## Course duration

- 4 days

## Course Benefits

- Learn to cloudera Manager features that make managing your clusters easier, such as aggregated logging, configuration management, resource management, reports, alerts, and service management.
- Learn to configuring and deploying production-scale clusters that provide key Hadoop-related services, including YARN, HDFS, Impala, Hive, Spark, Kudu, and Kafka.
- Learn to determining the correct hardware and infrastructure for your cluster.
- Learn to proper cluster configuration and deployment to integrate with the data center.
- Learn to ingesting, storing, and accessing data in HDFS, Kudu, and cloud object stores such as Amazon S3.
- Learn to how to load file-based and streaming data into the cluster using Kafka and Flume.
- Learn to configuring automatic resource management to ensure service-level agreements are met for multiple users of a cluster.
- Learn to best practices for preparing, tuning, and maintaining a production cluster.
- Learn to troubleshooting, diagnosing, and solving cluster issues.

## Course Outline

1. The Cloudera Enterprise Data Hub
    1. Cloudera Enterprise Data Hub
    2. CDH Overview
    3. Cloudera Manager Overview
    4. Hadoop Administrator Responsibilities
2. Installing Cloudera Manager and CDH
    1. Cluster Installation Overview
    2. Cloudera Manager Installation
    3. CDH Installation
    4. CDH Cluster Services
3. Configuring a Cloudera Cluster
    1. Overview
    2. Configuration Settings
    3. Modifying Service Configurations
    4. Configuration Files
    5. Managing Role Instances
    6. Adding New Services
    7. Adding and Removing Hosts

12. Cluster Maintenance
    1. Checking HDFS Status
    2. Copying Data Between Clusters
    3. Rebalancing Data in HDFS
    4. HDFS Directory Snapshots
    5. Upgrading a Cluster
13. Monitoring Clusters
    1. Cloudera Manager Monitoring Features
    2. Health Tests
    3. Events and Alerts
    4. Charts and Reports
    5. Monitoring Recommendations
14. Cluster Troubleshooting
    1. Overview
    2. Troubleshooting Tools
    3. Misconfiguration Examples
    4. Essential Points
15. Installing and Managing Hue
    1. Overview
    2. Managing and Configuring Hue
    3. Hue Authentication and Authorization
16. Security
    1. Hadoop Security Concepts
    2. Hadoop Authentication Using Kerberos
    3. Hadoop Authorization
    4. Hadoop Encryption
    5. Securing a Hadoop Cluster
17. Apache Kudu
    1. Kudu Overview
    2. Architecture
    3. Installation and Configuration
    4. Monitoring and Management Tools
18. Apache Kafka
    1. What Is Apache Kafka?
    2. Apache Kafka Overview
    3. Apache Kafka Cluster Architecture
    4. Apache Kafka Command Line Tools
    5. Using Kafka with Flume
19. Object Storage in the Cloud
    1. Object Storage
    2. Connecting Hadoop to Object Storag

## Class Materials

Each student will receive a comprehensive set of materials, including course notes and all the class examples.

Class Prerequisites

Experience in the following *is required* for this Hadoop class:

- Basic Linux experience.