



# Big Data Storage & Processing

## MSc in Data Analytics - Sept 2023 cohort

### Module Introduction

**CCT College Dublin**  
**Ireland**

# Introduction



- **Lecturer:** Dr. Muhammad Iqbal\*
- **Experience:** Data Analytics, Numerical Modelling & Simulations, Structured & Object-Oriented Programming, Data Structures & Algorithms, Scalable Systems Programming (Python, R, Matlab, etc..).
- **E-mail:** [miqbal@cct.ie](mailto:miqbal@cct.ie)
- **Contact:** Use CCT email address for contact along with your Studentid, Module and Course name.

# Module Information



- **Contact hours:**
  - 2.5 hours lecture and tutorials
  - More than 7 hours weekly
- **Continuous Assessments**
  - 100% Continuous Assessment
- **Machine Requirements**
  - Windows 10/ 11 machine, At least 8 GB Ram and 512 GB hard disk, Core i5 or higher Micro processor

# Objectives



- The underlying concepts of Big Data Storage and Processing are mentioned below
  1. Fundamentals of Big Data storage and data management paradigms
  2. Underlying principles of parallel and distributed computing
  3. Current solutions for retrieving, integrating and processing of Big Data
  4. Big Data programming models and their efficient usage at scales
  5. Big Data Streams and their processing techniques

# Learning Outcomes



- **On successful completion of this module, the learner will be able to**
  1. Critically assess the **data storage and management requirements** of a given data project from a modern perspective and evaluate limitations of legacy approaches to Big Data. (Linked to PLO 3)
  2. Assess the design concepts and architectural patterns of distributed Big Data systems and analyse the components that form their technology stack. (Linked to PLO 1, PLO 2)
  3. Critically evaluate and select a Big data environment suitable for retrieving and processing a given Big Data set, perform data management and select appropriate analytic algorithms for the required scale and speed. (Linked to PLO 2, PLO 3)
  4. Assess the functional differences between common Big data environments and particularities of Big data and appropriate Graph Big data and processing stacks. (Linked to PLO 4)
  5. Implement the tools and technologies that facilitate the processing of Streaming Big Data to perform real-time analytics. (Linked to PLO 2)

# Module Contents



Content
Legacy Approaches <ul style="list-style-type: none"><li>• Traditional Computing Architecture &amp; Data Storage</li><li>• Relational DBMS(SQL) &amp; Data Silos</li><li>• Old Data (SQL) vs. Big Data</li></ul>
Distributed Systems and Data Management <ul style="list-style-type: none"><li>• Architectures</li><li>• Methodologies</li><li>• Scaling</li></ul>
Big Data Storage <ul style="list-style-type: none"><li>• Physical Storage</li><li>• Data Processing ETL/ELT</li><li>• Data Tiering</li><li>• File Formats, Compression and Security</li><li>• Disaster Recovery</li></ul>
No-SQL <ul style="list-style-type: none"><li>• Key-value (e.g., Couchbase, Redis)</li><li>• Document (e.g., MongoDB, CouchDB)</li><li>• Columnar (e.g., Big Table, Cassandra)</li><li>• Graph (e.g., Neo4j)</li><li>• Spatial (e.g., OGC-compliant)</li></ul>

Big Data Platforms <ul style="list-style-type: none"><li>• Apache Hadoop and HDFS</li><li>• MapReduce</li><li>• YARN (resource management)</li><li>• Apache Spark</li></ul>
Big Data Programming <ul style="list-style-type: none"><li>• Apache Hive (SQL-like queries)</li><li>• Apache Pig (high-level scripts that run on Apache Hadoop)</li><li>• Apache Mahout (machine learning algorithms on Apache Hadoop)</li><li>• Spark MLlib (scalable and easy machine learning library on Apache Spark)</li></ul>
Streaming Big Data <ul style="list-style-type: none"><li>• Spark Streaming</li><li>• Kafka</li></ul>
Graph Big Data <ul style="list-style-type: none"><li>• Apache Giraph (Graph processing on Graph Big Data)</li></ul>

# CCT Resources



- SupportHub (<https://moodle.cct.ie/course/view.php?id=91>)
- CCT ARC (<https://arc.cct.ie/>)
- For technical support, contact with Mr. Juan Murguey ([jmurguey@cct.ie](mailto:jmurguey@cct.ie))

**cct**

College Dublin  
Computing • IT • Business



# Questions?



# Books and eBooks



- Spark: The Definitive Guide Big Data Processing Made Simple, Bill Chambers, Matei Zaharia, O'Reilly Media, Inc., February 2018.
- Hadoop: The Definitive Guide, 4th Edition, Tom White, O'Reilly Media, Inc., April 2015.
- Big Data at Work: Dispelling the Myths, Uncovering the Opportunities, Thomas H. Davenport, Harvard Business Review Press, ISBN: 979-1422168165, 2014.
- Big Data: Concepts, Technology, and Architecture, Balamurugan Balusamy, Nandhini Abirami R, Seifedine Kadry, Amir H. Gandomi, Mar 2021.
- Data Engineering with Apache Spark, Delta Lake, and Lakehouse, Manoj Kukreja, Packt, ISBN: 9781801077743, 2021.