# Pig Ingestion Using XML

R Prasannavenkatesh
*IMT2019063*
*prasannavenkatesh.ramkumar@iiitb.ac.in*

Channamsetty Sivani
*IMT2019020*
*channamsetty.sivani@iiitb.ac.in*

*Abstract*—The project focuses on creating a data ingestion pipeline that can handle XML data with a predefined schema. The pipeline begins by taking in the XML script, which is then parsed using Python. The parsed XML data is used to generate a basic Pig script, which is a high-level scripting language designed for large-scale data processing. The Pig script is then executed on a Hadoop cluster or other distributed computing framework to process and analyze the data. The output of the Pig script and the Pig script itself will be returned to the user, allowing them to easily view and manipulate the data.

*Index Terms*—xml, pig, react, fastapi, parser, abstraction

## I. INTRODUCTION

The ability to process and analyze large amounts of data has become increasingly important in today's world, with organizations of all sizes and across all industries relying on data-driven insights to make informed decisions. However, handling vast amounts of data can be a challenging and time-consuming task, particularly when dealing with data in a variety of formats.

This project aims to address these challenges by creating a data ingestion pipeline that can handle XML data with a predefined schema. The pipeline will use Python to parse the XML data and generate a Pig script, which is a high-level scripting language designed for large-scale data processing. The Pig script will then be executed on a distributed computing framework such as Hadoop, enabling the processing and analysis of large amounts of data in a scalable and efficient manner.

By automating the process of generating Pig scripts from XML data, this project aims to improve the efficiency and accuracy of large-scale data processing. The data ingestion pipeline and associated frontend and backend will allow users to quickly and easily process and analyze their data, saving time and reducing errors.

## II. PROBLEM DEFINITION

The problem statement of the project is to create a data ingestion pipeline that automates the process of generating Pig scripts from XML data, enabling users to process and analyze their data quickly and easily. The pipeline will use Python to parse the XML data and generate a Pig script, which will be executed on a distributed computing framework such as Hadoop.

## III. APPROACH

### A. XML file

In the XML file we write a pig script workflow with multiple processes defined. Each process corresponds to a specific operation or transformation performed on the input data. We first specify XML version and encoding used. Below are the tags definitions, like for purpose we are using a particular tag which helps user to write xml files.

- *<processes>* : The root element that contains all the process definitions.
- *<inputfile>*./yago_test.tsv*</inputfile>* : Specifies the input datafile path. Here yago_test.tsv is dataset and we are importing it from current directory.
- *<columns>*: Defines the column names and their types for the input data.
- *<col>*: Specifies a column name of data.
- *<type>*: Specifies the type of the corresponding column.
- *<delimiter>*space*</delimiter>*: Specifies the delimiter used in the input file.
- *<process>*: Represents a specific operation or transformation process.
- *<name>*: Specifies the name of the process.
- *<task>*: Specifies the task to be performed, such as filter, group, distinct, cross, join, limit, order, or foreach.
- *<column>*: Specifies the column used in the process.
- *<column1>*: Specifies the column from the second table used in the join operation or can use for another column.
- *<condition>*: Specifies the condition for the filter operation.
- *<clause>*: Specifies the value or condition for the task.
- *<table>*: Specifies the table or input data used in the process.
- *<table1>*: Specifies the second table or input data used in the process (for cross or join operations).
- *<variable>*: Specifies a variable name for the output of the process or can use for any other variable present.
- *<variable1>*: Specifies if any other variable is present.
- *<dumpvar>*: Specifies the variable to be dumped or stored in an output file.
- *<outputfile>*: Specifies the output file path for the dumped data.

Each process represents a step in the Pig script, and the commented sections at the end show examples of how to dump the results to an output file. After we parse this xml file to

parser, pigscript will be generated as shown in the following figure.



Fig. 1.   xml file



Fig. 2.   xml file



Fig. 3.   Generated pif script

## B. XML parser

The python script performs XML parsing and generates a Pig script based on the parsed data. The script imports the required packages, including os, json, and xml.etree.ElementTree for XML parsing. It defines a class named XMLParser that takes an XML file as input and initializes various attributes and data structures. The parseXML method is responsible for parsing the XML file and extracting relevant information such as the dataset, column definitions, and processing instructions. The writeScript method generates the Pig script based on the parsed data. It constructs script lines for loading the dataset, applying filter

conditions, performing grouping, distinct operations, cross products, joins, limits, ordering, and foreach operations. The generated Pig script is written to a file named "pig_script.pig." The executeScript method executes the generated Pig script using the os.system function. In the main section, an instance of the XMLParser class is created with the path to the XML file as an argument. The XML file is parsed, and the dataset, column definitions, and processes arPython script that performs XML parsing and generates a Pig script based on the parsed data.

This script provides a framework for parsing an XML file, extracting relevant information, and generating a Pig script based on the parsed data. It demonstrates an automated approach for processing XML data using Pig.e printed. Finally, the writeScript method is called to generate and execute the Pig script. It demonstrates an automated approach for processing XML data using Pig.



Fig. 4.   Xml parser



Fig. 5.   Xml parser

Fig. 6. Xml parser

## C. Platform

The project involves several platforms and technologies that work together to create an efficient and user-friendly system for processing large amounts of XML data with a predefined schema.

The data ingestion pipeline will be developed using Python, a popular programming language for data processing and analysis. Python will be used to parse the XML data and generate a Pig script, which is a high-level scripting language designed for large-scale data processing. The Pig script will be executed on a distributed computing framework such as Hadoop, allowing the processing of large amounts of data in a scalable and efficient manner.

To provide a user-friendly interface, the project will include a frontend and backend. The frontend will be developed using React, a popular JavaScript library for building web interfaces. React provides a flexible and powerful platform for building user interfaces, allowing the creation of intuitive and responsive interfaces that can handle large amounts of data. The frontend will allow users to upload their XML files and view the generated Pig scripts and output, making it easy to process and analyze data. The frontend in the initial state takes the following look.



Fig. 7. The initial state of the frontend

The frontend has a form object that takes in 2 corresponding form entities, namely the xml file and data file. Once the both

are taken, the submit button is pressed which then sends both of the files to the backend using the *ingest* route.



Fig. 8. The initial state of the frontend

The file storage, path generation, and uploads are handled in a cloud bucket as opposed to local for ease of development. We are using supabase, a cloud platform that supports file buckets. In supabase, a project is created which contains
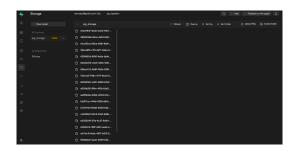


Fig. 9. The supabase bucket storage

The backend will be developed using FastAPI, a modern web framework for building APIs with Python. FastAPI provides a fast and efficient platform for building web applications, allowing for the rapid development of RESTful APIs that can handle large amounts of data. The backend will handle the communication between the frontend and the data ingestion pipeline, as well as the execution of the Pig scripts. The backend has an ingest route which will take in the file paths from the frontend, run the XML parser, generate the pig file, execute the pig file and return the results and the generated pig file back to the frontend.



Fig. 10. Backend Ingestion route

Fig. 11.   Frontend final state

Overall, the project involves a range of platforms and technologies that work together to create an efficient and user-friendly system for processing and analyzing large amounts of XML data with a predefined schema. By leveraging the strengths of each platform, the project aims to provide a powerful and flexible solution for data processing and analysis that can be easily used by organizations dealing with large amounts of data.

## IV. Evaluation and Results

The process of uploading input files from the frontend, parsing XML files, generating Pig scripts, and obtaining the output generally takes a comparable amount of time to running a regular Pig script. The additional steps involved in handling frontend interactions, XML parsing, and dynamic script generation are designed to be efficient and streamlined. Therefore, the overall execution time remains similar to the traditional Pig script execution. This allows for a seamless user experience without significant delays, ensuring that the workflow operates efficiently and delivers results in a timely manner. We also did this with large dataset successfully.

## V. Future

In this project, we utilized Supabase, an online database solution, which has a storage limit of 50MB. While Supabase served our needs for this project, it is essential to consider potential alternatives for future work or projects that may require larger storage capacity. Several options like Cloud-based Relational Databases, NoSql databases, serverless databases or self-hosted databases can be explored as potential replacements for Supabase based on specific requirements and considerations. And we can also do visualizations for ouput file.

## VI. Conclusion

In conclusion, the project aims to create a data ingestion pipeline for XML data with a predefined schema. The pipeline utilizes Python for parsing the XML data and generating Pig scripts, which are executed on a distributed computing framework. The goal is to enable efficient processing and analysis of large amounts of data in a scalable manner. By automating the process and providing an easy-to-use interface, the pipeline improves efficiency, accuracy, and accessibility for users working with XML data. The project addresses the challenges of handling diverse data formats and empowers organizations to gain valuable insights from their data quickly and easily.

## VII. Links

- Github : Nosql-Project

## References

[1] supabase, https://supabase.com/.
[2] fastapi, https://fastapi.tiangolo.com/lo/.
[3] React, https://react.dev/.
[4] XML, https://www.w3.org/standards/xml/core
[5] Apache Pig, https://pig.apache.org/
[6] Hadoop, https://hadoop.apache.org/