



**INTERNATIONAL QUALIFICATIONS  
AND ASSESSMENT CENTRE (IQAC)**



<b>Programme</b>	<b>Level 6 Diploma in Data Science</b>	
<b>Unit Number/ Unit Title</b>	UNIT 4 DATA ENGINEERING AND BIG DATA TECHNOLOGIES	
<b>Cohort Code:</b>	L06DEBD-U4	
<b>Unit Level</b>	Level 6	
<b>Total GLH</b>	Total qualification time 200/ Total Guided learning hours 90/ Self-guided learning hours 110	
<b>Credits</b>	20 CATS/ 10 ECTS	
<b>Lecturer</b>		
<b>Start Date</b>	<b>End Date</b>	

<b>Unit Aims</b>	This unit provides learners with practical knowledge of data engineering tools such as Apache Hadoop, Spark, Kafka, and modern ETL pipelines. Students will gain practical knowledge of essential data engineering tools, including Apache Hadoop, Spark, and Kafka, along with modern ETL (Extract, Transform, Load) pipelines. The unit will focus on hands-on experience in utilizing these technologies to manage and process large datasets efficiently. By the end of the course, students will be adept at implementing data workflows and optimizing data processing tasks, preparing them for real-world data engineering challenges.
<b>Differentiation Strategies</b> <i>(e.g. planned activities or support for individual learners according to their needs)</i>	The total number of students to be in the lesson is approximately 20. This is a multicultural group of students predominantly between the ages of 24 – 45, with numerous ethnic, gender, and creed background. These are UK academic level 5 students; hence it is assumed that they have practical, theoretical, or technological knowledge and understanding of a subject or field of work to find ways forward in broadly defined, complex contexts. These students must be able to generate information, evaluate, synthesise the use information from a variety of

	<p>sources. Various approaches to addressing the various identified students needs will be adopted throughout the lesson. Such will include:-</p> <ol style="list-style-type: none"> <li>1. Progressive tasks</li> <li>2. Digital resources</li> <li>3. Verbal support</li> <li>4. Variable outcomes</li> <li>5. Collaborative learning</li> <li>6. Ongoing assessment</li> <li>7. Flexible-pace learning</li> </ol>
<b>Equality &amp; Diversity</b>	Variety of teaching techniques will be employed to ensure that the needs of each individual learner are met.
<b>Safeguarding &amp; Prevent</b>	Safeguarding policies and the Prevent duty are strictly observed to ensure the safety, well-being, and inclusivity of all students and staff.
<b>Health &amp; Safety</b>	SIRM H&S policies will be maintained.
<b>Learning Resources</b>	<b>Teaching and Learning Materials</b>
	<ul style="list-style-type: none"> <li>• White, T. (2015). Hadoop: The Definitive Guide. O'Reilly Media.</li> <li>• Karau, H., &amp; Warren, R. (2017). High Performance Spark. O'Reilly Media.</li> <li>• Guller, M. (2020). Big Data Analytics with Spark. Apress.</li> </ul>

Learning Outcome	Assessment Criteria
<b>LO1.</b> 1. Understand the architecture of big data systems.	Written Report: <ul style="list-style-type: none"> <li>1.1 Explain batch vs stream processing.</li> <li>1.2 Describe the Hadoop ecosystem.</li> </ul>
<b>LO2.</b> 2. Use ETL tools to process large-scale datasets.	Programming Task: <ul style="list-style-type: none"> <li>2.1 Build ETL pipelines using Apache NiFi or Airflow.</li> <li>2.2 Transform data using Spark SQL or PySpark.</li> </ul>
<b>LO3.</b> 3. Apply streaming technologies for real time analytics.	Lab Work: <ul style="list-style-type: none"> <li>3.1 Use Apache Kafka or Flink.</li> <li>3.2 Monitor data streams for processing efficiency.</li> </ul>
<b>LO4.</b> 4. Ensure data quality and system reliability.	Portfolio: <ul style="list-style-type: none"> <li>4.1 Design validation rules for pipelines.</li> <li>4.2 Address fault tolerance in distributed environments.</li> </ul>

No	Learning Outcome / Topic	Learning and Teaching Activities	Which assessment criteria does the session relate to?	Day/month/year/ signature
1.	<b>Introduction to Big Data</b>	<b>Introduction to Big Data</b> 5 V's (Volume, Velocity, Variety, Veracity, Value)	LO1: Big Data System Architecture	
2.	<b>Batch vs. Stream Processing</b>	<b>Batch vs. Stream Processing</b> Lambda vs. Kappa architectures	LO1: Big Data System Architecture	
3.	<b>Hadoop Ecosystem</b>	<b>Hadoop Ecosystem</b> HDFS, YARN, MapReduce	LO1: Big Data System Architecture	
4.	<b>Modern Data Stack</b>	<b>Modern Data Stack</b> Data lakes (Delta Lake, Iceberg), lakehouses	LO1: Big Data System Architecture	
5.	<b>Cloud vs. On-Prem Solutions</b>	<b>Cloud vs. On-Prem Solutions</b> AWS EMR, Databricks, Cloudera	LO1: Big Data System Architecture	
6.	<b>ETL Fundamentals</b>	<b>ETL Fundamentals</b> Extract, Transform, Load workflows	LO2: ETL & Data Processing	
7.	<b>Batch ETL Tools</b>	<b>Batch ETL Tools</b> Apache Spark, AWS Glue, Talend	LO2: ETL & Data Processing	
8.	Half-Term Exam	<ul style="list-style-type: none"> <li>- Review of LO1 topics</li> <li>- Practice questions and mock assessment</li> <li>- <b>Half-term assessment</b> based on LO1 (theory)</li> </ul>	LO1 LO2	
9.	<b>Data Pipeline Orchestration</b>	<b>Data Pipeline Orchestration</b> Airflow, Prefect, Dagster	LO2: ETL & Data Processing	

10.	<b>Data Transformation Techniques</b>	<b>Data Transformation Techniques</b> SQL-based (dbt), PySpark	LO2: ETL & Data Processing	
11.	<b>Incremental Data Loading</b>	<b>Incremental Data Loading</b> CDC (Change Data Capture), SCD (Slowly Changing Dimensions)	LO2: ETL & Data Processing	
12.	<b>Stream Processing Basics</b>	<b>Stream Processing Basics</b> Event time vs. processing time, watermarks	LO3: Streaming & Real-Time Analytics	
13.	<b>Apache Kafka</b>	<b>Apache Kafka</b> Topics, partitions, producers/consumers	LO3: Streaming & Real-Time Analytics	
14.	Final Exam Preparation & Review	- Comprehensive review of all learning outcomes - Practice questions and revision of key topics		
15.	Final Exam	- <b>Final-term assessment</b> covering all learning outcomes (theory and practical elements)		
16.	Feedback & Reflection	- Review of final exam - Individual feedback on performance - Reflective discussion on key learning points		
17.	<b>Stream Processing Frameworks</b>	<b>Stream Processing Frameworks</b> Apache Flink, Spark Streaming	LO3: Streaming & Real-Time Analytics	
18.	<b>Real-Time Analytics</b>	<b>Real-Time Analytics</b> Windowed aggregations, stateful processing	LO3: Streaming & Real-Time Analytics	

19.	<b>Use Cases</b>	<b>Use Cases</b> Fraud detection, IoT monitoring	LO3: Streaming & Real-Time Analytics	
20.	<b>Data Quality Dimensions</b>	<b>Data Quality Dimensions</b> Accuracy, completeness, consistency	LO4: Data Quality & Reliability	
21.	<b>Data Validation Tools</b>	<b>Data Validation Tools</b> Great Expectations, Deequ	LO4: Data Quality & Reliability	
22.	<b>Monitoring &amp; Alerting</b>	<b>Monitoring &amp; Alerting</b> Prometheus, Grafana, custom dashboards	LO4: Data Quality & Reliability	
23.	Half-Term Exam	<b>Project</b> End-to-end data pipeline with quality checks		
24.	<b>Error Handling &amp; Recovery</b>	<b>Error Handling &amp; Recovery</b> Dead-letter queues, retry mechanisms	LO4: Data Quality & Reliability	
25.	<b>Data Governance</b>	<b>Data Governance</b> Metadata management, lineage tracking	LO4: Data Quality & Reliability	
26.	<b>Modern Data Platforms</b>	<b>Modern Data Platforms</b> Snowflake, BigQuery, Redshift	LO5: Capstone & Emerging Trends	
27.	<b>Data Mesh Principles</b>	<b>Data Mesh Principles</b> Domain-oriented ownership	LO5: Capstone & Emerging Trends	
28.	<b>ML Data Pipelines and Ethical Considerations</b>	<b>ML Data Pipelines</b> Feature stores (Feast, Tecton) <b>Ethical Considerations</b> Bias in big data, privacy preservation	LO5: Capstone & Emerging Trends	
29.	Final Exam Preparation & Review	LO1, LO2, LO3, LO4	LO1, LO2, LO3, LO4	
30.	Final Exam		LO1, LO2, LO3, LO4	