# HackerU

# Hadoop and Spark

## DM103

**32**
Academic Hours

# Hadoop and Spark

―――――

## Outline

The course delivers the key concepts and expertise participants need to ingest and process data on a Hadoop cluster using the most up-to-date tools and techniques. How to employ Hadoop ecosystem projects such as Spark, Hive, Flume, Sqoop, and Impala. Learning about the challenges faced by Hadoop developers. Participants learn to identify which tool is the right one to use in a given situation, and will gain hands-on experience in developing using those tools.
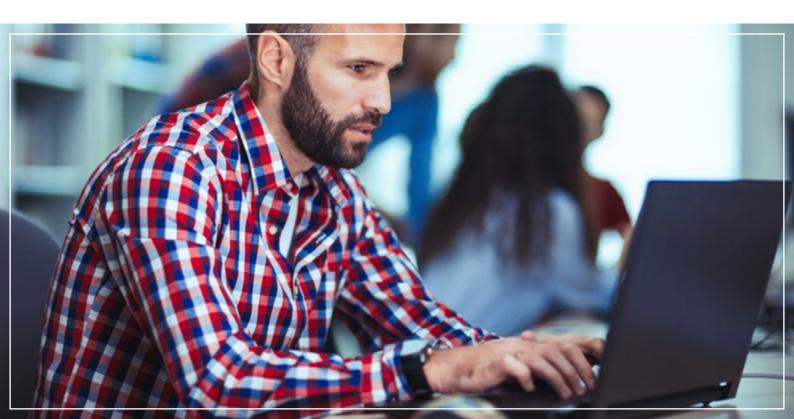
## Target Audience

I  Software developers

I  Software architects

I  Professionals willing to start working with Big Data

## Prerequisites

Basic knowledge of database concepts and development environments

# Content

"

Hadoop and Spark-
**Learning about
the challenges**
faced by Hadoop
developers"

# The HackerU
# Advantage

We have unparalleled experience in building advanced training programs for companies and organizations around the world — Talk to one of our experts and find out why.

## 01
**Handcrafted Training Programs**

## 02
**State-Of-The-Art Learning Materials**

## 03
**Israel's Premier Training Center**

## 04
**Fueled by Industry Leading Experts**

## 05
**Over 20 Years of Proven IT-Education Success**

info@hackerupro.com

www.hackerupro.com