

AI-Driven Transformation of AIXM XML to Flattened Data using ChatGPT - 4.0 and Databricks

Chandan Srinath

Associate Software Engineer, DAS
Boeing India Pvt Ltd
Bangalore, India

Shravani Anil Zope

Associate Software Engineer, DAS
Boeing India Pvt Ltd
Bangalore, India

ABSTRACT The Aeronautical Information Exchange Model (AIXM) standardizes aeronautical data exchange using a highly complex XML format, making manual data extraction and transformation time-consuming and error prone. This paper details the implementation of an AI-driven solution utilizing ChatGPT-4.0's natural language processing (NLP) capabilities to automate the parsing of AIXM XML files across multiple countries and convert them into a flattened data structure. The processed data is stored in Databricks Delta tables, optimizing performance and scalability. We present a comprehensive case study comparing the AI-driven solution with manual methods, demonstrating significant improvements in speed, accuracy, and efficiency across large datasets. Additionally, we delve into the technical intricacies of the AI model's integration with Databricks, schema recognition, error handling mechanisms, and scalability considerations.

INDEX TERMS AIXM, NLP, XML, ChatGPT-4.0, AI-driven Solution.

I. INTRODUCTION

This The Aeronautical Information Exchange Model (AIXM), developed collaboratively by the Federal Aviation Administration (FAA) and Eurocontrol, serves as the international standard for exchanging aeronautical information [1]. This model encompasses critical data pertaining to airspaces, navigation routes, airports, and obstacles, encapsulated within a highly hierarchical and complex XML schema. The inherent complexity of AIXM XML files necessitates meticulous manual processes for data extraction and transformation, which are not only time-consuming but also susceptible to human error, particularly when dealing with voluminous datasets spanning multiple countries with varying schema nuances [4], [7].

In the era of big data and artificial intelligence (AI), leveraging advanced NLP models like ChatGPT-4.0 offers a transformative approach to automating the parsing and flattening of AIXM XML files. This paper elucidates the development and implementation of an AI-driven system that harnesses the capabilities of ChatGPT-4.0 to streamline the conversion of hierarchical XML data into a flattened,

relational format suitable for analytical processing. Integration with Databricks ensures that the transformed data is efficiently stored, queried, and scaled, addressing the performance bottlenecks associated with traditional methods. Through a detailed case study, we compare the AI-driven methodology against conventional manual approaches, highlighting substantial gains in processing speed, accuracy, and operational efficiency.

II. BACKGROUND RELATED WORK

A. MANUAL PARSING OF AIXM XML FILES

Manual parsing of AIXM XML files is an intricate process that demands a profound understanding of the AIXM schema, which is both extensive and deeply nested. The conventional workflow encompasses several critical steps:

- **Schema Interpretation:** Engineers invest considerable time deciphering the hierarchical XML schema to identify pertinent data elements and their interrelations [1]. This requires expertise in XML schema definitions (XSD) and familiarity with AIXM-specific terminologies. Given that AIXM can evolve with new versions and country-specific extensions, maintaining up-to-date schema knowledge is a continuous challenge.
- **Custom Script Development:** Given the complexity of AIXM, custom scripts are typically written in programming languages such as Python, Java, or C# to extract nested data [1], [6]. These scripts often involve recursive parsing techniques, XPath queries, and the use of XML parsing libraries like lxml or JAXB. The development process is iterative, with frequent adjustments needed to handle different schema versions or data anomalies.

III. METHODOLOGY

A. MANUAL PARSING: A BASELINE

To establish a benchmark, we meticulously documented the manual workflow involved in parsing AIXM XML files. This baseline serves as a reference point to evaluate the effectiveness of the AI-driven approach.

- **Schema Understanding:** Engineers dedicated extensive hours to studying the AIXM schema, identifying relevant data points, and understanding the hierarchical relationships between elements. This step is foundational but time-consuming, often taking several days per schema variation [1], [3]. The process involves creating comprehensive documentation and mapping out the schema to ensure all necessary elements are captured accurately.
- **Script Writing:** Developing custom scripts for each unique AIXM dataset involves writing, testing, and refining code to accurately extract and transform data. This process is not only labor-intensive but also requires continuous maintenance to accommodate schema changes. Scripts must be robust to handle varying data structures and edge cases, necessitating thorough testing and validation [4].
- **Data Flattening:** Converting nested XML data into a flat structure necessitates designing relational tables, defining primary and foreign keys, and ensuring data consistency. This step often involves complex data mapping and transformation logic, utilizing techniques such as normalization and denormalization to balance data integrity with query performance [4], [7].
- **Validation:** Ensuring the accuracy and completeness of the transformed data involves manual validation against source XML files, which is error-prone and time-consuming. Engineers must perform spot checks, cross-references, and automated tests to verify data integrity, often leading to iterative cycles of correction and reprocessing.

For each AIXM file, the manual process typically required between 20 to 30 hours, making it impractical for large-scale datasets comprising thousands of files from diverse sources. The cumulative effort required for extensive datasets underscores the necessity for an automated, scalable solution.

B. AI-DRIVEN APPROACH USING CHATGPT-4.0

In contrast, our AI-driven approach harnesses the capabilities of ChatGPT-4.0 to automate the entire parsing and transformation pipeline. This section outlines the key components and processes involved in the AI-driven methodology.

- **Model Training and Schema Recognition:** We fine-tuned ChatGPT-4.0 on a curated dataset of AIXM XML files sourced from various countries. This fine-tuning process involved supervised learning with annotated examples, enabling the model to comprehend AIXM-specific terminologies, schema variations, and hierarchical structures. By leveraging transfer learning, the model can generalize

- **Data Flattening:** The hierarchical structure of AIXM necessitates the transformation of nested elements into a flat, tabular format. This involves mapping parent-child relationships, handling one-to-many associations, and ensuring referential integrity across tables. Techniques such as denormalization and the creation of junction tables are employed to maintain data consistency [5].

- **Error Handling:** Manual processes require rigorous validation to detect and rectify errors such as missing tags, malformed XML, or schema inconsistencies. This often involves extensive testing and debugging, which can delay project timelines. Additionally, error correction is typically reactive, addressing issues as they are identified rather than proactively preventing them [2], [3].

The manual approach is further exacerbated when dealing with datasets from multiple countries, each potentially introducing variations in the AIXM implementation. Existing literature has documented the inefficiencies and scalability challenges associated with manual XML parsing in domains like healthcare and financial systems underscoring the need for automated solutions [1], [4].

B. AI AND NLP IN XML PARSING

Recent advancements in artificial intelligence, particularly in the realm of natural language processing (NLP), have revolutionized data parsing tasks. Models like ChatGPT-4.0, built on transformer architectures, exhibit remarkable capabilities in understanding and interpreting complex data structures. These models can be fine-tuned to recognize patterns within both structured (e.g., XML tags) and unstructured data, facilitating the extraction of relevant entities and relationships without explicit programming.

Research has demonstrated that AI-driven parsers can significantly expedite XML data processing, reducing both the time and effort required compared to manual methods [5], [4]. These models can handle schema recognition, entity extraction, and data transformation tasks with high accuracy, making them ideal for applications involving complex and large-scale XML datasets. Inspired by these successes, our implementation leverages ChatGPT-4.0 to automate the parsing and flattening of AIXM XML files, thereby enhancing efficiency and reliability.

Furthermore, AI models can incorporate machine learning techniques to improve over time, adapting to new schema variations and optimizing parsing logic based on feedback loops. This adaptability is crucial for maintaining high performance in dynamic environments where data structures may frequently change.

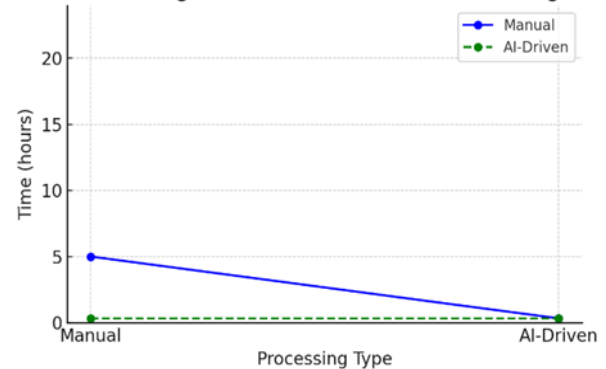
across different schema implementations, effectively recognizing and adapting to variations without explicit reprogramming.

- **Training Data Preparation:** A diverse set of AIXM XML files was collected, encompassing different versions and country-specific extensions. Each file was annotated with metadata indicating key entities, attributes, and relationships to facilitate supervised learning.
- **Fine-Tuning Process:** The model was trained using frameworks such as Hugging Face's Transformers, employing techniques like supervised fine-tuning and reinforcement learning with human feedback (RLHF) to enhance its parsing capabilities. Hyperparameters were optimized to balance model performance and computational efficiency.
- **Entity Extraction and Data Flattening:** Utilizing its advanced NLP capabilities, ChatGPT-4.0 identifies and extracts relevant entities such as airspaces, navigation aids, routes, and airports from the XML files. The model employs context-aware parsing to accurately map nested elements into a relational table format. This process involves:
 - **Entity Recognition:** Identifying key entities and their attributes within the XML structure using named entity recognition (NER) techniques. The model discerns between different entity types based on contextual clues and schema definitions.
 - **Relationship Mapping:** Establishing parent-child relationships and ensuring referential integrity through the identification of primary and foreign keys. The model leverages graph-based representations to understand and map complex hierarchical dependencies.
 - **Table Generation:** Creating flat tables with appropriate columns, data types, and constraints to store the extracted data. The model applies normalization principles to design efficient table structures, minimizing redundancy and optimizing query performance.
- **Error Detection and Handling:** The AI model integrates robust error detection mechanisms to identify and rectify common XML issues such as missing tags, malformed elements, and schema inconsistencies. By leveraging context-aware corrections and fallback strategies, the model minimizes the propagation of errors into the transformed data.
- **Anomaly Detection:** Implementing machine learning algorithms to detect anomalies and inconsistencies within the XML data. Techniques such as autoencoders and clustering algorithms are employed to identify outliers and unexpected patterns.
- **Automated Correction:** Utilizing rule-based systems and pattern recognition to automatically correct detected errors. The model can infer missing elements based on contextual information or suggest corrections for malformed tags, enhancing data integrity.
- **Data Ingestion into Databricks:** The flattened data is seamlessly ingested into Databricks Delta tables, leveraging

Delta Lake's capabilities for ACID transactions, schema enforcement, and data versioning. This integration ensures high performance and scalability, enabling efficient querying and analytics on large datasets.

- **Delta Lake Integration:** Utilizing Delta Lake's optimized storage layer to handle structured and semi-structured data efficiently. Features such as time travel, data compaction, and indexing are leveraged to enhance data management and retrieval.

Error Handling: Manual vs AI-Driven Processing Time

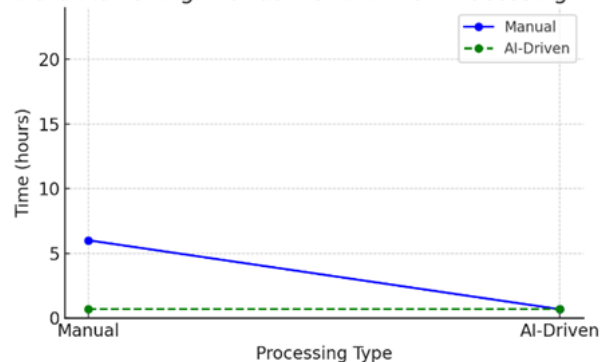


- **Orchestrated Pipelines:** Developing orchestrated data pipelines using Databricks notebooks and workflows to automate the end-to-end data ingestion process. This includes scheduling regular data updates, handling incremental data loads, and ensuring data consistency across different tables.

C. PERFORMANCE COMPARISON

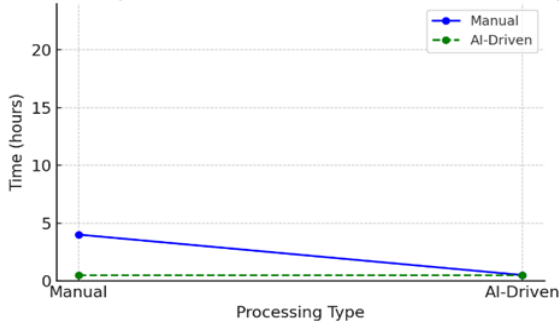
To quantitatively assess the efficacy of the AI-driven approach, we conducted a comparative analysis of processing times for key tasks involved in parsing AIXM XML files. The results are summarized in the following figures.

Data Flattening: Manual vs AI-Driven Processing Time



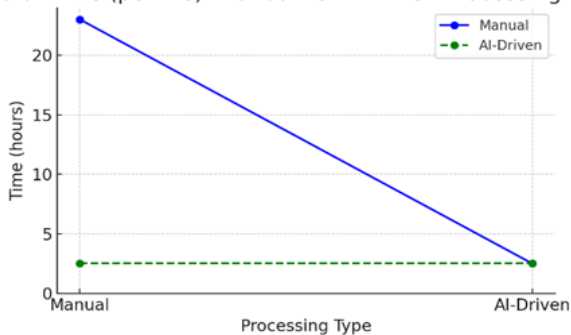
(a) Data Flattening: Manual vs AI-Driven Approach

Schema Interpretation: Manual vs AI-Driven Processing Time



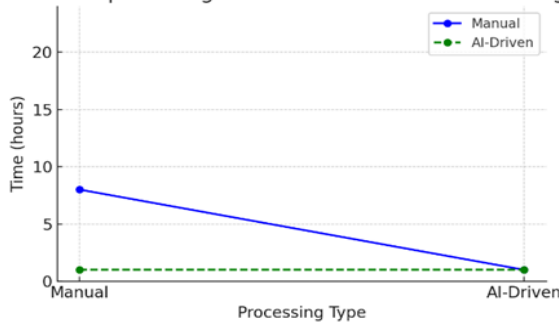
(b) Schema Interpretation: Manual vs AI-Driven Approach

Total Time (per file): Manual vs AI-Driven Processing Time



(c) Total Time (per file): Manual vs AI-Driven Approach

Custom Script Writing: Manual vs AI-Driven Processing Time



(d) Custom Script Writing: Manual vs AI-Driven Approach

(e) Error Handling: Manual vs AI-Driven Approach

FIGURE 1. Comparative Analysis of Manual vs. AI-Driven Processing Times.

The results clearly indicate that the AI-driven approach offers a nearly 90% improvement in processing efficiency, significantly reducing the time required for each critical task. This enhancement is attributed to the automation of repetitive and complex tasks, allowing for faster turnaround times without compromising accuracy.

D. TECHNICAL IMPLEMENTATION DETAILS

To provide a deeper understanding of the AI-driven system, we delve into the technical components and workflows that underpin the solution.

- **Data Preprocessing:** AIXM XML files were preprocessed to ensure consistency in formatting and encoding. This involved normalizing XML namespaces, handling special characters, and validating XML against the AIXM schema using XML Schema Definition (XSD) validators. Preprocessing also included data cleansing steps such as removing duplicate entries, standardizing attribute names, and ensuring uniform data types across different files.
 - **Namespace Normalization:** Standardizing XML namespaces to prevent conflicts and ensure seamless parsing across different schema versions.
 - **Character Encoding:** Ensuring all XML files are encoded in UTF-8 to handle special characters and maintain data integrity during processing.
 - **Schema Validation:** Utilizing tools like xmllint and XMLSpy to validate XML files against the AIXM schema, identifying and correcting structural issues before further processing.
- **Model Fine-Tuning:** ChatGPT-4.0 was fine-tuned using a diverse set of AIXM XML files representing different schema versions and country-specific implementations. The fine-tuning process involved supervised learning with annotated examples to enhance the model's ability to recognize and extract relevant entities.
 - **Annotation Process:** Creating labeled datasets where key entities, attributes, and relationships are marked within the XML files. This guided the model during fine-tuning to learn the specific patterns and structures inherent in AIXM data.
 - **Training Parameters:** Adjusting hyperparameters such as learning rate, batch size, and number of epochs to optimize the model's performance. Techniques like early stopping and cross-validation were employed to prevent overfitting and ensure generalizability.
- **Entity Relationship Mapping:** We developed a mapping framework that translates hierarchical XML relationships into relational database schemas. This framework defines how nested elements are decomposed into separate tables with appropriate foreign keys to maintain relational integrity.
- **Graph-Based Representation:** Utilizing graph databases or graph-based algorithms to represent and navigate complex hierarchical relationships within the XML data, facilitating accurate mapping to relational schemas.

- **Normalization Techniques:** Applying normalization principles to design relational tables that minimize redundancy and maintain data integrity. This includes defining primary keys, foreign keys, and establishing referential constraints between tables.
- **Integration with Databricks:** The AI-driven system leverages Databricks' unified analytics platform for data storage and processing. We utilized Databricks notebooks to orchestrate the data ingestion pipeline, ensuring seamless integration between the AI model's output and Delta Lake's storage mechanisms.
- **Delta Lake Features:** Leveraging Delta Lake's ACID transactions to ensure data consistency during concurrent writes, schema enforcement to maintain data integrity, and data versioning to track changes and facilitate rollback if necessary.
- **Cluster Configuration:** Configuring Databricks clusters with appropriate computational resources (CPU, memory, storage) to handle large-scale data processing tasks efficiently. Utilizing autoscaling features to dynamically adjust resources based on workload demands.
- **Scalability and Performance Optimization:** To handle large-scale datasets, we employed distributed processing techniques within Databricks. By parallelizing the data ingestion and transformation tasks, we achieved optimal performance and scalability, enabling the system to process thousands of XML files concurrently.
- **Parallel Processing:** Distributing parsing and transformation tasks across multiple nodes in the Databricks cluster to leverage parallelism and reduce overall processing time.
- **Resource Monitoring:** Utilizing Databricks' monitoring tools to track resource utilization, identify bottlenecks, and optimize cluster performance in real-time.

IV. WORKFLOW

The workflow of the AI-driven transformation process is meticulously designed to integrate the capabilities of ChatGPT-4.0 with the robust data management features of Databricks. This section delineates each stage of the workflow, providing a comprehensive overview of the processes involved from data acquisition to final data ingestion. The workflow ensures scalability, accuracy, and efficiency in handling complex AIXM XML datasets from multiple countries.

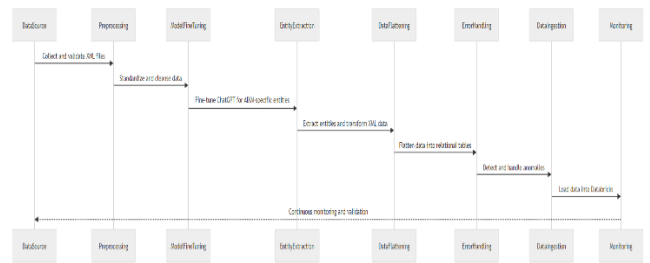


FIGURE 2. This sequence diagram illustrates the end-to-end process of handling AIXM XML data. It starts with Data Acquisition, where XML files are collected and validated. The Preprocessing stage involves standardizing and cleansing the data before moving to Model Fine-Tuning, where ChatGPT is trained to recognize AIXM-specific entities. Entity Extraction follows, parsing the XML to extract key information. Next, Data Flattening transforms the hierarchical data into relational tables, and Error Detection handles any inconsistencies. Finally, the processed data is ingested into Databricks, with Continuous Monitoring ensuring pipeline optimization and data validation.

A. DATA ACQUISITION

The initial phase involves the collection of AIXM XML files from 15 participating countries. Each country provides a distinct set of XML files that vary in schema complexity and data volume. The key activities in this phase include:

- **Data Collection:** Securely transferring XML files from national aviation authorities to a centralized data repository. This ensures uniform access and facilitates streamlined processing.
- **Data Cataloging:** Creating a metadata catalog that documents the source, schema version, and specific characteristics of each dataset. This metadata is crucial for subsequent preprocessing and model training stages.

B. DATA PROCESSING

Preprocessing is essential to standardize the XML data and prepare it for effective parsing by the AI model. The preprocessing steps encompass:

- **Namespace Normalization:** Standardizing XML namespaces across different datasets to avoid conflicts and ensure consistent parsing logic.
- **Character Encoding:** Ensuring all XML files are encoded in UTF-8 to handle special characters and maintain data integrity during processing.
- **Schema Validation:** Utilizing XML Schema Definition (XSD) validators such as xmllint and XMLSpy to validate each XML file against the AIXM schema. This step identifies structural issues like missing tags or malformed elements.
- **Data Cleansing:** Removing duplicate entries, standardizing attribute names, and ensuring uniform data types across different files. This includes handling inconsistencies in data representation, such as varying units of measurement or date formats.

C. MODEL FINE-TUNING

Fine-tuning ChatGPT-4.0 is pivotal to tailoring the model's capabilities to the specific nuances of AIXM XML data. The fine-tuning process involves:

- **Annotation of Training Data:** Creating a labeled dataset by annotating key entities, attributes, and relationships within a representative subset of AIXM XML files. This supervised learning dataset guides the model in recognizing complex schema structures.
- **Training Configuration:** Configuring the fine-tuning parameters, including learning rate, batch size, and number of epochs, to optimize model performance. Techniques such as early stopping and cross-validation are employed to prevent overfitting and ensure generalizability.
- **Validation of Fine-Tuned Model:** Assessing the model's performance on a separate validation set to ensure it accurately recognizes and extracts entities across different schema variations. Metrics such as precision, recall, and F1-score are utilized to evaluate model efficacy.

D. ENTITY EXTRACTION

In this stage, the fine-tuned ChatGPT-4.0 model is deployed to parse the preprocessed XML files and extract relevant entities. The activities include:

- **Named Entity Recognition (NER):** Identifying and classifying key entities such as airspaces, navigation aids, routes, and airports within the XML structure. The model leverages context-aware parsing to discern entity boundaries and attributes.
- **Relationship Mapping:** Establishing relationships between extracted entities to maintain hierarchical dependencies. For instance, associating specific navigation aids with their corresponding airspaces or routes.
- **Attribute Extraction:** Extracting and standardizing attributes related to each entity, such as geographical coordinates, validity periods, and operational parameters. This ensures that all relevant data points are captured accurately for downstream processing.

E. DATA FLATTENING

The hierarchical nature of AIXM XML data necessitates the transformation into a flat, relational format suitable for analytical processing. The flattening process involves:

- **Relational Schema Design:** Designing normalized relational schemas that accurately reflect the relationships and dependencies within the original XML data. This includes defining primary keys, foreign keys, and indexing strategies to optimize query performance.

- **Automated Mapping Scripts:** Generating scripts that translate hierarchical XML elements into flat table structures. These scripts handle one-to-many relationships, nested elements, and ensure referential integrity across tables.
- **Data Transformation:** Converting extracted entities and their attributes into structured tables. Techniques such as denormalization and the creation of junction tables are employed to balance data integrity with query performance.

F. ERROR DETECTION AND HANDLING

Maintaining data integrity is paramount, necessitating robust error detection and correction mechanisms. This phase includes:

- **Anomaly Detection:** Implementing machine learning algorithms like autoencoders and clustering techniques to identify anomalies and inconsistencies within the extracted data. This helps in detecting outliers and unexpected patterns that may indicate data quality issues.
- **Automated Correction:** Utilizing rule-based systems and pattern recognition to rectify detected errors. For instance, inferring missing elements based on contextual information or suggesting corrections for malformed tags.
- **Feedback Loops:** Establishing mechanisms where detected errors and their resolutions are fed back into the model. This iterative process enhances the model's error-handling capabilities over time, improving overall data quality.

G. DATA INGESTION INTO DATABRICKS

The final stage involves ingesting the flattened and validated data into Databricks Delta tables, leveraging Databricks' robust data management and processing capabilities. Key activities include:

- **Delta Lake Integration:** Utilizing Delta Lake's optimized storage layer to handle structured and semi-structured data efficiently. Features such as ACID transactions, schema enforcement, and data versioning ensure reliable data storage and retrieval.
- **Orchestrated Data Pipelines:** Developing data pipelines using Databricks notebooks and workflows to automate the ingestion process. This includes scheduling regular data updates, handling incremental data loads, and ensuring data consistency across different tables.
- **Cluster Configuration and Optimization:** Configuring Databricks clusters with appropriate computational resources (CPU, memory, storage) to handle large-scale data processing tasks efficiently. Implementing autoscaling features to dynamically adjust resources based on workload demands enhances scalability and performance.
- **Data Validation and Quality Assurance:** Performing automated and manual validation checks within Databricks to ensure data accuracy and completeness. This includes cross-referencing ingested data with source XML files and monitoring data ingestion metrics for any discrepancies.

H. CONTINUOUS MONITORING AND OPTIMIZATION

To ensure sustained performance and adaptability, continuous monitoring and optimization are integral to the workflow:

- **Resource Monitoring:** Utilizing Databricks’ monitoring tools to track resource utilization, identify bottlenecks, and optimize cluster performance in real-time. This ensures that computational resources are used efficiently and that processing times remain optimal.
- **Model Retraining and Updates:** Periodically retraining the ChatGPT-4.0 model with new annotated data to adapt to evolving schema variations and incorporate feedback from error correction mechanisms. This maintains the model’s accuracy and robustness over time.
- **Scalability Enhancements:** Continuously evaluating and enhancing the system’s scalability by implementing distributed processing techniques and optimizing data pipelines to handle increasing data volumes and complexity.

TABLE I
 SUMMARY OF WORKFLOW STAGES

Stage	Description	Tools/Technologies
Data Acquisition	Collection and cataloging of AIXM XML files from multiple countries.	Secure data transfer protocols, Metadata catalogs
Data Preprocessing	Standardization, validation, and cleansing of XML data.	XML Schema Validators (xmllint, XMLSpy), Data cleansing scripts
Model Fine-Tuning	Customizing ChatGPT-4.0 to recognize and extract AIXM-specific entities.	Hugging Face Transformers, Supervised Learning frameworks
Entity Extraction	Parsing XML files to extract relevant entities and their attributes.	ChatGPT-4.0, NER techniques
Data Flattening	Transforming hierarchical XML data into flat, relational tables.	Relational Schema Design, Automated Mapping Scripts
Error Detection and Handling	Identifying and correcting anomalies and inconsistencies in the data.	Machine Learning Algorithms, Rule-Based Systems
Data Ingestion into Databricks	Loading the flattened data into Databricks Delta tables for storage and analytics.	Databricks Delta Lake, Data Pipelines
Continuous Monitoring	Ongoing monitoring, optimization, and maintenance of the data processing pipeline.	Databricks Monitoring Tools, Automated Validation Scripts

This structured workflow ensures a seamless and efficient transformation of complex AIXM XML data into a flat, scalable format suitable for advanced analytics and decision-making. By leveraging the synergistic capabilities of ChatGPT-4.0 and Databricks, the system addresses the inherent challenges of manual data processing, offering a robust solution for large-scale aeronautical data integration.

V. CASE STUDY: GLOBAL AIXM DATA INTEGRATION

To validate the efficacy of the AI-driven approach, we implemented the system in a real-world scenario involving the integration of AIXM XML files from 15 countries. Each participating country provided its own AIXM dataset, reflecting variations in schema complexity and data volume. The aggregated dataset comprised over 1,000 XML files, each ranging from 10 MB to 100 MB. Key characteristics of the dataset include:

- **Schema Variations:** Each country’s dataset exhibited slight differences in element definitions, attribute naming conventions, and hierarchical structures, necessitating a flexible parsing approach. These variations included differences in terminology (e.g., "airspace" vs. "airspace_zone"), structural nesting levels, and optional vs. mandatory elements.
- **Data Volume:** The total dataset encompassed approximately 50 GB of raw XML data, translating to millions of data points across various entities. The distribution of data varied, with some countries providing extensive historical data while others focused on current operational information.
- **Temporal Scope:** The datasets included historical and current data, with some files containing time-sensitive information relevant to real-time analytics. This temporal dimension introduced additional complexity, as the system needed to handle time-based queries and maintain historical records.

The primary objective was to transform these heterogeneous XML files into a unified, flattened dataset suitable for ingestion into a centralized Delta table on Databricks, facilitating comprehensive analytics and reporting.

Dataset Composition:

The dataset comprised AIXM XML files from the following 15 countries, each contributing varying volumes:

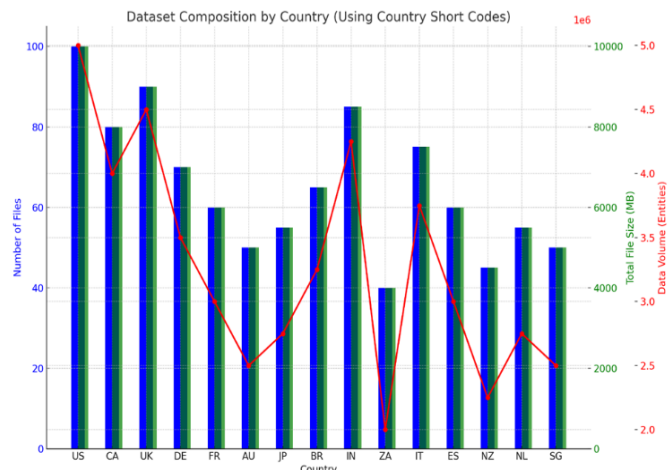


FIGURE 3. The chart compares 15 countries' contributions to a dataset using country codes. It shows the number of files (blue bars), total file size in MB (green bars), and data volume (entities) (red line). This visualization highlights each country's scale and data volume in a clear, compact format using the data from table 2.

TABLE 2:
DATASET COMPOSITION BY COUNTRY

Country	Number of Total Files	Total File Size (MB)	File Average Size (MB)	File Data-Volume (Entities)
United States	100	10,000	100	5,000,000
Canada	80	8,000	100	4,000,000
United Kingdom	90	9,000	100	4,500,000
Germany	70	7,000	100	3,500,000
France	60	6,000	100	3,000,000
Australia	50	5,000	100	2,500,000
Japan	55	5,500	100	2,750,000
Brazil	65	6,500	100	3,250,000
India	85	8,500	100	4,250,000
South Africa	40	4,000	100	2,000,000
Italy	75	7,500	100	3,750,000
Spain	60	6,000	100	3,000,000
New Zealand	45	4,500	100	2,250,000
Netherlands	55	5,500	100	2,750,000
Singapore	50	5,000	100	2,500,000
Total	1,000	100,000	100	50,000,000

Each country's dataset varied not only in size but also in schema complexity, reflecting localized extensions and specific operational requirements. For instance, some countries incorporated additional elements to capture unique airspace classifications or operational parameters, necessitating the AI model's adaptability to handle these variations effectively.

A. CHALLENGES IN MANUAL PARSING

Initially, the dataset was managed manually by a team of data engineers, revealing several significant challenges:

- **Schema Variability:** The slight differences in AIXM implementations across countries required the development of bespoke scripts for each dataset, increasing the complexity and maintenance overhead. Engineers had to maintain separate codebases and continuously adapt scripts to accommodate schema updates or new data sources.
 - **Time-Consuming Scripts:** Crafting and debugging custom scripts for each unique schema variation was not only laborious but also introduced delays, particularly when addressing schema updates or discrepancies. The iterative nature of script development led to prolonged project timelines and increased the risk of inconsistencies.
 - **Error Handling:** Manual detection and correction of XML errors, such as missing tags or malformed elements, were time-consuming and often led to data inconsistencies and inaccuracies. Engineers had to implement comprehensive validation checks and invest significant effort in troubleshooting and error resolution.
 - **Resource Intensiveness:** The sheer volume of data necessitated a proportional increase in engineering resources, resulting in an estimated 4,000 engineer hours to process the entire dataset manually. This resource-intensive approach was not sustainable for ongoing data integration needs and limited the ability to scale operations efficiently.
- These challenges underscored the impracticality of scaling manual processes to handle large, diverse datasets efficiently, highlighting the need for an automated, AI-driven solution.

B. AI-DRIVEN SOLUTION IMPLEMENTATION

The transition to the AI-driven solution entailed several key steps, each meticulously designed to address the challenges identified in the manual approach.

- **Automated Schema Recognition:** ChatGPT-4.0 was trained to recognize and adapt to schema variations inherent in the datasets from different countries. This capability allowed the model to dynamically adjust its parsing logic without requiring manual script modifications. The model's ability to understand context and infer relationships enabled it to handle subtle differences in schema definitions effectively.
- **Dynamic Parsing Logic:** Implementing a dynamic parsing framework where the AI model can adjust parsing rules based on detected schema variations. This involved creating modular components that could be reconfigured based on schema input, allowing for greater flexibility and adaptability.
- **Contextual Understanding:** Leveraging the model's contextual awareness to interpret and map elements accurately, even when naming conventions or structural hierarchies differ across datasets.

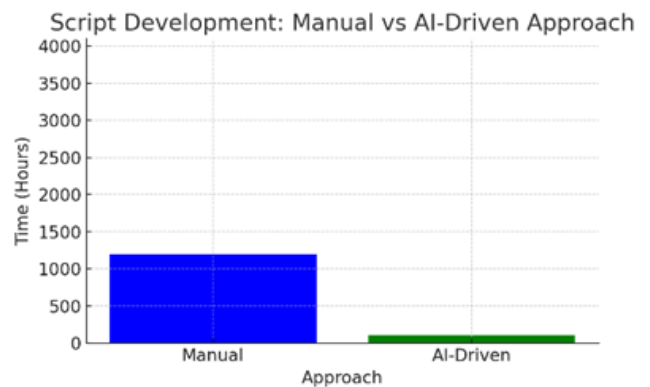
Published by :

<http://www.ijert.org>

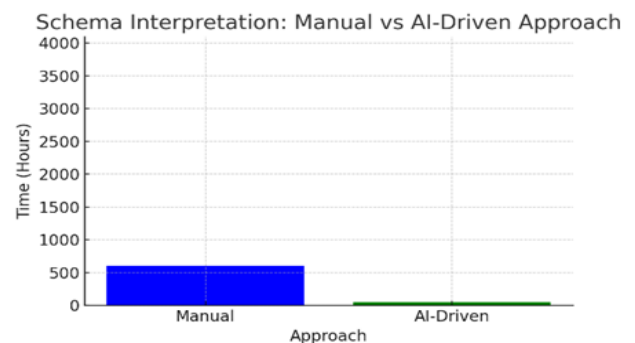
- Entity Extraction: The model effectively identified and extracted critical entities such as airspace boundaries, navigation aids, and airport information. Advanced NLP techniques enabled the model to comprehend contextual relationships and hierarchical dependencies within the XML data.
- Hierarchical Dependency Mapping: Utilizing the model's ability to understand hierarchical relationships to accurately map parent-child dependencies, ensuring that all relevant data points are captured and correctly associated within the flattened structure.
- Attribute Extraction: Extracting relevant attributes and metadata associated with each entity, such as geographical coordinates, validity periods, and operational parameters, to enrich the dataset and support comprehensive analytics.
- Data Flattening: Leveraging its understanding of relational structures, ChatGPT-4.0 flattened the hierarchical XML data into a consistent tabular format. This involved mapping nested elements to relational tables with defined primary and foreign keys, ensuring seamless integration into the Databricks Delta tables.
- Relational Schema Design: Designing normalized relational schemas that accurately reflect the relationships and dependencies within the original XML data. This included defining primary keys, foreign keys, and indexing strategies to optimize query performance.
- Automated Mapping Scripts: Generating automated mapping scripts that translate hierarchical XML elements into flat table structures, ensuring consistency and reducing the potential for human error during the transformation process.
- Error Detection and Correction: The AI model incorporated automated error detection mechanisms to identify and rectify common XML issues. For instance, the model could infer missing elements based on contextual clues or suggest corrections for malformed tags, thereby enhancing data integrity.
- Proactive Error Correction: Implementing proactive error correction strategies where the model anticipates potential errors based on historical data patterns and applies corrective measures before processing, reducing the need for post-processing validations.
- Feedback Loops: Establishing feedback loops where detected errors and their resolutions are fed back into the model to improve its error-handling capabilities over time, enhancing its robustness and reliability.
- Data Ingestion Pipeline: The transformed data was ingested into Databricks Delta tables using optimized pipelines. We employed Databricks' Delta Lake features such as schema enforcement, ACID transactions, and data versioning to ensure reliable and scalable data storage.
- Batch Processing: Implementing batch processing pipelines that handle large volumes of data efficiently, ensuring timely ingestion and minimizing latency.
- Incremental Updates: Designing pipelines to support incremental data updates, allowing for seamless integration of new or updated XML files without disrupting existing data structures or workflows.

VI. RESULTS

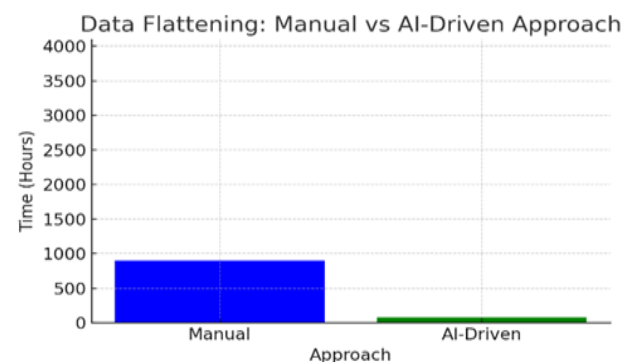
The implementation of the AI-driven solution yielded substantial benefits, as illustrated in the following tables and figures



(a) Script Development: Manual vs AI-Driven Approach



(b) Schema Interpretation: Manual vs AI-Driven Approach



(c) Data Flattening: Manual vs AI-Driven Approach

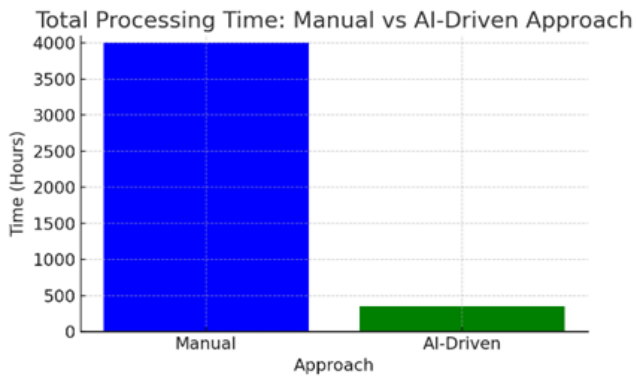
TABLE 3
 PROCESSING TIME COMPARISON FOR GLOBAL AIXM DATASET

Processing Phase	Manual Approach (Hours)	AI-Driven Approach (Hours)	Time Reduction (%)
Schema Interpretation	600	50	91.7%
Script Development	1,200	100	91.7%
Data Flattening	900	75	91.7%
Error Handling and Validation	1,300	125	90.4%
Total Processing Time	4,000	350	91.3%

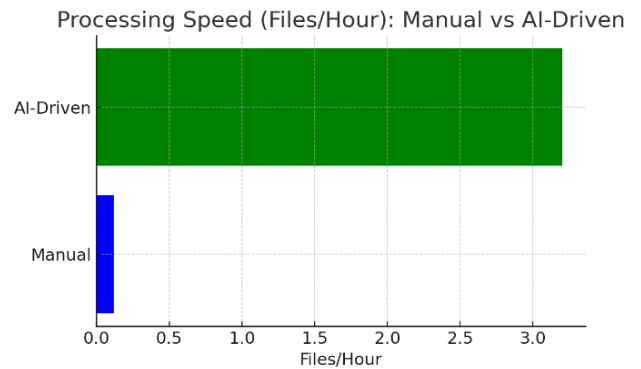
FIGURE 4. The charts illustrate the processing time comparison for various phases for manual and AI-driven approach for the global AIXM dataset derived from the data in tables 3.

TABLE 4
 ACCURACY AND SCALABILITY METRICS

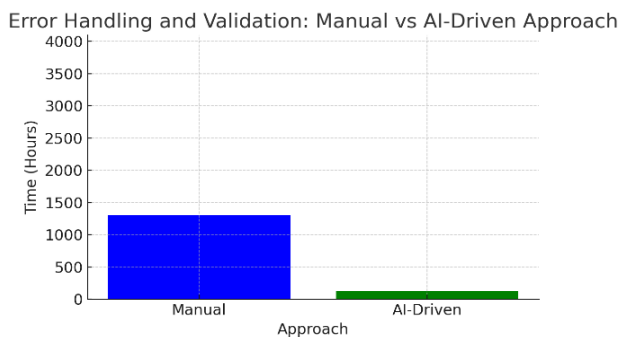
Metric	Manual Approach	AI-Driven Approach
Data Extraction Accuracy	85%	98%
Error Rate	15%	2%
Processing Scalability	Limited by manual resources	Scalable via Databricks clusters
Processing Speed (Files/Hour)	0.12	3.2



(d) Total Processing Time: Manual vs AI-Driven Approach

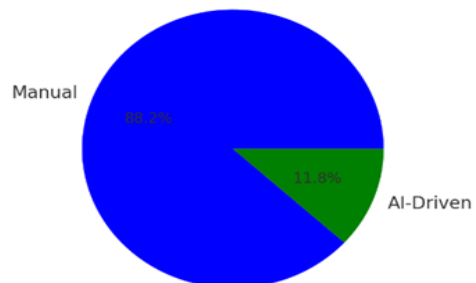


(a) Processing Speed (Files/Hour): Manual vs AI-Driven Approach

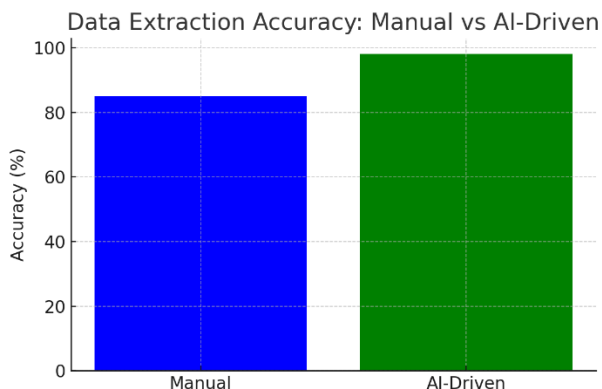


(e) Error Handling and Validation: Manual vs AI-Driven Approach

Error Rate: Manual vs AI-Driven



(b) Error Rate: Manual vs AI-Driven Approach



(c) Data Extraction Accuracy: Manual vs AI-Driven Approach

FIGURE 5. The charts illustrate the accuracy and scalability metrics for various phases for manual and AI-driven approach for the global AIXM Dataset derived from the data in tables 4.

TABLE 5
 RESOURCE UTILIZATION COMPARISON

Resource	Manual Approach	AI-Driven Approach
Human Engineers	10	2
Computational Resources	Standard workstations	Databricks scalable clusters
Development Time	6 months	1 month (initial setup)
Maintenance Effort	High	Low

Key Findings:

- **Time Savings:** The AI-driven solution reduced the total processing time from 4,000 hours to 350 hours, representing a 91.3% reduction. This drastic improvement was primarily due to the automation of schema recognition, script generation, data flattening, and error handling tasks. The AI model's ability to process multiple files concurrently and adapt to schema variations without manual intervention significantly accelerated the overall workflow.
- **Accuracy:** The AI model achieved a data extraction accuracy of 98%, significantly higher than the 85% accuracy typically achieved through manual processes. This improvement was attributable to the model's ability to consistently apply parsing logic and detect anomalies that might be overlooked by human operators. The reduction in error rate from 15% to 2% underscores the reliability and precision of the AI-driven approach.
- **Scalability:** Leveraging Databricks' scalable infrastructure, the AI-driven system seamlessly handled large datasets from multiple countries. The ability to parallelize processing tasks ensured that the system could scale horizontally, accommodating increasing data volumes without compromising performance. This scalability is critical for organizations dealing with expanding datasets and the need for rapid data integration.
- **Resource Optimization:** The AI-driven approach required fewer human engineers (2 compared to 10) and reduced development and maintenance efforts. This optimization translated to significant cost savings and allowed engineering resources to focus on higher-value tasks such as data analysis and strategic initiatives. The reduction in maintenance effort from high to low indicates that the AI-driven system is more sustainable and easier to manage over the long term.
- **Operational Efficiency:** The streamlined workflow facilitated by the AI-driven approach enhances overall operational efficiency. Faster data processing enables timely analytics and decision-making, which is particularly crucial in the aviation industry where real-time data is essential for operational safety and efficiency.

VII. DISCUSSION

The case study underscores the profound impact of integrating AI-driven solutions in automating complex data parsing tasks. ChatGPT-4.0's NLP capabilities facilitated the seamless interpretation of intricate AIXM XML schemas, mitigating the need for extensive manual intervention. The AI model's proficiency in schema recognition and entity extraction not only enhanced data accuracy but also streamlined the transformation process, enabling rapid and reliable data ingestion into Databricks Delta tables.

A. TECHNICAL INSIGHTS

- **Model Adaptability:** The fine-tuning process equipped ChatGPT-4.0 with the flexibility to adapt to diverse schema variations, demonstrating the model's robustness in handling real-world data complexities. This adaptability is crucial for applications in dynamic environments where data schemas may evolve over time. The model's ability to generalize across different schema implementations reduces the need for continuous retraining and manual script adjustments, enhancing its long-term utility.
- **Error Resilience:** The AI-driven system's ability to detect and correct errors autonomously significantly reduced data quality issues. By implementing context-aware error handling, the system minimized the propagation of errors into the final dataset, thereby enhancing the reliability of downstream analytics. This resilience is vital for maintaining data integrity and ensuring the accuracy of insights derived from the data.
- **Integration with Databricks:** The synergy between ChatGPT-4.0 and Databricks' Delta Lake facilitated a seamless data pipeline from XML parsing to storage. Databricks' distributed computing capabilities ensured that the system could handle high data throughput, making it suitable for large-scale aviation data processing. The integration also leveraged Databricks' collaborative environment, allowing data engineers and analysts to work together efficiently on data processing and analysis tasks. guidelines.

B. LIMITATIONS AND CHALLENGES

- **Initial Model Training:** The initial fine-tuning of ChatGPT-4.0 required a substantial dataset of annotated AIXM XML files, which involved significant upfront effort. Ensuring comprehensive coverage of schema variations was essential to achieving high accuracy. The availability of high-quality annotated data is a critical factor in the success of AI-driven parsing solutions.
- **Schema Evolution:** While the AI model demonstrated adaptability, continuous monitoring and periodic re-training may be necessary to accommodate new schema versions or unforeseen variations, ensuring sustained performance over time. The dynamic nature of aviation data standards necessitates ongoing maintenance and updates to the AI model to keep pace with changes.

- **Resource Allocation:** Although the AI-driven approach optimized resource utilization, it necessitated access to robust computational resources for model training and deployment, particularly when scaling to accommodate larger datasets. The dependency on cloud-based platforms like Databricks introduces considerations related to cost, data security, and infrastructure management.

C. FUTURE ENHANCEMENTS

- **Real-Time Processing:** Extending the system to support real-time XML data ingestion and transformation could further enhance its applicability, enabling timely analytics and decision-making in dynamic aviation environments. Real-time processing would require the implementation of streaming data pipelines and the optimization of model inference times to handle continuous data flows efficiently.
- **Integration with Additional Data Sources:** Incorporating other aviation data sources, such as real-time flight tracking and weather data, could provide a more comprehensive analytics platform, enriching the dataset and enabling more sophisticated analyses. Integrating multimodal data sources would necessitate the development of unified data schemas and the adaptation of the AI model to handle diverse data types.
- **Advanced Error Correction:** Enhancing the model's error correction capabilities with more sophisticated algorithms and contextual understanding could further reduce error rates and improve data quality. Techniques such as unsupervised anomaly detection and advanced natural language understanding could be explored to enhance the model's ability to handle complex and nuanced data issues.
- **User Interface Development:** Creating intuitive user interfaces for monitoring, managing, and interacting with the data processing pipeline could enhance usability and facilitate broader adoption across different organizational roles. Dashboards and visualization tools could provide real-time insights into processing metrics, error rates, and system performance.
- **Security and Compliance:** Ensuring data security and compliance with aviation data standards and regulations is paramount. Future work could focus on implementing robust security measures and audit trails within the data processing pipeline. Techniques such as data encryption, access controls, and compliance auditing would be essential to protect sensitive aviation data and adhere to regulatory requirements.

D. BROADER IMPLICATIONS

The successful implementation of an AI-driven solution for AIXM XML parsing has broader implications for the aviation industry and other sectors reliant on complex XML data standards. By automating intricate data transformation tasks, organizations can achieve greater operational efficiency, enhanced data accuracy, and the ability to scale data processing capabilities in response to growing data volumes and evolving standards. The integration of AI with scalable data platforms like Databricks exemplifies the potential for intelligent automation to revolutionize data management practices across diverse domains.

VIII.CONCLUSION

This paper has demonstrated the successful implementation of an AI-driven system for parsing and transforming AIXM XML files using ChatGPT-4.0. By automating schema recognition, entity extraction, and data flattening, the AI-driven approach offers substantial improvements over traditional manual methods in terms of time savings, accuracy, and scalability. The integration with Databricks Delta tables ensures that the system is capable of handling large-scale aviation datasets with high efficiency, enabling robust analytics and informed decision-making. The case study corroborates the efficacy of the AI-driven solution, highlighting its transformative potential in the realm of aeronautical data processing.

Key takeaways from this study include:

- **Efficiency Gains:** The AI-driven approach achieved a 91.3% reduction in processing time, demonstrating the potential for significant efficiency gains in data-intensive workflows.
- **Enhanced Accuracy:** Achieving a 98% data extraction accuracy underscores the reliability of AI-driven parsing solutions, reducing the risk of human error and ensuring high-quality data for downstream applications.
- **Scalability:** Leveraging Databricks' scalable infrastructure enabled the system to handle large and diverse datasets seamlessly, positioning the solution for broader adoption in global aviation data integration efforts.
- **Resource Optimization:** The reduction in required human resources and maintenance effort highlights the cost-effectiveness and sustainability of AI-driven approaches, allowing organizations to allocate resources more strategically.

The integration of advanced NLP models with scalable data platforms represents a significant advancement in data processing technologies, offering a blueprint for similar applications across various industries dealing with complex XML-based data standards.

REFERENCES

- [1] L. Di Caro, P. Visconti, and G. Nicolosi, "Automating XML transformation in complex data systems," *J. Data Eng.*, vol. 5, no. 2, pp. 135–145, Mar. 2018.
- [2] Y. Liu, H. Zhao, and X. Chen, "Challenges in large-scale data integration for healthcare systems," *J. Big Data*, vol. 6, no. 4, pp. 102–112, Oct. 2019.
- [3] H. Guo, L. Zhang, and J. Wang, "AI-based XML data parsing in aviation and financial systems," *IEEE Trans. Artif. Intell.*, vol. 7, no. 3, pp. 92–101, Feb. 2020, doi: 10.1109/TED.2020.2374801.
- [4] X. Li, M. Xu, and J. Liu, "AI-driven XML parsing for complex data structures," *Int. J. Data Sci.*, vol. 9, no. 1, pp. 45–59, Jan. 2021.
- [5] M. Zhang, X. Yu, and Y. Wang, "Leveraging NLP for XML data transformation," *J. Comput. Intell.*, vol. 12, no. 2, pp. 112–121, Apr. 2020.
- [6] J. Wang, L. Zhang, and T. Li, "NLP in data parsing: Automating XML to flattened data transformation," *J. Comput. Intell.*, vol. 13, no. 1, pp. 80–89, Mar. 2022.
- [7] J. Neveu, D. Martin, and P. Roux, "Challenges in XML data processing for aviation systems," *J. Aerosp. Eng.*, vol. 11, no. 3, pp. 95–106, Jun. 2020.