

Course Title:	Big data framework & technologies
Course Code:	CSE523/1
Program:	Computer science Engineering
Department:	Computer Engineering
Course coordinator:	Dhekra CHERMITI
Institution:	Private Higher School of Engineers of Gafsa (ESIP)

#### A. Course Identification

1. Credit hours: 3 (1.5-1.5-0)		
2. Course type		
a. College Department Others		
b. Fundamental Transversal Optional		
<b>3.</b> Level/year at which this course is offered: 3.1/3		
4. Pre-requisites for this course (if any): CSE131, CSE323		

#### 1. Mode of Instruction (mark all that apply)

No	Mode of Instruction	Contact Hours	Self- study	Total workload
1	Traditional classroom			
2	Blended	45		
3	E-learning		33	78
4	Distance learning			
5	Other ()			

## 2. Contact Hours (based on academic semester)

No	Activity	<b>Contact Hours</b>
1	Lecture <b>Frivee de Gais</b>	30
2	Laboratory/Studio	15
3	Tutorial	-
4	Others (specify)	-
	Total	45



#### **B.** Course Objectives and Learning Outcomes

#### 1. Course Description

This course introduces students to Apache Spark and its core components. It covers the different versions of Spark (Scala, Python, Java), programming with Resilient Distributed Datasets (RDDs), working with Spark SQL, and handling graphs with GraphX. Students will also explore cluster architectures (Standalone, Apache Mesos, Hadoop YARN) and learn how to develop applications using SparkML and MLlib for machine learning.

#### 2. Course Main Objective

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

- ✓ Understand the fundamental concepts of Apache Spark.
- ✓ Explore Spark's different modules and compare them with the Apache Hadoop environment.
- ✓ Develop distributed applications using Apache Spark.
- ✓ Handle large datasets using SQL queries in Spark.
- $\checkmark$  Set up and configure a Spark cluster for real-world applications.
- ✓ Process real-time data streams using Spark Streaming.
- ✓ Manipulate and analyze graphs using the GraphX API.

#### 1. Course Learning Outcomes

	CLOs	Aligned PLOs
1	Knowledge and Understanding	
1.1	✓ Acquire the fundamental concepts of spark.	PLO.K1
1.2	✓ Explore the different modules of Spark and its relationship with the	PLO.K2
	Apache Hadoop environment.	
2	Skills	
2.2	✓ Handle datasets using SQL queries in Spark.	PLO.S2
	✓ Use the GraphX API to manipulate and analyze graph data.	
2.3	✓ Evaluate and analyze real-time data streams with Spark Streaming.	PLO.S3
2.6	✓ Develop distributed applications using Apache Spark.	PLO. S6

### C. Course Content Privée de Gafsa

No	List of Topics	<b>Contact Hours</b>
1	<ul> <li>Chapter 1: Introduction to Apache Spark &amp; Framework Overview</li> <li>1. What is Apache Spark?</li> <li>2. History of the Framework and its evolution.</li> <li>3. Versions of Spark: Scala, Python, and Java.</li> <li>4. Overview of Spark Modules: Spark Core, Spark SQL, Spark Streaming, MLlib, and GraphX.</li> </ul>	4



	<ol> <li>Comparison with Apache Hadoop: Differences in architecture, speed, and efficiency.</li> </ol>	
2	<ul> <li>Chapter 2: Working with Resilient Distributed Datasets (RDDs)</li> <li>1. Understanding RDDs: What they are and how they work.</li> <li>2. Creating, transforming, and reusing RDDs.</li> <li>3. Using Accumulators and Broadcast Variables for distributed computing.</li> </ul>	3
3	<ul> <li>Chapter 3: Processing Structured Data with Spark SQL</li> <li>1. Introduction to Spark SQL, DataFrames, and Datasets.</li> <li>2. Converting between RDDs and DataFrames.</li> <li>3. Connecting Spark to different data sources (JSON, Parquet, CSV, JDBC, NoSQL).</li> <li>4. Optimizing performance when querying large datasets.</li> </ul>	6
4	<ul> <li>Chapter 4: Deploying Spark Applications on a Cluster</li> <li>1. Setting up different cluster architectures: Standalone, Apache Mesos, Hadoop YARN.</li> <li>2. Using Spark-submit to deploy applications.</li> <li>3. Configuring cluster resources: memory, CPU, and workload balancing.</li> </ul>	6
5	<ol> <li>Chapter 5: Real-Time Data Processing with Spark Streaming</li> <li>How Spark Streaming processes real-time data.</li> <li>Understanding Discretized Streams (DStreams).</li> <li>Working with Kafka, Flume, HDFS, and TCP sockets for streaming data.</li> <li>Comparing Spark Streaming and Apache Storm for real-time analytics.</li> </ol>	5
Ecolo 6	<ul> <li>Chapter 6: Analyzing Graph Data with GraphX</li> <li>1. What is GraphX, and where is it used?</li> <li>2. Building graphs in Spark using Vertex and Edge RDDs.</li> <li>3. Applying graph algorithms like PageRank and Connected Components</li> </ul>	nieur
7	<ul> <li>Chapter 7: Introduction to Machine Learning with Spark MLlib</li> <li>1. Overview of machine learning in Spark.</li> <li>2. Understanding classification, regression, and clustering.</li> <li>3. Using MLlib for scalable machine learning with big data</li> </ul>	3
	Total	30



#### **C.1Practical work Content**

No	List of Topics	<b>Contact Hours</b>
1	Lab 1: Installing Apache Spark & Setting Up Environment	3
2	Lab 2: Programming with RDDs: Creating, Transforming & Using Broadcast Variables	3
3	Lab 3: Manipulating Structured Data with Spark SQL	3
4	Lab 4: Handling Graphs with GraphX: Building Graphs & Applying Graph Algorithms	3
5	Lab 5: Introduction to SparkML & MLlib	3
	Total	15

#### **D.** Teaching and Assessment

### 1. Alignment of Course Learning Outcomes with Teaching Strategies and Assessment Methods

Code	<b>Course Learning Outcomes</b>	Teaching Strategies	Assessment Methods
1.0	Knowledge and Understanding		
PLO.K1	<ul> <li>Acquire the fundamental concepts of spark.</li> </ul>		Assignments Quizzes
PLO.K2	<ul> <li>Explore the different modules of Spark and its relationship with the Apache Hadoop environment.</li> </ul>	Lecturing	, Exams,
2.0	Skills		
PLO.S2	<ul> <li>✓ Handle datasets using SQL queries in Spark.</li> <li>✓ Use the GraphX API to manipulate and analyze graph data.</li> </ul>		Assignments, Quizzes , Exams,
PLO.S3	<ul> <li>✓ Evaluate and analyze real-time data streams with Spark Streaming.</li> </ul>	Lecturing	aónioun
PLO. S6	<ul> <li>Develop distributed applications using Apache Spark.</li> </ul>		gemeur

#### 2. Assessment Tasks for Students

#	Assessment task*	Week Due	Percentage of Total Assessment Score
1	Practical Work (written or oral)	Weekly	25%
2	Quizzes, Homework assignments	-	-
3	First mid Term	8	25%



#	Assessment task*	Week Due	Percentage of Total Assessment Score
4	Final Exam	16	50%

#### E. Student Academic Counselling and Support

Arrangements for availability of faculty and teaching staff for individual student consultations and academic advice:

- Office hours \_
- Blackboard interface -

#### F. Learning Resources and Facilities

1. Learning Resources

Required Textbooks	<ol> <li>Jules S. Damji, Brooke Wenig, Tathagata Das, Denny Lee         <ul> <li>Learning Spark: Lightning-Fast Data Analytics, 2nd Edition, O'Reilly Media, 2020.</li> </ul> </li> <li>Holden Karau, Andy Konwinski, Patrick Wendell, Matei Zaharia – High-Performance Spark: Best Practices for Scaling and Optimizing Apache Spark, O'Reilly Media, 2017</li> <li>Ben Lorica, Mike Loukides – What You Need to Know About Apache Spark, 2015</li> </ol>	
Essential References Materials	Apache Spark	
Electronic Materials	<ol> <li>Coursera &amp; Udemy – Apache Spark &amp; Big Data Processing Courses</li> <li>MIT OpenCourseWare (OCW) – Big Data Analytics with Apache Spark</li> <li>GitHub Repositories</li> </ol>	
Other Learning Materials	NA	
2 Englistics Dequired	rieure d'Ingenieurs	

#### 2. Facilities Required

Item	Resources	
Accommodation	ut Galsa	
(Classrooms, laboratories, demonstration	classroom board software	
rooms/labs, etc.)		
Technology Resources	data show;	
(AV, data show, Smart Board, software, etc.)		

#### **G.** Course Quality Evaluation



<b>Evaluation Areas/Issues</b>	Evaluators	<b>Evaluation Methods</b>
Effectiveness of teaching and	Students, Faculty, Program Leaders, Peer	Direct/Indirect
assessment.	Reviewer	
Extent of achievement of	Faculty Program Landars Door Paviawar	Direct, Indirect
course learning outcomes.	Faculty, Flogram Leaders, Feel Reviewer	
Quality of Learning resources	Faculty, Program Leaders, Peer Reviewer	Direct, Indirect
Teaching and learning quality	Students, Faculty Program Leaders, Peer Reviewer Direct, Indirect	
and effectiveness.		

#### H. Specification Approval Data

Council / Committee	Computer Engineering Council
Date	11/09/2023

# Ecole Supérieure d'Ingénieurs Privée de Gafsa