

**University of Djilali Bounaama - Khemis Miliana**  
**Faculty of Sciences of Matter and Computer Science**

**Second year of Master's degree**

**Full Name:** .....

**Exam in Big Data and Cloud Computing**

**Grade:**

**Duration: 1h30min**

.....

**University Year: 2024-2025**

**Exercise 1 (15pt):** Choose the correct answer. Each question has only one correct answer. Correct answer: +1, Incorrect answer: 0.

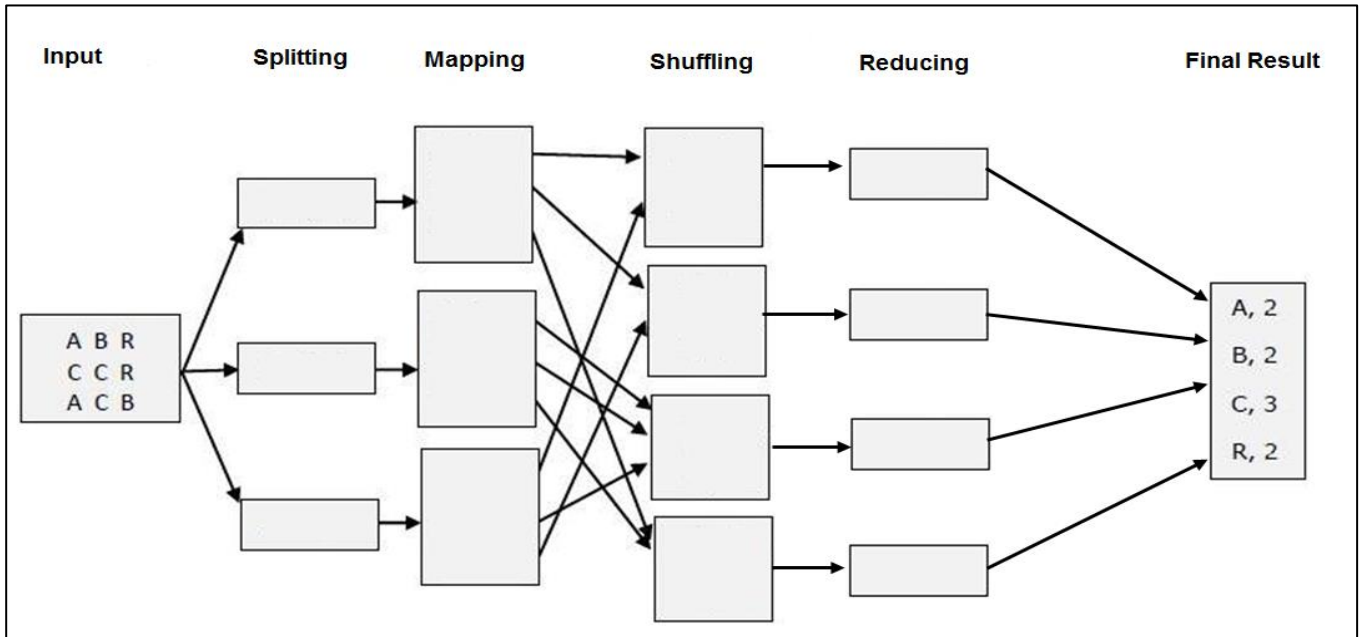
<p>1- Why is structuring Big Data important?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It increases the amount of raw data stored without any processing.</li> <li><input type="checkbox"/> It eliminates the need for data governance and security measures.</li> <li><input type="checkbox"/> It improves data management, enhances analysis, and ensures scalability.</li> </ul>	<p>5- What is the purpose of data cleansing in Big Data structuring?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> To encrypt data for better security.</li> <li><input type="checkbox"/> To remove duplicates, errors, and irrelevant information.</li> <li><input type="checkbox"/> To divide data into smaller chunks for parallel processing.</li> </ul>
<p>2- What is governance in Big Data?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Policies, procedures, and standards ensuring data is managed, protected, and used effectively and ethically.</li> <li><input type="checkbox"/> Techniques for compressing and storing large datasets without losing quality.</li> <li><input type="checkbox"/> A system for automating all data analysis tasks without human intervention.</li> </ul>	<p>6- What is a single-core processor?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> A CPU with one processing unit that executes tasks sequentially.</li> <li><input type="checkbox"/> A single-core processor can perform high-level data analysis and machine learning tasks.</li> <li><input type="checkbox"/> Many tasks can be performed by this type of processor.</li> </ul>
<p>3- What is the primary advantage of a data lake?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It stores large volumes of raw, unprocessed data in its native format.</li> <li><input type="checkbox"/> A data lake ensures all data is stored in a structured and organized format.</li> <li><input type="checkbox"/> A data lake always completes the unprocessed data imported from the data warehouse.</li> </ul>	<p>7- What is metadata in the context of a data warehouse?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Raw data stored in its original format without processing.</li> <li><input type="checkbox"/> Encrypted data used for securing sensitive information.</li> <li><input type="checkbox"/> Information about the data, such as schema, relationships, and lineage.</li> </ul>
<p>4- What is the role of the ETL process in Big Data?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> To analyze and visualize data directly from its raw format.</li> <li><input type="checkbox"/> To encrypt and store data without any processing.</li> <li><input type="checkbox"/> To extract, transform, and load data into a staging area for consistency and integration.</li> </ul>	<p>8- What is data wrangling?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Data wrangling is the process of encrypting data to enhance security.</li> <li><input type="checkbox"/> Data wrangling is the process of cleaning, organizing, and preparing data so it can be analyzed properly.</li> <li><input type="checkbox"/> Data wrangling is the process of storing raw data without any modifications.</li> </ul>

<p>9- How does isolation benefit transactions in ACID theorem?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It allows transactions to be executed simultaneously, even if they conflict.</li> <li><input type="checkbox"/> It is the process of retrieving information after an operating system failure.</li> <li><input type="checkbox"/> It ensures that transactions do not interfere with each other.</li> </ul>	<p>15- What is the architecture of HDFS?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> A peer-to-peer architecture where all nodes have equal responsibilities.</li> <li><input type="checkbox"/> A client-server architecture with a single central server handling all data storage.</li> <li><input type="checkbox"/> A master/slave architecture with a NameNode and DataNodes.</li> </ul>
<p>10- Which of the following data is structured data?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Tables with rows and columns.</li> <li><input type="checkbox"/> Social media posts and videos.</li> <li><input type="checkbox"/> JSON or XML files.</li> </ul>	<p>16- What type of data is stored in relational databases?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Structured data.</li> <li><input type="checkbox"/> Unstructured data like videos and images.</li> <li><input type="checkbox"/> Semi-structured data such as JSON files.</li> </ul>
<p>11- What is an example of data variety in Big Data?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Storing all data in a single, uniform format for easier processing.</li> <li><input type="checkbox"/> Combining structured data from databases with unstructured data like social media posts.</li> <li><input type="checkbox"/> Using only numerical data for analysis and ignoring text or images.</li> </ul>	<p>17- What is a cluster in Big Data storage?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> A cluster is a collection of servers (nodes) connected via a network to work as a single unit.</li> <li><input type="checkbox"/> A cluster is a single high-performance computer used for processing data.</li> <li><input type="checkbox"/> A cluster is a database used exclusively for storing structured data.</li> </ul>
<p>12- Why is velocity critical in Big Data?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It guarantees that all data is stored in a single location for better security.</li> <li><input type="checkbox"/> It focuses on reducing the size of data to save storage space.</li> <li><input type="checkbox"/> It ensures that fast-moving data is processed in real time for timely insights.</li> </ul>	<p>18- What resources does each node in a cluster have?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Each node has its own memory, processor, and hard drive.</li> <li><input type="checkbox"/> Each node shares all resources equally with other nodes.</li> <li><input type="checkbox"/> Each node has only a processor and relies on the network for memory and storage.</li> </ul>
<p>13- What is the main advantage of replication?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It provides fault tolerance and data availability.</li> <li><input type="checkbox"/> It reduces the overall storage requirements by storing data in a single location.</li> <li><input type="checkbox"/> It increases network traffic and reduces data access speed.</li> </ul>	<p>19- How does a cluster complete tasks efficiently?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> By using a single powerful server to handle all tasks.</li> <li><input type="checkbox"/> By dividing tasks into manageable chunks and distributing them among nodes.</li> <li><input type="checkbox"/> By storing all data in one location to reduce processing time.</li> </ul>
<p>14- Why is unstructured data harder to process?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Because it lacks a predefined structure.</li> <li><input type="checkbox"/> Because it is always encrypted and requires decryption before use.</li> <li><input type="checkbox"/> Because it is stored in a relational database.</li> </ul>	<p>20- What does Hadoop Common provide?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> A user interface for managing and monitoring Hadoop clusters.</li> <li><input type="checkbox"/> Tools and libraries for other Hadoop components.</li> <li><input type="checkbox"/> A distributed file system for storing large datasets across multiple machines.</li> </ul>

<p>21- What is Hadoop YARN?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It stores and processes large datasets across distributed clusters.</li> <li><input type="checkbox"/> It manages resources and coordinates the execution of applications across the system.</li> <li><input type="checkbox"/> It provides a user interface for monitoring and managing Hadoop jobs.</li> </ul>	<p>26- What are Data Marts in the context of data warehousing?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Data Marts are large, centralized repositories of all company data.</li> <li><input type="checkbox"/> Data Marts are the same as Data Lakes.</li> <li><input type="checkbox"/> Data Marts are smaller subsets of the data warehouse focused on specific business domains.</li> </ul>
<p>22- How are the write and read operations applied in the master-slave model?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Writing and reading are both done on the slave side to reduce master load.</li> <li><input type="checkbox"/> Both write and read operations are handled by the master, with no distinction between them.</li> <li><input type="checkbox"/> Writing is done at the master level, while the read operation is performed on the slave.</li> </ul>	<p>27- What does atomicity ensure in ACID theorem?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> That the transaction can be partially completed, with some operations applied and others not.</li> <li><input type="checkbox"/> That all operations in a transaction are applied in sequence, but without any guarantee of success.</li> <li><input type="checkbox"/> That all operations in a transaction are completed successfully, or none are applied.</li> </ul>
<p>23- What is the purpose of the Analysis Layer in a data warehouse?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> The analysis layer is used for storing backup copies of data.</li> <li><input type="checkbox"/> The analysis layer is where the data is collected from external sources.</li> <li><input type="checkbox"/> The analysis layer is where users interact with the data warehouse to extract insights.</li> </ul>	<p>28- What does the NameNode do?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> It determines the mapping of blocks to DataNodes.</li> <li><input type="checkbox"/> It stores and manages the actual data in the Hadoop cluster.</li> <li><input type="checkbox"/> It processes the data stored in the HDFS and runs MapReduce jobs.</li> </ul>
<p>24- How does horizontal scaling benefit Big Data systems?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> By adding more resources to increase capacity and performance.</li> <li><input type="checkbox"/> By upgrading the existing hardware to improve processing power.</li> <li><input type="checkbox"/> By reducing the number of nodes in the system to simplify management.</li> </ul>	<p>29- What is the role of DataNodes in HDFS?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> DataNodes manage the HDFS namespace and handle file system metadata.</li> <li><input type="checkbox"/> DataNodes coordinate communication between clients and the NameNode.</li> <li><input type="checkbox"/> They handle storage and serve read/write requests.</li> </ul>
<p>25- Can a system achieve all three aspects of the CAP theorem simultaneously?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> No, it can only achieve two out of the three at a time.</li> <li><input type="checkbox"/> Yes, it can achieve consistency, availability, and partition tolerance simultaneously.</li> <li><input type="checkbox"/> No, but in a business intelligence system, all aspects can be applied.</li> </ul>	<p>30- How is the relationship between the operating system and the file system?</p> <ul style="list-style-type: none"> <li><input type="checkbox"/> File systems are used by operating systems to store and retrieve data for programs.</li> <li><input type="checkbox"/> The operating system uses the file system only for temporary data storage.</li> <li><input type="checkbox"/> The file system controls the operating system's memory allocation.</li> </ul>

**Exercise 2 (2 pt):**

- The image above illustrates the stages of the **MapReduce** process to calculate the **frequency of each letter** in the original data (input). Each stage has a specific function in processing data.
- Your task is to **complete the missing information**.



**Exercise 3 (3 pt):**

Parallel processing is a method of computation where multiple tasks are executed simultaneously across multiple processors or cores. So:

Give me the **advantages** of parallel Data processing:

- .....
- .....
- .....

Give me the **disadvantages** Of Parallel Data Processing:

- .....
- .....
- .....