

# Deep Dive Into Spark in 5 Days

Developer to Developer Training Series





## Course Outline

Apache Spark is the next generation successor to MapReduce. Spark is a powerful, open source processing engine for data in the Hadoop cluster, optimized for speed, ease of use, and sophisticated analytics. The Spark framework supports streaming data processing and complex, iterative algorithms, enabling applications to run up to 100x faster than traditional Hadoop MapReduce programs.

The 5 day Spark course is aimed at developers who are encountering Spark for the first time and want to understand how to build Big Data Products with Spark. The course would enable participants to build complete, unified Big Data applications combining batch, streaming, and interactive analytics on all their data.

Developers would be able to write sophisticated parallel applications to execute faster decisions, better decisions, and realtime actions, applied to a wide variety of use cases, architectures, and industries.

The course has a practical focus, mixing presentation with in depth hands on labs and exercises.



## Proposed Structure

Day 1	Day 2
<b>First brush</b>	<b>RDDs</b>
<ul style="list-style-type: none"><li>◆ Big Data Why and What?</li><li>◆ Introduction to Spark</li><li>◆ Spark Installation and Modes of Operation</li><li>◆ Spark shell</li></ul>	<ul style="list-style-type: none"><li>◆ RDD API In Detail</li><li>◆ Types of RDD (Pair RDD, Numeric RDD, JDBC RDD, KeyValue etc)</li><li>◆ Creating RDD From Different File Formats (Parquet, Avro, JSON, JDBC)</li></ul>
<b>Programming with Spark</b>	<b>Parallel Programming</b>
<ul style="list-style-type: none"><li>◆ Spark Fundamentals</li><li>◆ Role of Spark Context</li></ul>	<ul style="list-style-type: none"><li>◆ Partitions and Data Locality</li><li>◆ Executing parallel operations</li></ul>

<ul style="list-style-type: none"> <li>◆ Map Reduce in Spark</li> </ul>	
<b>RDD Fundamentals</b>	<b>Caching and Persistence</b>
<ul style="list-style-type: none"> <li>◆ Transformations in RDD</li> </ul>	<ul style="list-style-type: none"> <li>◆ Caching Overview</li> </ul>
<ul style="list-style-type: none"> <li>◆ Actions in RDD</li> </ul>	<ul style="list-style-type: none"> <li>◆ Distributed Persistence</li> </ul>
<b>Day 3</b>	<b>Advanced Concepts of RDD</b>
<b>Spark SQL</b>	<ul style="list-style-type: none"> <li>◆ Accumulators and Broadcast Variables</li> </ul>
<ul style="list-style-type: none"> <li>◆ Overview</li> </ul>	<ul style="list-style-type: none"> <li>◆ RDD Internals</li> </ul>
<ul style="list-style-type: none"> <li>◆ Role of SQLContext</li> </ul>	<ul style="list-style-type: none"> <li>◆ RDD Lineage</li> </ul>
<ul style="list-style-type: none"> <li>◆ Running Spark SQL in Spark shell</li> </ul>	
<b>DataFrames</b>	<b>Day 4</b>
<ul style="list-style-type: none"> <li>◆ Introduction to Data Frames</li> </ul>	<b>Spark Streaming</b>
<ul style="list-style-type: none"> <li>◆ Creating Data Frames</li> </ul>	<ul style="list-style-type: none"> <li>◆ Overview</li> </ul>
<ul style="list-style-type: none"> <li>◆ Transformations and Operations on Data Frames</li> </ul>	<ul style="list-style-type: none"> <li>◆ Role of StreamingContext</li> </ul>
<ul style="list-style-type: none"> <li>◆ Interoperating with RDDs</li> </ul>	<ul style="list-style-type: none"> <li>◆ Receivers</li> </ul>
	<ul style="list-style-type: none"> <li>◆ Streaming Applications</li> </ul>
<b>Datasets</b>	<b>DStreams</b>
<ul style="list-style-type: none"> <li>◆ Overview</li> </ul>	<ul style="list-style-type: none"> <li>◆ Introduction</li> </ul>
<ul style="list-style-type: none"> <li>◆ Creating Datasets</li> </ul>	<ul style="list-style-type: none"> <li>◆ Operations in DStreams</li> </ul>
<ul style="list-style-type: none"> <li>◆ Difference between Data Frames and Data Sets.</li> </ul>	<ul style="list-style-type: none"> <li>◆ Sliding Window Operations</li> </ul>
<ul style="list-style-type: none"> <li>◆ Conversion from Data Frame to Dataset and vice versa.</li> </ul>	<ul style="list-style-type: none"> <li>◆ Performance Tuning of DStreams</li> </ul>
	<ul style="list-style-type: none"> <li>◆ Stateful and Stateless Transformations in DStreams</li> </ul>

<b>Spark Schedulers</b>	<b>Spark MLlib</b>
◆ Overview	◆ Data Types
◆ Scheduling Across Applications	◆ Basic Statistics
◆ Scheduling Within Application	◆ Classification
	◆ Clustering
	◆ Pipelining
<b>Day 5</b>	
<b>Clustering</b>	
◆ Standalone	
◆ Configuration	
<b>Monitoring</b>	
◆ Web UI	
◆ REST API	
<b>Tuning and Debugging</b>	
◆ Data Serialization	
◆ Memory Management	
◆ Broadcasting Large Variables	
<b>Security</b>	
◆ Event Logging	
◆ Encryption	
◆ SSL Configuration	
◆ Standalone mode	
<b>Deployment</b>	
◆ Submitting Applications	
◆ Spark Standalone	
◆ Amazon EC2	
◆ Logging	
<b>Q&amp;A</b>	



## Course Prerequisites

To benefit from this course you should have programming experience with Scala or with Python. The language of instruction is Scala. Basic Linux knowledge is expected.

- ◆ For more information on the course or a discussion on your custom need, send a mail to [info@knoldus.com](mailto:info@knoldus.com)