

# Data Engineering on Google Cloud v2.5.x

## Instructor Materials

---

*Last updated: March 25, 2022*

Data Engineering on Google Cloud 4-Day ILT courseware:

[Day 1: Modernizing Data Lakes and Data Warehouses with Google Cloud](#)

[Day 2: Building Batch Data Pipelines on Google Cloud](#)

[Day 3: Building Resilient Streaming Analytics Systems on Google Cloud](#)

[Day 4: Smart Analytics, Machine Learning and AI on Google Cloud](#)

[Timing Guide](#)

[Logging Issues](#)

[Changelog](#)

Day 1: Modernizing Data Lakes and Data Warehouses with Google Cloud					
#	Module	Trainer	Student	Lab	Demo
1.0	<a href="#">Introduction</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-
1.1	<a href="#">Introduction to Data Engineering</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Using BigQuery to do Analysis	<a href="#">Finding PII in your dataset with the DLP API</a>
1.2	<a href="#">Building a Data Lake</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Loading Taxi Data into Google Cloud SQL	-

1.3	<a href="#">Building a Data Warehouse</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	<p>Loading Data into BigQuery</p> <p>Working with JSON and Array Data in BigQuery</p>	<p><a href="#">Query TB+ of data in seconds</a></p> <p><a href="#">Exploring BigQuery Public Datasets with SQL using INFORMATION_SCHEMA</a></p> <p><a href="#">Nested and repeated fields in BigQuery</a></p>
1.4	<a href="#">Summary</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-

Day 2: Building Batch Data Pipelines on Google Cloud					
#	Module	Trainer	Student	Lab	Demo
2.0	<a href="#">Introduction</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-
2.1	<a href="#">Introduction to Building Batch Data Pipelines</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	<a href="#">ELT to improve data quality in BigQuery</a>
2.2	<a href="#">Executing Spark on Dataproc</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Running Apache Spark jobs on Dataproc	-
2.3	<a href="#">Serverless Data Processing with Dataflow</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	<p>Serverless Data Analysis with Dataflow: A Simple Dataflow Pipeline (Python/Java)</p> <p>Serverless Data Analysis with Dataflow: MapReduce in Dataflow (Python/Java)</p> <p>Serverless Data</p>	-

				Analysis with Dataflow: Side Inputs (Python/Java)	
2.4	<a href="#">Manage Data Pipelines with Cloud Data Fusion and Cloud Composer</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Building and Executing a Pipeline Graph in Cloud Data Fusion  [HOMEWORK IF NECESSARY] An Introduction to Cloud Composer	[OPTIONAL] <a href="#">Event-triggered Loading of data with Cloud Composer, Cloud Functions, Cloud Storage, and BigQuery</a>
2.5	<a href="#">Summary</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-

### Day 3: Building Resilient Streaming Analytics Systems on Google Cloud

#	Module	Trainer	Student	Lab	Demo
3.1	<a href="#">Introduction</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-
3.2	<a href="#">Serverless Messaging with Pub/Sub</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Streaming Data Processing: Publish Streaming Data into Pub/Sub	-
3.3	<a href="#">Dataflow Streaming Features</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Streaming Data Processing: Streaming Data Pipelines	-
3.4	<a href="#">High-Throughput BigQuery and Bigtable Streaming Features</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Streaming Data Processing: Streaming Analytics and Dashboards  Streaming Data Processing: Streaming Data Pipelines into	-

				Bigtable	
3.5	<a href="#">Advanced BigQuery Functionality and Performance</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Optimizing your BigQuery Queries for Performance  [OPTIONAL] Partitioned Tables in BigQuery	<a href="#">Mapping Fastest Growing Zip Codes with BigQuery GeoViz</a>
3.6	<a href="#">Summary</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-

Day 4: Smart Analytics, Machine Learning and AI on Google Cloud					
#	Module	Trainer	Student	Lab	Demo
4.0	<a href="#">Introduction</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-
4.1	<a href="#">Analytics and AI: Introduction</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-
4.2	<a href="#">Prebuilt ML Model APIs for Unstructured Data</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Using the Natural Language API to Classify Unstructured Text	-
4.3	<a href="#">Big Data Analytics with Notebooks</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	BigQuery in JupyterLab on Vertex AI	-
4.4	<a href="#">Production ML Pipelines with Kubeflow</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Running ML Pipelines on Kubeflow	-
4.5	<a href="#">Custom Model building with SQL in BigQuery ML</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	Predict Bike Trip Duration with a Regression Model in BigQuery ML  Movie Recommendations	<a href="#">Train a model with BigQuery ML to predict NYC taxi fares</a>

				in BigQuery ML	
4.6	<a href="#">Custom Model Building with AutoML</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-
4.7	<a href="#">Summary</a>	<a href="#">PDF</a>	<a href="#">PDF</a>	-	-

## Timing Guide

Access the recommended Timing Guide [here](#).

## Logging Issues

We are committed to ensuring a positive experience for both yourself and your students. You are encouraged to use the channels available to log issues as follows:

- **Customer Issues** allows you to log issues that compromise the learning experience. This can relate to functionality or content. Please use the **Priority** settings provided to allow us to prioritize remedial action.
  - [Internal CI link](#)
  - [External ATP CI link](#)
- **Feature Requests** allows you to log suggestions for the improvement of the overall learning experience.
  - [Internal FR link](#)
  - [External ATP FR link](#)

**NOTE:** When logging issues as an ATP, please do not delete the issue title prefix that is generated by the template. Kindly add a descriptor after the prefix.

## Changelog

Access the Changelog for this course [here](#).