
Data Engineering on Google Cloud Platform

Duration: 4 Days Course Code: GO5975

Overview:

This four-day instructor-led Google Cloud Platform class provides participants a hands-on introduction to designing and building data processing systems on Google Cloud Platform. Through a combination of presentations, demos, and hands-on labs, participants will learn how to design data processing systems, build end-to-end data pipelines, analyze data, and carry out machine learning. The course covers structured, unstructured, and streaming data.

Target Audience:

This class is intended for experienced developers who are responsible for managing big data transformations including: Extracting, Loading, Transforming, cleaning, and validating data Designing pipelines and architectures for data processing Creating and maintaining machine learning and statistical models Querying datasets, visualizing query results, and creating reports

Objectives:

- This course teaches participants the following skills:
 - Design and build data processing systems on Google Cloud Platform
 - Process batch and streaming data by implementing autoscaling data pipelines on Cloud Dataflow
 - Derive business insights from extremely large datasets using Google BigQuery
 - Train, evaluate, and predict using machine learning models using Tensorflow and Cloud ML
 - Leverage unstructured data using Spark and ML APIs on Cloud Dataproc
 - Enable instant insights from streaming data
-

Prerequisites:

To get the most out of this course, participants should have

- Completed Google Cloud Basics: Great Machine and Data Learning course OR have equivalent experience
 - Basic knowledge of the most common query language, such as SQL
 - Experience in data modeling, extraction, transformation, loading activities
 - Application development using a common programming language such as Python
- Familiarity with machine learning and/or statistics
-

Content:

| | | |
|--|--|---|
| Module 1: Introduction to Data Engineering | Optimizing with Partitioning and Clustering | Cloud Pub/Sub |
| Explore the role of a data engineer | Demo: Partitioned and Clustered Tables in BigQuery | ■ Lab: Publish Streaming Data into Pub/Sub |
| Analyze data engineering challenges | Preview: Transforming Batch and Streaming Data | Module 10: Cloud Dataflow Streaming Features |
| Intro to BigQuery | Module 4: Introduction to Building Batch Data Pipelines | Cloud Dataflow Streaming Features |
| Data Lakes and Data Warehouses | EL, ELT, ETL | ■ Lab: Streaming Data Pipelines |
| Demo: Federated Queries with BigQuery | Quality considerations | Module 11: High-Throughput BigQuery and Bigtable Streaming Features |
| Transactional Databases vs Data Warehouses | How to carry out operations in BigQuery | BigQuery Streaming Features |
| Website Demo: Finding PII in your dataset with DLP API | Demo: ELT to improve data quality in BigQuery | Lab: Streaming Analytics and Dashboards |
| Partner effectively with other data teams | Shortcomings | Cloud Bigtable |
| Manage data access and governance | ■ ETL to solve data quality issues | ■ Lab: Streaming Data Pipelines into Bigtable |
| Build production-ready pipelines | Module 5: Executing Spark on Cloud Dataproc | Module 12: Advanced BigQuery Functionality and Performance |
| Review GCP customer case study | The Hadoop ecosystem | Analytic Window Functions |
| Lab: Analyzing Data with BigQuery | Running Hadoop on Cloud Dataproc | Using With Clauses |
| Module 2: Building a Data Lake | GCS instead of HDFS | GIS Functions |
| Introduction to Data Lakes | Optimizing Dataproc | Demo: Mapping Fastest Growing Zip Codes with BigQuery GeoViz |
| Data Storage and ETL options on GCP | ■ Lab: Running Apache Spark jobs on Cloud Dataproc | Performance Considerations |
| Building a Data Lake using Cloud Storage | Module 6: Serverless Data Processing with Cloud Dataflow | Lab: Optimizing your BigQuery Queries for Performance |
| Optional Demo: Optimizing cost with Google Cloud Storage classes and Cloud Functions | Cloud Dataflow | ■ Optional Lab: Creating Date-Partitioned Tables in BigQuery |
| Securing Cloud Storage | Why customers value Dataflow | Module 13: Introduction to Analytics and AI |
| Storing All Sorts of Data Types | Dataflow Pipelines | What is AI? |
| Video Demo: Running federated queries on | | From Ad-hoc Data Analysis to Data Driven Decisions |
| | | ■ Options for ML models on GCP |

| | | |
|--|---|--|
| Parquet and ORC files in BigQuery | Lab: A Simple Dataflow Pipeline (Python/Java) | Module 14: Prebuilt ML model APIs for Unstructured Data |
| Cloud SQL as a relational Data Lake | Lab: MapReduce in Dataflow (Python/Java) | Unstructured Data is Hard |
| Lab: Loading Taxi Data into Cloud SQL | Lab: Side Inputs (Python/Java) | ML APIs for Enriching Data |
| Module 3: Building a Data Warehouse | Dataflow Templates | ■ Lab: Using the Natural Language API to Classify Unstructured Text |
| The modern data warehouse | ■ Dataflow SQL | Module 15: Big Data Analytics with Cloud AI Platform Notebooks |
| Intro to BigQuery | Module 7: Manage Data Pipelines with Cloud Data Fusion and Cloud Composer | What's a Notebook |
| Demo: Query TB+ of data in seconds | Building Batch Data Pipelines visually with Cloud Data Fusion | BigQuery Magic and Ties to Pandas |
| Getting Started | Components | ■ Lab: BigQuery in Jupyter Labs on AI Platform |
| Loading Data | UI Overview | Module 16: Production ML Pipelines with Kubeflow |
| Video Demo: Querying Cloud SQL from BigQuery | Building a Pipeline | Ways to do ML on GCP |
| Lab: Loading Data into BigQuery | Exploring Data using Wrangler | Kubeflow |
| Exploring Schemas | Lab: Building and executing a pipeline graph in Cloud Data Fusion | AI Hub |
| Demo: Exploring BigQuery Public Datasets with SQL using INFORMATION_SCHEMA | Orchestrating work between GCP services with Cloud Composer | ■ Lab: Running AI models on Kubeflow |
| Schema Design | Apache Airflow Environment | Module 17: Custom Model building with SQL in BigQuery ML |
| Nested and Repeated Fields | DAGs and Operators | BigQuery ML for Quick Model Building |
| Demo: Nested and repeated fields in BigQuery | Workflow Scheduling | Demo: Train a model with BigQuery ML to predict NYC taxi fares |
| Lab: Working with JSON and Array data in BigQuery | Optional Long Demo: Event-triggered Loading of data with Cloud Composer, Cloud Functions, Cloud Storage, and BigQuery | Supported Models |
| | Monitoring and Logging | Lab Option 1: Predict Bike Trip Duration with a Regression Model in BQML |
| | ■ Lab: An Introduction to Cloud Composer | ■ Lab Option 2: Movie Recommendations in BigQuery ML |
| | Module 8: Introduction to Processing Streaming Data | Module 18: Custom Model building with Cloud AutoML |
| | ■ Processing Streaming Data | Why Auto ML? |

Further Information:

For More information, or to book your course, please call us on 0800/84.009

info@globalknowledge.be

www.globalknowledge.com/en-be/