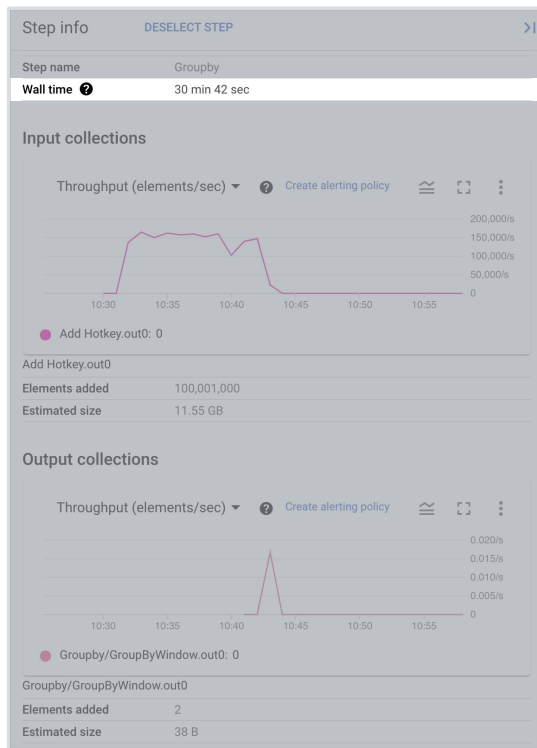# Performance problems



Step info  DESELECT STEP

| Step name | Groupby |
|---|---|
| **Wall time** ? | 30 min 42 sec |

**Input collections**

Throughput (elements/sec) ▾ ? Create alerting policy

200,000/s
150,000/s
100,000/s
50,000/s
0

10:30 10:35 10:40 10:45 10:50 10:55

● Add Hotkey.out0: 0

Add Hotkey.out0
| Elements added | 100,001,000 |
| Estimated size | 11.55 GB |

**Output collections**

Throughput (elements/sec) ▾ ? Create alerting policy

0.020/s
0.015/s
0.010/s
0.005/s
0

10:30 10:35 10:40 10:45 10:50 10:55

● Groupby/GroupByWindow.out0: 0

Groupby/GroupByWindow.out0
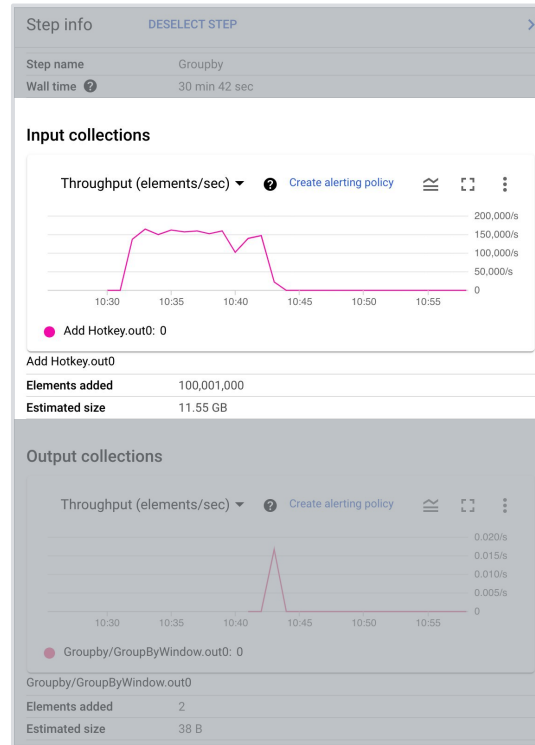| Elements added | 2 |
| Estimated size | 38 B |

# Performance problems

# Performance problems

# Performance problems

# Performance problems

# Thank you!

Questions?