

COURSE DESCRIPTION – ACADEMIC YEAR 2022/2023

Course title	Real-Time Big Data Processing
Course code	73033
Scientific sector	INF/01
Degree	Master in Computational Data Science (LM-18)
Semester	2
Year	1 and 2
Credits	6
Modular	No
Total lecturing hours	40
Total lab hours	20
Attendance	Generally, attendance is not compulsory, but non-attending students have to contact the lecturer at the start of the course to agree on the modalities of the independent study.
Prerequisites	<ul style="list-style-type: none"> • Programming (Java or Python – knowledge of one is enough) • Databases (SQL – basic knowledge) • Web development (HTML, JavaScript – basic knowledge)
Course page	https://ole.unibz.it/
Specific educational objectives	<p>The course belongs to the type "caratterizzanti – discipline informatiche" for the study path "no curriculum".</p> <p>The course aims at teaching both scientific foundations and practical aspects of real-time big data processing technologies. The students will learn the basic concepts of such systems and how to use them to solve concrete problems, including real-time data analyses, applications of machine learning and complex event processing over streaming data. Moreover, students will be trained to evaluate the advantages and disadvantages of such technologies in different application contexts.</p>
Lecturer	Francesco Corcoglioniti
Contact	Piazza Domenicani 3, francesco.corcoglioniti@unibz.it
Scientific sector of lecturer	ING-INF/05
Teaching language	English
Office hours	Arrange beforehand by email.
Lecturing assistant (if any)	--
Contact LA	--
Office hours LA	--
List of topics	<ul style="list-style-type: none"> • Reactive Streams • Messaging System (e.g., Apache Kafka) • Stateful Stream Processing (e.g., Apache Flink) • Micro-batch Stream Processing (e.g., Apache Spark) • Applications over Stream Processing • Semantic Technologies for streaming data

<p>Teaching format</p>	<p>Frontal lectures and hands-on labs (not evaluated).</p> <p>The lectures present the basic concepts, their realization in the open source systems studied in the course (e.g., Kafka, Flink, Spark), and their practical use with concrete examples.</p> <p>The labs permit students to practice the technologies of the course, by solving small tasks as part of guided tutorials often involving complete applications (from data ingestion to web front end) showcasing the use of these technologies.</p>
<p>Learning outcomes</p>	<p>Knowledge and understanding:</p> <ul style="list-style-type: none"> • D1.1 - Knowledge of the key concepts and technologies of data science disciplines • D1.3 - Knowledge of principles, methods and techniques for processing data in order to make them usable for practical purposes, and understanding of the challenges in this field • D1.4 - Sound basic knowledge of storing, querying and managing large amounts of data and the associated languages, tools and systems • D1.5 - Knowledge of principles and models for the representation, management and processing of complex and heterogeneous data <p>Applying knowledge and understanding:</p> <ul style="list-style-type: none"> • D2.1 - Practical application and evaluation of tools and techniques in the field of data science • D2.2 - Ability to address and solve a problem using scientific methods <p>Making judgments</p> <ul style="list-style-type: none"> • D3.2 - Ability to autonomously select the documentation (in the form of books, web, magazines, etc.) needed to keep up to date in a given sector <p>Communication skills</p> <ul style="list-style-type: none"> • D4.1 - Ability to use English at an advanced level with particular reference to disciplinary terminology • D4.3 - Ability to structure and draft scientific and technical documentation <p>Learning skills</p> <ul style="list-style-type: none"> • D5.1 - Ability to autonomously extend the knowledge acquired during the course of study • D5.2 - Ability to autonomously keep oneself up to date with the developments of the most important areas of data science • D5.3 - Ability to deal with problems in a systematic and creative way and to appropriate problem solving techniques.
<p>Assessment</p>	<p>The assessment of the course is based on a project which is done during the semester and requires students to solve a concrete problem by using methods and technologies taught in the course (100% of the mark).</p> <p>The project verifies whether the student is able to apply advanced data management techniques to solve concrete problems. The project is assessed through the submission of the solution source</p>

	<p>code and an accompanying written report, followed by an oral exam where the student defends the project with a short presentation including slides and live demo.</p> <p>The exam modalities are the same for attending and non-attending students.</p>
Assessment language	English
Assessment Typology	Monocratic
Evaluation criteria and criteria for awarding marks	<p>The final exam grade is the project mark (100%).</p> <p>Criteria for the evaluation of the project: correctness of the solution, complexity of the project, technologies used in the solution, quality of the report and the presentation.</p>
Required readings	<p>There is no single textbook that covers the entire course. The course material is collected from various textbooks, research papers and online documentation, including the following books:</p> <ul style="list-style-type: none"> • G. Shapira, T. Palino, R. Sivaram, K. Petty. "Kafka: The Definitive Guide". 2nd edition. O'Reilly Media, Inc. November 2021. • V. Kalavri, F. Hueske. "Stream Processing with Apache Flink". 1st edition. O'Reilly Media, Inc. April 2019. • M. Armbrust. "Learning Spark". 2nd edition. O'Reilly Media, Inc. July 2020. • M. Kleppmann: "Designing Data-Intensive Applications". 1st edition. O'Reilly Media, Inc. March 2017. <p>Subject Librarian: David Gebhardi, David.Gebhardi@unibz.it</p>
Supplementary readings	Additional sources will be announced during the course.
Software used	<p>Languages: Java or Python (students may also use Scala for their projects), SQL, HTML, JavaScript (or TypeScript)</p> <p>Software: Apache Kafka, Apache Flink, Apache Spark, ReactiveX, Docker and Docker Compose</p>