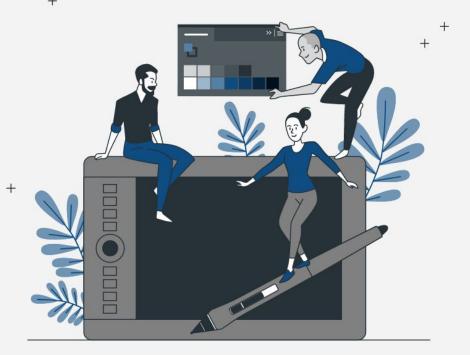
GOOD READS . * ETL PIPELINE

JUSTICE OHENE AMOFA



07 **ADAPTABILITY** Adaptability to other approaches and results. **TECHNICAL TOOLS AND** 08 Tools and mediums used for the project **TECHNOLOGY CLOSING STATEMENT** 09 **REFERENCES** 10

Project Overview



Goal: To Develop a real-time ETL pipeline for Goodreads data, enabling efficient data capture, storage, transformations, and analysis

Focus : Implementing an automated, scalable solution using Apache Spark, Apache Airflow, AWS S3 and Redshift



My Role

+ ----

Designing and Implementing the ETL pipeline.

04

+

Ensuring Data Quality and Integrity through validation checks.



02

Setting up and configuring AWS infrastructure (S3, Redshift, EMR)

05

Performing performance tuning and scalability testing.

03

+

Writing and scheduling ETL jobs in Apache Airflow.

06

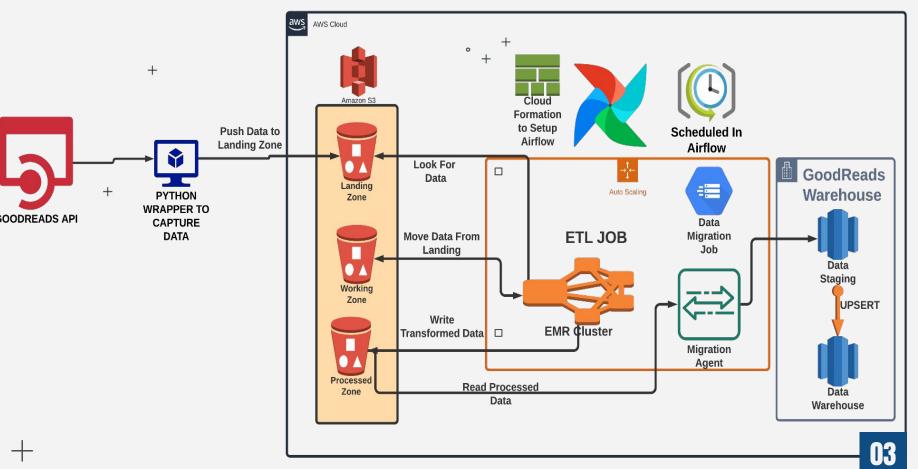
Identify and resolving technical issues within the ETL pipeline.







Concept of Solution



Technical Challenges

- Real Time Data Capture from the Goodreads API.
- Efficiently Managing Large Volumes of data in S3.
- Optimizing Sparks Jobs For performance
- Ensuring smooth data transfer between S3,
 EMR and Redshift.
- Implementing robust data quality checks and error handling.



Solution Approach



- Developed Modular ETL Jobs to handle each stage of data processing.
- Automated data transfer and transformation processes.
- Implemented data quality checks at multiple stages.
- Conducted Scalability testing using synthetic data.

Results and Impacts

Successful implementation of real-time ETL pipeline.



2

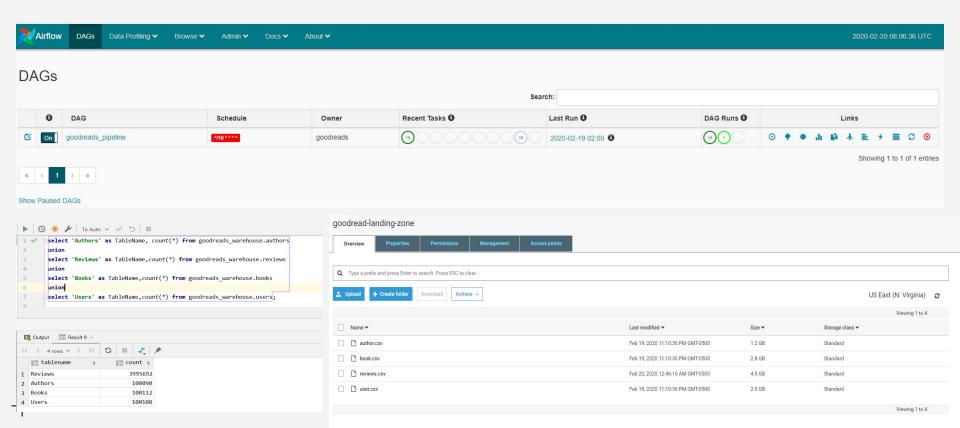
Efficient handling of large data volumes (up to 1.6 TB/day).

3

Enhanced data analysis with fast query performance in Redshift.



Results



Lesson Learned

- Importance of thorough planning for the ETL processes.
- Need for robust error handling and data quality checks.
- Value of automation in managing large-scale operations.
- Continuous performance tuning and monitoring are essential for maintaining efficiency.





















Goodreads API: Real time Data capture

AWS S3: Scalable and data storage

Apache Spark: Data Transformation and

processing.

Apache Airflow: ETL job scheduling and

Orchestration.

AWS Redshift: Data Warehousing and Analytics.

Python: Scripting and automation

Boto3: AWS SDK for python used for S3 operations.

Psycopg2: PostgreSQL adapter for python used for

Redshift.

Technical Tools And Requirements

Closing Statement

- Demonstrated ability to handle complex data engineering challenges.
- Proven expertise in using a variety of tools and technologies to build scalable, efficient ETL pipelines.
- Effective collaboration and adaptability in evolving project requirements.
- Commitment to ensuring data quality and integrity throughout the process.





Sponsors and Acknowledgements



Implemented by:

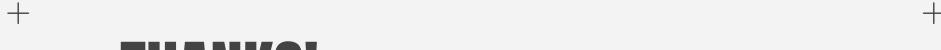








Special Thanks To Facilitators and Organizers: Derek Degbedzui ,Wonders Aggor and All other Staffs



THANKS!

Do you have any questions? Follow the project updates

<u>Joamofa@st.ug.edu.gh</u>

https://github.com/iamamofa/Final_project_Justice_goodreads_etl_pipeline

