

## Big Data Essentials Bootcamp

### Course Summary

#### Description

Big Data needs proper tools and skills, and this workshop brings you "from zero to hero," that is, provides the student with the necessary knowledge of Hadoop, Spark, and NoSQL. With these three fundamentals, you will be able to build systems processing massive amounts of data, in archival, batch, interactive and finally real-time manner. The workshop also lays foundations for proper analytics, allowing to extract insights from data.

#### Objectives

By the end of this course, students will be able to learn:

- Hadoop: HDFS, MapReduce, Pig, Hive
- Spark: Spark core, SparkSQL, Spark Java API, Spark Streaming
- NoSQL: Cassandra/HBase Architecture, Java API, Drivers, Data Modeling

#### Topics

- Hadoop
- Introduction to Hadoop
- SparkSpark Basics
- RDDs In Depth
- Partitions
- Spark API programming
- Introduction to Spark API / RDD API
- Spark Streaming
- NoSQL
- Cassandra Basics
- Cassandra Drivers
- Data Modeling – Part 1
- Data Modeling – Part 2
- Data Modeling Labs : Group Design Sessions

#### Audience

This course was designed for Developers.

#### Prerequisites

Before taking this course, students should be:

- Comfortable with Java programming language (most programming exercises are in Java)
- Comfortable in Linux environment (be able to navigate Linux command line, edit files using vi / nano)

#### Duration

Five days

## Big Data Essentials Bootcamp

### Course Outline

- I. Hadoop**
- II. Introduction to Hadoop**
  - A. Hadoop history, concepts
  - B. ecosystem
  - C. distributions
  - D. High-level architecture
  - E. Hadoop myths
  - F. Hadoop challenges
  - G. hardware / softwareHDFS Overview
  - H. concepts (horizontal scaling, replication, data locality, rack awareness)
  - I. architecture (Namenode, Secondary NameNode, DataNode)
  - J. data integrity
  - K. future of HDFS : Namenode HA, Federation
  - L. lab exercisesMapReduce Overview
  - M. MapReduceee concepts
  - N. phases : driver, mapper, shuffle/sort, reducer
  - O. thinking in MapReduce
  - P. future of mapreduce (yarn)
  - Q. lab exercisesPig
  - R. pig vs java vs MapReduce
  - S. pig latin language
  - T. user defined functions
  - U. understanding pig job flow
  - V. basic data analysis with Pig
  - W. complex data analysis with Pig
  - X. multi datasets with Pig
  - Y. advanced concepts
  - Z. lab exercisesHive
  - AA. hive concepts
  - BB. architecture
  - CC. data types
  - DD. Hive data management
  - EE. hive vs sql
  - FF. lab exercises
- III. SparkSpark Basics**
  - A. Background and history
  - B. Spark and hadoop
  - C. Spark concepts and architecture
  - D. Spark eco system (core, spark sql, mllib, streaming)
  - E. First look at Spark
  - F. Spark in local mode
  - G. Spark web UI
  - H. Spark shell
  - I. Analyzing dataset – part 1
  - J. Inspecting RDDs
- IV. RDDs In Depth**
  - A. Partitions
  - B. RDD Operations / transformations
  - C. RDD types
  - D. MapReduce on RDD
  - E. Caching and persistence
  - F. Sharing cached RDDs
- V. Spark API Programming**
  - A. Introduction to Spark API / RDD API
  - B. Submitting the first program to Spark
  - C. Debugging / logging
  - D. Configuration properties
- VI. Spark Streaming**
  - A. Streaming overview
  - B. Streaming operations
  - C. Sliding window operations
  - D. Writing spark streaming applications
- VII. NoSQL**
  - A. Introduction to Big Data / NoSQL
  - B. NoSQL overview
  - C. CAP theorem
  - D. When is NoSQL appropriate
  - E. NoSQL ecosystem
  - F. Cassandra Basics
  - G. Cassandra nodes, clusters, datacenters
  - H. Keyspaces, tables, rows and columns
  - I. Partitioning, replication, tokens
  - J. Quorum and consistency levels
  - K. Labs

## Big Data Essentials Bootcamp

### Course Outline (cont'd)

#### VIII. Cassandra Drivers

- A. Introduction to Java driver
- B. CRUD (Create / Read / Update, Delete) operations using Java client
- C. Asynchronous queries
- D. Labs

#### IX. Data Modeling – Part 1

- A. introduction to CQL
- B. CQL Datatypes
- C. creating keyspaces & tables
- D. Choosing columns and types
- E. Choosing primary keys
- F. Data layout for rows and columns
- G. Time to live (TTL), create, insert, update
- H. Querying with CQL
- I. CQL updates
- J. Labs

#### X. Data Modeling – Part 2

- A. Creating and using secondary indexes
- B. Denormalization and join avoidance
- C. composite keys (partition keys and clustering keys)
- D. Time series data
- E. Best practices for time series data
- F. Counters
- G. Lightweight transactions (LWT)

#### XI. Data Modeling Labs : Group Design Sessions

- A. Multiple use cases from various domains are presented
- B. Students work in groups to come up designs and models
- C. Discuss various designs, analyze decisions
- D. Lab : implement 'Netflix' data models, generate data