# Hadoop for Administrators

Course Outline:

- **Introduction**

  o   Hadoop history, concepts

  o   Ecosystem

  o   Distributions

  o   High level architecture

  o   Hadoop myths

  o   Hadoop challenges (hardware / software)

  o   Labs: discuss your Big Data projects and problems

- **Planning and installation**

  o   Selecting software, Hadoop distributions

  o   Sizing the cluster, planning for growth

  o   Selecting hardware and network

  o   Rack topology

  o   Installation

  o   Multi-tenancy

  o   Directory structure, logs

  o   Benchmarking

  o   Labs: cluster install, run performance benchmarks

- **HDFS operations**

  o Concepts (horizontal scaling, replication, data locality, rack awareness)

  o Nodes and daemons (NameNode, Secondary NameNode, HA Standby NameNode, DataNode)

  o Health monitoring

  o Command-line and browser-based administration

  o Adding storage, replacing defective drives

  o Labs: getting familiar with HDFS command lines

- **Data ingestion**

  o Flume for logs and other data ingestion into HDFS

  o Sqoop for importing from SQL databases to HDFS, as well as exporting back to SQL

  o Hadoop data warehousing with Hive

  o Copying data between clusters (distcp)

  o Using S3 as complementary to HDFS

  o Data ingestion best practices and architectures

  o Labs: setting up and using Flume, the same for Sqoop

- **MapReduce operations and administration**

  o Parallel computing before mapreduce: compare HPC vs Hadoop administration

  o MapReduce cluster loads

  o Nodes and Daemons (JobTracker, TaskTracker)

  o MapReduce UI walk through

  o Mapreduce configuration

  o Job config

ASM Educational Center (ASM)

- o Optimizing MapReduce

- o Fool-proofing MR: what to tell your programmers

- o Labs: running MapReduce examples

- **YARN: new architecture and new capabilities**

  - o YARN design goals and implementation architecture

  - o New actors: ResourceManager, NodeManager, Application Master

  - o Installing YARN

  - o Job scheduling under YARN

  - o Labs: investigate job scheduling

- **Advanced topics**

  - o Hardware monitoring

  - o Cluster monitoring

  - o Adding and removing servers, upgrading Hadoop

  - o Backup, recovery and business continuity planning

  - o Oozie job workflows

  - o Hadoop high availability (HA)

  - o Hadoop Federation

  - o Securing your cluster with Kerberos

  - o Labs: set up monitoring

- **Optional tracks**

  - o Cloudera Manager for cluster administration, monitoring, and routine tasks; installation, use. In this track, all exercises and labs are performed within the Cloudera distribution environment (CDH5)

o   Ambari for cluster administration, monitoring, and routine tasks; installation, use. In this track, all

   exercises and labs are performed within the Ambari cluster manager and Hortonworks Data

   Platform (HDP 2.0)