# Cognixia®

**COURSE CONTENT**

## COURSE NAME

**Data Engineering with Databricks**

## DURATION

**32 Hours**

## PREREQUISITES

- Beginner familiarity with basic cloud concepts (virtual machines, object storage, identity management)
- Ability to perform basic code development tasks (create compute, run code in notebooks, use basic notebook operations, import repos from git, etc.)
- Intermediate familiarity with basic SQL concepts (CREATE, SELECT, INSERT, UPDATE,
- DELETE, WHILE, GROUP BY, JOIN, etc.)
- Basic knowledge of Python programming, jupyter notebook interface, and PySpark fundamentals

## COURSE OUTLINE

### Module 1: Introduction to Data Engineering and Data Ingestion Strategies

- What is Data Engineering
- Role of a Data Engineer in modern data architecture
- Conceptual, Logical, and Physical Data Models
- OLTP vs. OLAP Systems
- ETL vs. ELT in Modern Data Architectures
- Data Ingestion Methods
- Batch Processing vs. Streaming Processing
- Extracting Data from APIs, Databases, Object Storage (S3, HDFS, ADLS, GCS)
- Working with Data Formats: CSV, JSON, XML, Avro, Parquet, ORC, Protobuf, Thrift
- Real-time Data Ingestion with Apache Kafka, AWS Kinesis, Google Pub/Sub
- Data Lake Architecture (Raw, Processed, Curated Layers)
- Incremental Load & Change Data Capture (CDC) Strategies
- Log-Based CDC
- Trigger-Based CDC
- Watermarking & Checkpointing in Streaming Pipelines
- Data Validation & Quality Checks
- Handling Corrupt Records in Large Datasets
- Schema Evolution & Enforcement in Avro & Parquet

### Module 2: Introduction to Apache Spark & PySpark

- What is Apache Spark? Why is it Used for Big Data Processing?
- Spark Core Components & Execution Flow

- RDDs (Resilient Distributed Datasets) vs. DataFrames vs. Datasets
- Transformations vs. Actions in Spark
- Understanding PySpark
- PySpark vs. Pandas vs. Dask
- Schema Inference & Explicit Schema Definition in PySpark
- Reading & Writing Data in PySpark (CSV, JSON, Avro, Parquet, ORC, Delta)
- Creating Spark DataFrames from Multiple Sources (relevant will be covered)
- RDDs, Lists, Databases, CSV, JSON, Avro, Parquet, ORC, Google BigQuery, Snowflake
- Essential DataFrame Operations

## Module 3: Databricks Lakehouse & Medallion Architecture

- Introduction to Databricks Lakehouse
  - Lakehouse vs. Data Lake vs. Data Warehouse
  - Key Features and Benefits

- Medallion Architecture
  - Bronze Layer (Raw Data Ingestion)
  - Silver Layer (Data Cleaning & Transformations)
  - Gold Layer (Aggregated Data for Analytics)

- Performance Optimization in Medallion Architecture
  - Best Practices for Data Organization
  - Incremental Data Processing for Each Layer

- Lab:
  - Ingest raw data into Bronze, clean it in Silver, and aggregate it in Gold.

## Module 4: Managing Data with Delta Lake

- Delta Lake Overview
  - What is Delta Lake?
  - Delta vs. Traditional Data Lakes
  - ACID Transactions & Schema Enforcement
  - Delta Lake in Medallion Architecture

- Schema Evolution & Enforcement
  - Schema-on-Read vs. Schema-on-Write
  - Handling Evolving Data Schemas
  - Auto Evolution vs. Explicit Schema Management

- Time Travel & Data Versioning
  - How to Query Previous Versions of Data
  - Implementing Rollback & Recovery

o Tracking Data Modifications for Audit Purposes

- Optimizing Delta Tables
  o Optimize & Compact Small Files
  o Auto Optimize, Vacuum & Retention Policies
  o Handling Large-Scale Data Growth with OPTIMIZE

- Handling Late Arriving Data
  o Strategies to Handle Late Events
  o Watermarking & Checkpointing in Streaming Pipelines
  o MERGE INTO for Handling Late Data Updates

- Lab:
  o Implement a Delta Lake-based Data Pipeline with Time Travel & Schema Evolution
  o Optimize queries using Z-Ordering and Auto Optimize

- Scalable Data Pipelines & Job Orchestration

## Module 5: Delta Live Table Basics

- Introduction to Delta Live Tables
  o What is Delta Live Tables (DLT)?
  o Why Use DLT for Data Pipelines?
  o DLT vs. Standard Data Pipelines

- Key Features of DLT
  o Continuous vs. Triggered Processing
  o Data Quality Enforcements & Expectations
  o Understanding the Event Log in DLT

## Module 6: Deploy Workloads with Databricks Workflows

- Introduction to Databricks Workflows
  o What are Workflows & Why Use Them?
  o Jobs vs. Notebooks Execution
  o Best Practices for Productionizing Data Pipelines

- Job Orchestration Strategies
  o Triggering Workflows
  o Managing Job Dependencies
  o Handling Job Failures & Retries

- Advanced Workflow Features
  o Multi-Task Workflows
  o Task Parameterization & Job Chaining

- o   Notifications & Monitoring in Workflows

- • Lab:
  - o   Automate an ELT Pipeline using Databricks Workflows
  - o   Schedule a Pipeline & Monitor Jobs

## Module 7: Data Governance with Unity Catalog

- • Introduction to Unity Catalog
  - o   Multi-Cloud Data Governance in Databricks
  - o   Key Benefits of Unity Catalog
  - o   Comparison with Other Governance Tools
- • Access Control & Permissions
  - o   Role-Based Access Control (RBAC)
  - o   Managing Users, Groups, and Roles
  - o   Row-Level Security (RLS) & Column-Level Security (CLS)

- • Data Lineage & Auditing
  - o   Tracking Data Lineage
  - o   Using Unity Catalog for Audit Logging
  - o   Compliance Considerations (GDPR, HIPAA, SOC 2)

- • Lab:
  - o   Configure Unity Catalog & Set Access Controls
  - o   Implement Row & Column Level Security Policies

## Module 8: Performance Tuning & Cost Optimization

- • Optimizing Spark Performance
  - o   Understanding the Catalyst Optimizer
  - o   Predicate Pushdown & Join Optimization
  - o   Caching Strategies for Performance Gains

- • Databricks Cluster Optimization
  - o   Choosing the Right Cluster for Workloads
  - o   Auto Scaling vs. Fixed Cluster Sizing
  - o   Optimizing Costs with Spot Instances

- • Optimizing Delta Lake Queries
  - o   Z-Ordering vs. Clustering
  - o   Auto Optimize & File Compaction Strategies
  - o   Impact of Large File Sizes on Query Performance

- • Lab:
  - o   Run performance benchmarks on optimized vs. non-optimized queries

o   Tune cluster configuration for cost and performance efficiency

**Module 9: Capstone Project - Real-World Implementation**

**Note: The Capstone project will take place post-training and falls outside the designated training hours.**

- End-to-End Data Pipeline Implementation
  o   Define Business Use Case & Data Requirements
  o   Design and Implement an ETL Pipeline with Delta Live Tables
  o   Use Workflows to Automate & Optimize the Pipeline
  o   Implement Security & Governance using Unity Catalog

- Final Review & Certification Preparation
  o   Mock Test for Databricks Data Engineer Associate Certification
  o   Best Practices for Databricks Deployment in Production

• Capstone Project:
  o   Build a Full Data Pipeline from S3 to Delta Lake to Gold Tables in Databricks
  o   Optimize, Secure, and Monitor the Entire Data Pipelin