

# **Big Data Analytics**

## **Assignment-02 (Spring 2024)**

### **Submission to be done on portal**

**Due Date: Friday, May 3, 2024, before 11:59 PM**

- **Use Assignment pages and scan your assignment. (Hand written Only)**
- **Upload .pdf file on portal**
- **[Note: Title page must have Student Name, Registration Number, Section, and Date of Submission.**
- **Don't use register pages**
- **Submission will be on portal only**

### **Assignment 02**

Q #01:	<i>How does the implementation of erasure coding (EC) in HDFS, address the inherent storage overhead of replication while maintaining durability guarantees, and what are the key design considerations and optimizations made to seamlessly integrate EC into the existing distributed storage system of HDFS, encompassing modifications to the NameNode, DataNode, and client read and write paths, what future developments are anticipated regarding support for varied data layouts and advanced EC algorithms?</i>	5 Marks
--------	---	---------

Q #02:	<i>How do replication schemes like RAID-1 and erasure coding (EC), particularly Reed-Solomon (RS), differ in their approaches to tolerating disk failures, and how do their configurations impact data durability, storage efficiency, and storage overhead, considering factors such as the number of tolerated simultaneous failures and the ratio of data cells to parity cells? (Hint: You can use following table to write answer of question #02)</i>	5 Marks
--------	---	---------

	Data Durability	Storage Efficiency
Single replica	0	100%
Three-way replication	2	33%
XOR with six data cells	1	86%
RS(6,3)	3	67%
RS(10,4)	4	71%