

End-to-End IPL Data Engineering Pipeline using Apache Spark & Databricks

This document details the construction of a robust data engineering pipeline for Indian Premier League (IPL) cricket data. Leveraging the power of Apache Spark on Databricks, this project demonstrates a comprehensive approach from raw data ingestion to insightful visualization. Designed for aspiring data engineers and hiring managers, it showcases essential skills in distributed processing, data warehousing, and analytics, making it a resume-ready showcase of practical expertise.

Project Overview & Objectives

Our primary objective was to build an end-to-end data pipeline to process and analyze IPL cricket match data. This involved several key stages: ingesting raw CSV files, transforming them into a clean, queryable format, and then generating actionable insights. The project emphasizes best practices in data engineering, including schema management, data quality checks, and efficient distributed processing.

Ingestion

Automated loading of raw IPL match data.

Transformation

Cleaning, enriching, and structuring data for analysis.

Analysis

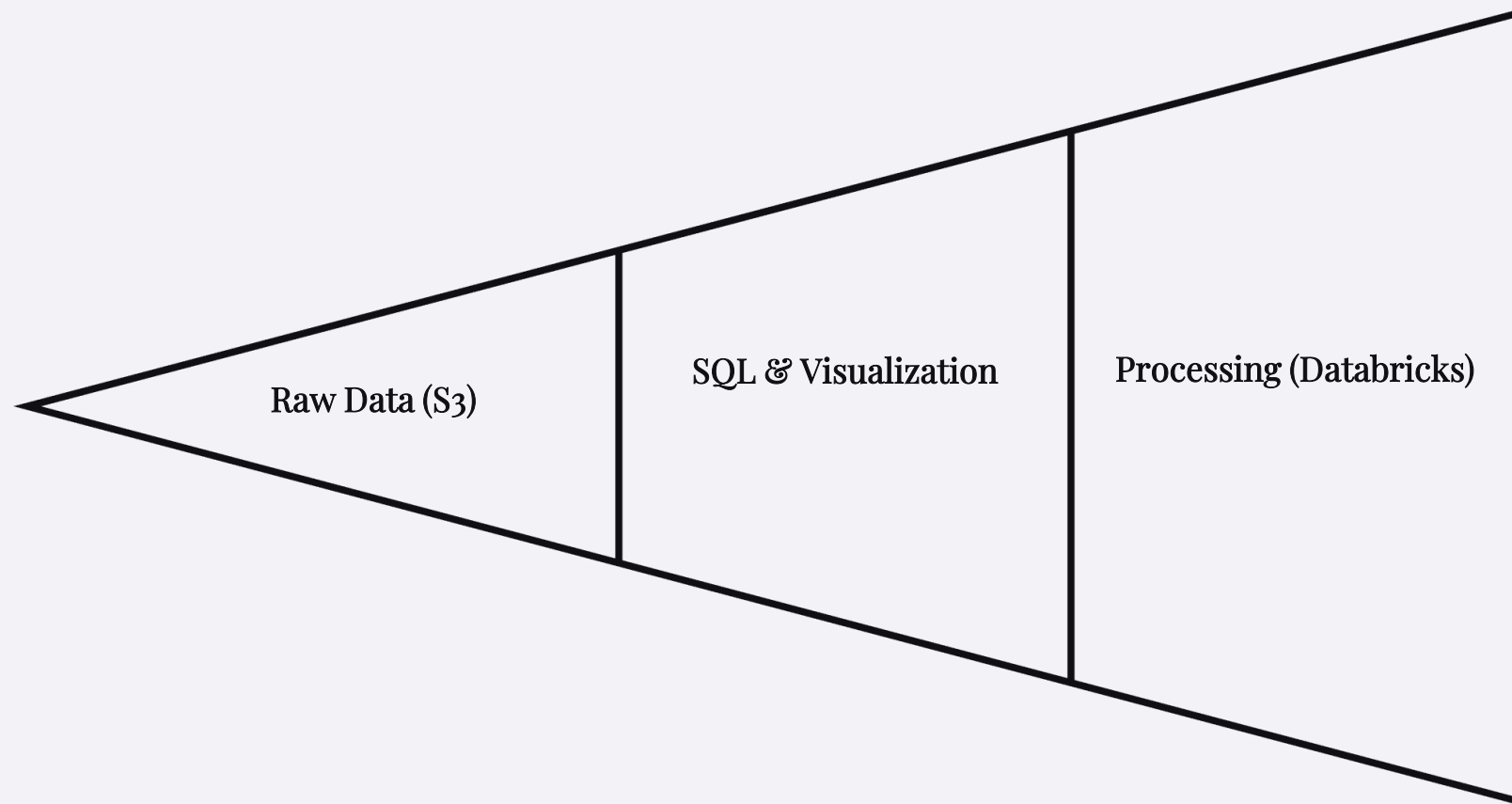
Deriving key statistics and performance metrics.

Visualization

Presenting insights through interactive dashboards.

Architecture Diagram: From Raw Data to Insights

The pipeline architecture is designed for scalability and efficiency, leveraging cloud-native services and distributed computing. Raw data resides in AWS S3, is processed by Apache Spark on Databricks, stored in a data warehouse (Delta Lake), and finally consumed by visualization tools.



This robust setup ensures data integrity and allows for rapid iteration on analytical requirements.



Dataset Description & Grain

The dataset comprises comprehensive IPL match information, including ball-by-ball details, player statistics, team performance, and match outcomes. Understanding the data grain is crucial for accurate analysis and preventing aggregation errors.

Core Data Sources

- Match-level details (venue, date, teams)
- Ball-by-ball events (batsman, bowler, runs, wickets)
- Player information (roles, teams)



Data Grain

The finest grain of our data is at the **ball-by-ball level**. This allows for granular analysis of individual deliveries, batting partnerships, and bowling spells, providing deep insights into match dynamics.

📋 **Grain definition:** The lowest level of detail in a dataset. Identifying the grain is essential before any aggregation or transformation.

Apache Spark Architecture: Distributed Processing

Apache Spark is central to our pipeline, offering powerful distributed processing capabilities. Its architecture, composed of a Driver and Executors within a Cluster, enables parallel computation across vast datasets.



Data Ingestion Process

Our data ingestion process is designed to be robust and efficient. Raw CSV files, representing individual IPL matches, are loaded from S3 into Databricks using Spark. We employed a structured approach to ensure data quality from the outset.



S3 Storage

Raw IPL CSV files stored in a designated S3 bucket.



Databricks Mount

S3 bucket mounted to Databricks for direct access.



Spark DataFrames

Raw data read into Spark DataFrames with defined schemas.



Delta Lake Write

Data written to Delta Lake tables for ACID properties.

Transformation Logic: Filters, Aggregations, Windows

Once ingested, the data undergoes a series of transformations to clean, enrich, and prepare it for analytical queries. These steps are crucial for creating a reliable "gold layer" of data.



Filters

Removing irrelevant records or handling missing values (e.g., incomplete match data).



Aggregations

Summarizing data (e.g., total runs scored per match, wickets taken per bowler).



Window Functions

Performing calculations over a defined set of rows (e.g., cumulative runs, moving average of scores).

These transformations ensure that the data is optimized for performance and ready for complex analytical queries.

SQL Analytics Use Cases

With our clean and transformed data in Delta Lake tables, we can perform advanced SQL analytics to uncover key performance indicators and patterns. Databricks SQL endpoints provide efficient access for these queries.



- **Top Player Performances**
Identify leading run-scorers, highest wicket-takers, and best strike rates.
- **Team Strategy Analysis**
Evaluate team strengths and weaknesses based on historical match data.
- **Match Impact Analysis**
Determine critical overs, partnerships, and moments that sway game outcomes.

Key Learnings & Outcomes

This project provided invaluable hands-on experience in building a production-ready data pipeline. From technical challenges to architectural decisions, several key learnings emerged, culminating in a robust and insightful analytics solution.

Technical Proficiency

- Mastered Spark DataFrame operations, including complex transformations.
- Gained expertise in Delta Lake for reliable data warehousing.
- Deepened understanding of Databricks ecosystem and best practices.

Project Summary

Developed an end-to-end data pipeline for IPL cricket data using Apache Spark on Databricks, ingesting raw S3 data, performing complex transformations, and enabling SQL analytics for visualization. Implemented schema enforcement and optimized Spark jobs for performance, showcasing expertise in distributed data processing and cloud data warehousing.

DATA ENGINEERING

APACHE SPARK

DATABRICKS

DELTA LAKE

AWS S3

SQL

Tools & Technologies Used

- **Cloud Platform:** Data Brick (Volume)
- **Compute & Orchestration:** Databricks
- **Distributed Processing:** Apache Spark
- **Data Warehousing:** Delta Lake
- **Query Language:** SQL
- **Programming Language:** Python
- **Version Control:** Git

