

ENHANCING CLOUD DATA PIPELINES WITH DATABRICKS AND APACHE SPARK FOR OPTIMIZED PROCESSING

Rajkumar Kyadasu¹, Rahul Arulkumar², Krishna Kishor Tirupati³, Prof. (Dr) Sandeep Kumar⁴, Prof. (Dr) MSR Prasad⁵ & Prof. (Dr) Sangeet Vashishtha⁶

¹Rivier University, South Main Street Nashua, NH 03060

²University At Buffalo, New York, Srinagar Colony, Hyderabad, India

³International Institute of Information Technology Bangalore, India

⁴Department of Computer Science and Engineering Koneru Lakshmaiah Education Foundation Vadeshawaram, A.P., India

⁵Department of Computer Science and Engineering Koneru Lakshmaiah Education Foundation Vadeshawaram, A.P., India

⁶IIMT University, Meerut, India

ABSTRACT

The growing complexity of data ecosystems necessitates robust and scalable solutions to efficiently manage and process vast amounts of data in real-time. Cloud data pipelines, powered by advanced technologies like Databricks and Apache Spark, have emerged as key enablers for optimized data processing across industries. This paper explores how the integration of Databricks and Apache Spark enhances cloud data pipelines by enabling high-performance, distributed processing and real-time analytics. Databricks offers a collaborative environment that simplifies the orchestration of large-scale data workflows, while Apache Spark provides the core engine for executing data transformations with speed and scalability. By leveraging these platforms, organizations can improve the performance, cost-efficiency, and flexibility of their data pipelines, ensuring faster insights from their data. Additionally, this research examines best practices for deploying Databricks and Spark in cloud environments, highlighting their role in reducing operational complexity and optimizing resources for large-scale data operations. The study concludes with an analysis of real-world use cases demonstrating the effectiveness of this combination in various sectors, including finance, healthcare, and e-commerce.

KEYWORDS: *Cloud Data Pipelines, Databricks, Apache Spark, Distributed Processing, Real-Time Analytics, Data Transformation, Scalability, Optimized Data Workflows, Performance, Resource Optimization*

Article History

Received: 09 Feb 2020 | Revised: 11 Feb 2020 | Accepted: 12 Feb 2020
