



# Serverless Data Processing with Dataflow

This training is intended for big data practitioners who want to further their understanding of Dataflow in order to advance their data processing applications.

Beginning with foundations, this training explains how Apache Beam and Dataflow work together to meet your data processing needs without the risk of vendor lock-in. The section on developing pipelines covers how you convert your business logic into data processing applications that can run on Dataflow. This training culminates with a focus on operations, which reviews the most important lessons for operating a data application on Dataflow, including monitoring, troubleshooting, testing, and reliability.

## DURATION

3 days

## LEVEL

Advanced

## FORMAT

Instructor led  
On-demand

## What you'll learn

- Demonstrate how Apache Beam and Dataflow work together to fulfill your organization's data processing needs.
- Summarize the benefits of the Beam Portability Framework and enable it for your Dataflow pipelines.
- Enable Shuffle and Streaming Engine, for batch and streaming pipelines respectively, for maximum performance.
- Enable Flexible Resource Scheduling for more cost-efficient performance.
- Select the right combination of IAM permissions for your Dataflow job.
- Implement best practices for a secure data processing environment.
- Select and tune the I/O of your choice for your Dataflow pipeline.
- Use schemas to simplify your Beam code and improve the performance of your pipeline.
- Develop a Beam pipeline using SQL and DataFrames.
- Perform monitoring, troubleshooting, testing and CI/CD on Dataflow pipelines.

Overview	21 Modules · 21 Labs · 81 Videos · 18 Quizzes
Who this course is for	<ul style="list-style-type: none"><li>• Data Engineer</li><li>• Data Analysts and Data Scientists aspiring to develop Data Engineering skills</li></ul>
Products	Dataflow, Cloud Operations
Prerequisite	<ul style="list-style-type: none"><li>• Completed “Building Batch Data Pipelines”</li><li>• Completed “Building Resilient Streaming Analytics Systems”</li></ul>

## Module 01 **Introduction**

Topics	<ul style="list-style-type: none"><li>• Course Introduction</li><li>• Beam and Dataflow Refresher</li></ul>
Objectives	<ul style="list-style-type: none"><li>• Introduce the course objectives.</li><li>• Demonstrate how Apache Beam and Dataflow work together to fulfill your organization’s data processing needs.</li></ul>
Activities	–

## Module 02 **Beam Portability**

Topics	<ul style="list-style-type: none"><li>• Beam Portability</li><li>• Runner v2</li><li>• Container Environments</li><li>• Cross-Language Transforms</li></ul>
Objectives	<ul style="list-style-type: none"><li>• Summarize the benefits of the Beam Portability Framework.</li><li>• Customize the data processing environment of your pipeline using custom containers.</li><li>• Review use cases for cross-language transformations.</li><li>• Enable the Portability framework for your Dataflow pipelines.</li></ul>
Activities	Quiz

## Module 03 **Separating Compute and Storage with Dataflow**

Topics	<ul style="list-style-type: none"><li>• Dataflow</li></ul>
--------	------------------------------------------------------------

---

<b>Topics</b>	<ul style="list-style-type: none"><li>• Dataflow Shuffle Service</li><li>• Dataflow Streaming Engine</li><li>• Flexible Resource Scheduling</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Enable Shuffle and Streaming Engine, for batch and streaming pipelines respectively, for maximum performance.</li><li>• Enable Flexible Resource Scheduling for more cost-efficient performance.</li></ul>
<b>Activities</b>	Quiz

---

#### **Module 04 IAM, Quotas, and Permissions**

<b>Topics</b>	<ul style="list-style-type: none"><li>• IAM</li><li>• Quota</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Select the right combination of IAM permissions for your Dataflow job.</li><li>• Determine your capacity needs by inspecting the relevant quotas for your Dataflow jobs.</li></ul>
<b>Activities</b>	Quiz

---

#### **Module 05 Security**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Data Locality</li><li>• Shared VPC</li><li>• Private IPs</li><li>• CMEK</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Select your zonal data processing strategy using Dataflow, depending on your data locality needs.</li><li>• Implement best practices for a secure data processing environment.</li></ul>
<b>Activities</b>	Hands-on lab and quiz

---

#### **Module 06 Beam Concepts Review**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Beam Basics</li><li>• Utility Transforms</li><li>• DoFn Lifecycle</li></ul>
<b>Objectives</b>	Review main Apache Beam concepts (Pipeline, PCollections, PTransforms, Runner, reading/writing, Utility PTransforms, side inputs), bundles and DoFn Lifecycle.
<b>Activities</b>	Hands-on lab and quiz

**Module 07 Windows, Watermarks, Triggers**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Windows</li><li>• Watermarks</li><li>• Triggers</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Implement logic to handle your late data.</li><li>• Review different types of triggers.</li><li>• Review core streaming concepts (unbounded PCollections, windows).</li></ul>
<b>Activities</b>	Hands-on lab and quiz

---

**Module 08 Sources and Sinks**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Sources and Sinks</li><li>• Text IO and File IO</li><li>• BigQuery IO</li><li>• PubSub IO</li><li>• Kafka IO</li><li>• Bigable IO</li><li>• Avro IO</li><li>• Splittable DoFn</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Write the I/O of your choice for your Dataflow pipeline.</li><li>• Tune your source/sink transformation for maximum performance.</li><li>• Create custom sources and sinks using SDF.</li></ul>
<b>Activities</b>	Quiz

---

**Module 09 Schemas**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Beam Schemas</li><li>• Code Examples</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Introduce schemas, which give developers a way to express structured data in their Beam pipelines.</li><li>• Use schemas to simplify your Beam code and improve the performance of your pipeline.</li></ul>
<b>Activities</b>	Hands-on lab and quiz

**Module 10 State and Timers**

<b>Topics</b>	<ul style="list-style-type: none"><li>• State API</li><li>• Timer API</li><li>• Summary</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Identify use cases for state and timer API implementations.</li><li>• Select the right type of state and timers for your pipeline.</li></ul>
<b>Activities</b>	Quiz

---

**Module 11 Best Practices**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Schemas</li><li>• Handling unprocessable Data</li><li>• Error Handling</li><li>• AutoValue Code Generator</li><li>• JSON Data Handling</li><li>• Utilize DoFn Lifecycle</li><li>• Pipeline Optimizations</li></ul>
<b>Objectives</b>	Implement best practices for Dataflow pipelines.
<b>Activities</b>	Hands-on lab and quiz

---

**Module 12 Dataflow SQL and DataFrames**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Dataflow and Beam SQL</li><li>• Windowing in SQL</li><li>• Beam DataFrames</li></ul>
<b>Objectives</b>	Develop a Beam pipeline using SQL and DataFrames.
<b>Activities</b>	Hands-on lab and quiz

---

**Module 13 Beam Notebooks**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Beam Notebooks</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Prototype your pipeline in Python using Beam notebooks.</li><li>• Launch a job to Dataflow from a notebook.</li></ul>
<b>Activities</b>	Quiz

**Module 14** **Monitoring**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Job List</li><li>• Job Info</li><li>• Job Graph</li><li>• Job Metrics</li><li>• Metrics Explorer</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Navigate the Dataflow Job Details UI.</li><li>• Interpret Job Metrics charts to diagnose pipeline regressions.</li><li>• Set alerts on Dataflow jobs using Cloud Monitoring.</li></ul>
<b>Activities</b>	Quiz

---

**Module 15** **Logging and Error Reporting**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Logging</li><li>• Error Reporting</li></ul>
<b>Objectives</b>	Use the Dataflow logs and diagnostics widgets to troubleshoot pipeline issues.
<b>Activities</b>	Quiz

---

**Module 16** **Troubleshooting and Debug**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Troubleshooting Workflow</li><li>• Types of Troubles</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Use a structured approach to debug your Dataflow pipelines.</li><li>• Examine common causes for pipeline failures.</li></ul>
<b>Activities</b>	Hands-on lab and quiz

---

**Module 17** **Performance**

<b>Topics</b>	<ul style="list-style-type: none"><li>• Pipeline Design</li><li>• Data Shape</li><li>• Source, Sinks, and External Systems</li><li>• Shuffle and Streaming Engine</li></ul>
<b>Objectives</b>	<ul style="list-style-type: none"><li>• Understand performance considerations for pipelines.</li><li>• Consider how the shape of your data can affect pipeline performance.</li></ul>

---

Activities	Quiz
------------	------

---

**Module 18 Testing and CI/CD**

Topics	<ul style="list-style-type: none"><li>• Testing and CI/CD Overview</li><li>• Unit Testing</li><li>• Integration Testing</li><li>• Artifact Building</li><li>• Deployment</li></ul>
Objectives	<ul style="list-style-type: none"><li>• Testing approaches for your Dataflow pipeline.</li><li>• Review frameworks and features available to streamline your CI/CD workflow for Dataflow pipelines.</li></ul>
Activities	Hands-on labs and quiz

---

**Module 19 Reliability**

Topics	<ul style="list-style-type: none"><li>• Introduction to Reliability</li><li>• Monitoring</li><li>• Geolocation</li><li>• Disaster Recovery</li><li>• High Availability</li></ul>
Objectives	Implement reliability best practices for your Dataflow pipelines.
Activities	Quiz

---

**Module 20 Flex Templates**

Topics	<ul style="list-style-type: none"><li>• Classic Templates</li><li>• Flex Templates</li><li>• Using Flex Templates</li><li>• Google-provided Templates</li></ul>
Objectives	Using flex templates to standardize and reuse Dataflow pipeline code.
Activities	Hands-on labs and quiz

---

**Module 21 Summary**

Topics	Summary
--------	---------

**Objectives** Quick recap of training topics

**Activities** –