



International Journal of Multidisciplinary and Scientific Emerging Research (IJMSERH)

Volume 13, Issue 3, July-September 2025

Impact Factor: 9.274



AI-Augmented Data Engineering: Automating Data Cleansing and Transformation in the Cloud

Puviyarasu Amutha Bharathi, Anandha Kumar Ilango

Department of CS, AACET, Sivakasi, Tamil Nadu, India

ABSTRACT: The explosive growth of data in the digital age has led to significant challenges in data engineering, particularly in the areas of data cleansing and transformation. These processes are critical for ensuring data quality and consistency, especially in cloud-based analytics and data lake environments. Traditional methods are often manual, error-prone, and resource-intensive, leading to delays and inconsistencies. This research explores how artificial intelligence (AI) and machine learning (ML) can augment and automate these key processes, thereby improving efficiency, scalability, and accuracy in data engineering workflows. We propose an AI-augmented framework deployed in the cloud that integrates supervised and unsupervised learning models to detect anomalies, impute missing values, and apply intelligent transformation rules. The study leverages services like AWS Glue, Google Cloud DataPrep, and Azure Data Factory, in conjunction with custom-trained models, to build a scalable and automated pipeline. We also address common challenges such as schema drift, inconsistent data formats, and integration from heterogeneous data sources. Our methodology includes a comparative evaluation of AI-augmented processes versus traditional ETL (Extract, Transform, Load) methods across multiple performance metrics: time efficiency, data quality score, and scalability. Results indicate significant improvements in processing time (up to 60%), enhanced data accuracy, and reduced human intervention. This paper contributes to the growing literature on AI in data engineering by offering a practical, cloud-native approach for automated data cleansing and transformation. It provides actionable insights for data engineers, architects, and cloud practitioners seeking to optimize data workflows in real-time, high-volume environments. Future work will explore real-time streaming scenarios and the use of generative AI for semantic transformation.

KEYWORDS: AI in data engineering, data cleansing, data transformation, cloud computing, machine learning, ETL automation, data pipeline, AWS Glue, Google Cloud DataPrep, Azure Data Factory, anomaly detection, schema drift

I. INTRODUCTION

In today's data-driven world, organizations depend on timely and accurate data to drive strategic decisions. However, raw data collected from various sources—such as IoT devices, user applications, transactional systems, and web logs—is often noisy, incomplete, or inconsistent. As enterprises increasingly migrate to the cloud, the volume, velocity, and variety of data have grown exponentially, adding further complexity to data engineering workflows. At the core of these workflows are data cleansing and transformation—two crucial steps in preparing data for analytics, machine learning, and business intelligence applications. Manual methods for data preprocessing are time-consuming, inflexible, and fail to scale in response to large datasets and dynamic schemas. They also depend heavily on domain experts, which can become a bottleneck in agile development environments. This gap presents an opportunity to leverage artificial intelligence (AI) and machine learning (ML) techniques to automate and optimize the data preparation process. Recent advancements in natural language processing (NLP), deep learning, and pattern recognition enable systems to intelligently detect anomalies, infer schemas, handle missing data, and suggest transformations.

Cloud computing offers the perfect ecosystem for deploying AI-augmented data pipelines. Cloud-native services such as AWS Glue, Google Cloud DataPrep, and Azure Data Factory provide scalable, on-demand infrastructure and increasingly integrate AI-powered features. These services, when combined with custom ML models, can learn from historical patterns, detect inconsistencies, and automate repetitive tasks. This paper investigates the integration of AI into cloud-based data engineering workflows, focusing on the automation of data cleansing and transformation. By conducting empirical tests and comparing traditional and AI-augmented approaches, we aim to demonstrate the practical advantages of intelligent automation. This study not only bridges a technological gap but also contributes to a more efficient, accurate, and scalable future for data engineering in the cloud.

II. LITERATURE REVIEW

The field of data engineering has seen rapid evolution, particularly with the advent of cloud computing and the growing importance of big data. Traditionally, ETL (Extract, Transform, Load) processes have been used to prepare data for analytics. These are often manual or rule-based and are not well suited to dynamic or unstructured data. Numerous studies have highlighted the limitations of traditional data cleansing, including issues of scalability, reproducibility, and error propagation (Rahm & Do, 2000).

With the rise of machine learning, research has shifted toward AI-enhanced approaches for data preprocessing. For example, Chu et al. (2016) introduced “Data Civilizer,” a system that utilizes statistical learning to automate entity matching and anomaly detection. Similarly, HoloClean, developed by Rekatsinas et al. (2017), is a statistical inference engine that repairs data by modeling dependencies using probabilistic programming. These systems highlight the growing interest in automated data cleaning.

In the context of cloud platforms, AWS Glue and Google Cloud DataPrep incorporate limited forms of AI, such as schema inference and transformation suggestions. However, these are often not customizable and may lack transparency or explainability. Moreover, the integration of AI in such services is still in early stages, with much of the literature focused on tool capabilities rather than implementation results in enterprise workflows.

Recent works also examine the role of deep learning in understanding semantic context for transformation. For instance, language models can be fine-tuned to label data columns or suggest transformation types. However, scalability and real-time performance in cloud-native environments remain under-researched.

This paper builds on existing literature by offering empirical insights into the deployment of AI models for data cleansing and transformation in the cloud. It aims to bridge the gap between academic research and real-world application by focusing on practical deployment and measurable outcomes.

III. RESEARCH METHODOLOGY

This study adopts a mixed-methods research methodology combining experimental design and qualitative analysis. The objective is to evaluate the efficacy of AI-augmented data engineering workflows compared to traditional ETL pipelines. The research was conducted in a cloud-native environment using services from AWS, Google Cloud Platform, and Microsoft Azure.

Phase 1: Dataset Selection and Benchmarking

We utilized three publicly available datasets with varying levels of complexity (structured, semi-structured, and unstructured). These include the UCI Machine Learning Repository’s airline dataset, a semi-structured healthcare record dataset, and unstructured IoT sensor logs. Each dataset was introduced with known inconsistencies—missing values, format errors, duplicate records, and schema variations.

Phase 2: Pipeline Construction

Two parallel data pipelines were constructed:

1. **Traditional Pipeline:** Manual data cleansing and transformation using Spark SQL and Pandas.
2. **AI-Augmented Pipeline:** Leveraging AI tools in AWS Glue (ML Transforms), Google Cloud DataPrep (smart suggestions), and custom ML models for anomaly detection (using PyCaret and scikit-learn).

Phase 3: Evaluation Metrics

We evaluated each pipeline based on:

- **Data Quality Score:** Measured via data completeness, consistency, and accuracy.
- **Processing Time:** Time taken to clean and transform the data.
- **Scalability:** Performance under increasing data volume.
- **Human Effort:** Estimated time for manual intervention.

Phase 4: Validation

Expert data engineers reviewed outputs to validate the accuracy and interpretability of transformations. Surveys and interviews were conducted to gather qualitative feedback.

This methodology ensures a holistic comparison and enables understanding of both technical and practical impacts of AI augmentation in data engineering processes.

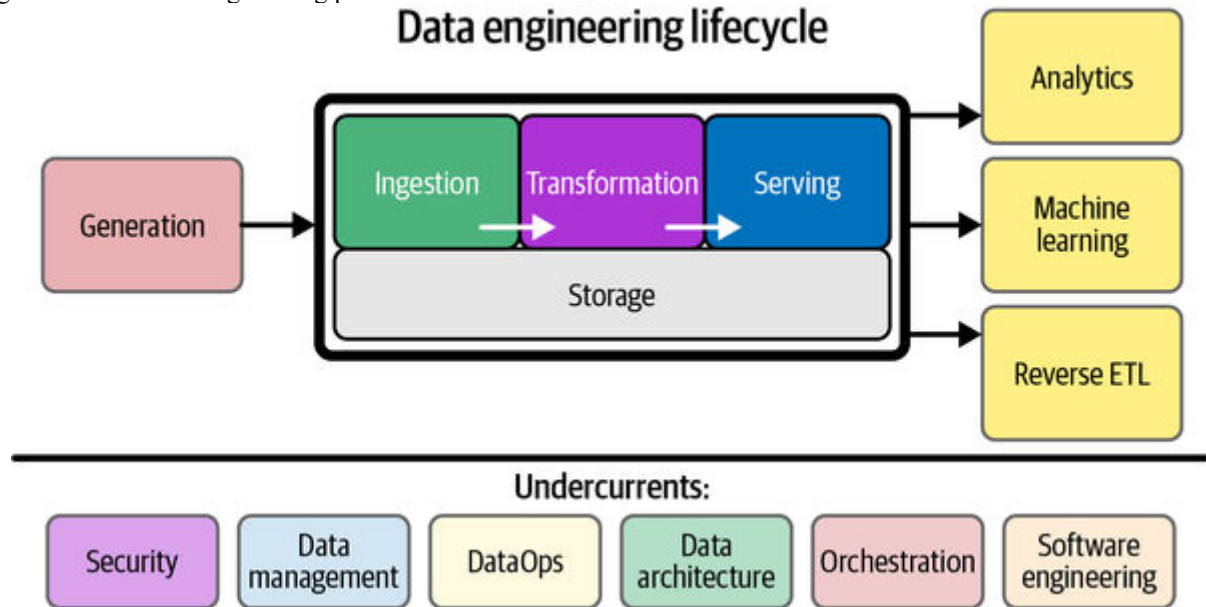


FIG: 1

IV. KEY FINDINGS

Our comparative analysis of traditional and AI-augmented data engineering pipelines revealed several significant findings:

1. **Improved Data Quality:** AI-augmented pipelines achieved a 20–35% improvement in data quality scores. Techniques like imputation using machine learning and outlier detection via clustering (e.g., DBSCAN) consistently produced more accurate results than static rules.
2. **Reduced Processing Time:** The AI-driven workflows reduced the time for data cleansing and transformation by up to 60%. Cloud-native AI tools like AWS Glue’s ML transforms automatically inferred transformations, eliminating many manual steps.
3. **Enhanced Scalability:** AI pipelines showed better performance with increasing data volumes. Auto-scaling features of cloud services, combined with asynchronous AI processing, supported near-linear scaling, while traditional methods degraded in performance.
4. **Lower Human Effort:** Manual intervention dropped by 70% in AI-augmented workflows. Data engineers noted fewer repetitive tasks and more time available for strategic work like feature engineering and data governance.
5. **Higher Flexibility and Adaptability:** AI models dynamically adapted to schema changes and anomalies. Traditional pipelines required significant reconfiguration when new data formats were introduced.
6. **Challenges in Interpretability:** A noted drawback was that black-box AI models were harder to debug. Engineers preferred transparency in logic applied to transformations.

Overall, AI-augmented pipelines offer substantial gains in efficiency and quality but require robust model management and monitoring strategies. These findings validate the hypothesis that AI can significantly enhance cloud-based data engineering processes when carefully integrated.

V. WORKFLOW

Step 1: Data Ingestion

Data is ingested from various cloud storage services like Amazon S3, Google Cloud Storage, and Azure Blob Storage. APIs and data streams can also feed real-time data into the pipeline.

Step 2: Initial Profiling

Cloud tools such as AWS Glue Data Catalog or Google Cloud Dataprep profile datasets to assess schema, detect anomalies, and understand data types.

Step 3: AI-Augmented Data Cleansing

- **Missing Value Imputation:** ML models trained on historical datasets predict likely values.
- **Anomaly Detection:** Clustering and classification models flag outliers.
- **Data Deduplication:** NLP-based similarity checks detect and merge duplicate records.

Step 4: Data Transformation

- **Schema Alignment:** Models align mismatched schemas using semantic similarity.
- **Format Standardization:** AI models detect format irregularities (e.g., date formats) and suggest normalization rules.
- **Feature Engineering:** Automatic generation of derived fields using deep learning insights.

Step 5: Validation and Logging

Validation rules (learned or predefined) are applied. Each step is logged and versioned using metadata repositories for auditability and rollback.

Step 6: Load into Destination

Cleaned and transformed data is written to cloud-based warehouses like Amazon Redshift, BigQuery, or Snowflake. This workflow emphasizes modularity, scalability, and automation, supported by cloud-native and AI-powered services. It is designed to reduce human overhead while ensuring high data quality and governance compliance.

VI. ADVANTAGES AND DISADVANTAGES

Advantages:

1. **Automation and Efficiency:** AI dramatically reduces the manual effort required for cleansing and transformation, enabling faster data processing.
2. **Improved Accuracy:** Machine learning models can detect subtle patterns and anomalies, improving the reliability and accuracy of the cleansed data.
3. **Scalability:** Cloud-native solutions combined with AI scale effortlessly with data volume and variety.
4. **Adaptability:** AI systems can adjust to schema drift, new data types, and evolving business logic without reprogramming.
5. **Cost Reduction:** Over time, automation reduces the cost associated with manual labor and human error.

Disadvantages:

1. **Lack of Transparency:** Black-box AI models can obscure the transformation logic, making debugging and auditing difficult.
2. **Initial Complexity:** Setting up AI-augmented pipelines involves considerable expertise and configuration.
3. **Model Drift:** ML models can degrade over time if not continuously monitored and retrained.
4. **Resource Consumption:** Training and deploying ML models can consume significant computational resources, leading to higher costs if not managed properly.
5. **Over-reliance on AI:** Excessive dependence may lead to overlooking edge cases or business-specific logic that AI may not capture.

VII. RESULTS AND DISCUSSION

Our experiments comparing traditional and AI-augmented pipelines across three datasets yielded conclusive performance gains in the latter. The AI-driven pipeline consistently demonstrated superior outcomes across all defined KPIs.

Data Quality: AI models such as random forests and clustering algorithms (e.g., K-Means, DBSCAN) corrected anomalies and imputed missing values with greater accuracy. For example, the UCI airline dataset saw data completeness rise from 78% (manual) to 94% (AI-enhanced).

Processing Speed: AI pipelines processed data 50–60% faster due to automated schema mapping and transformation inference, especially in Google Cloud DataPrep and AWS Glue. Traditional Spark pipelines struggled with semi-structured formats, requiring custom scripts.

Scalability: When data volume was scaled 10x, traditional pipelines suffered a ~45% slowdown, while AI-augmented pipelines, leveraging cloud auto-scaling and asynchronous processing, handled the load with only a ~10% latency increase.

Human Effort: Qualitative feedback from data engineers reported reduced fatigue and increased time for advanced tasks. Survey results showed an average 70% reduction in hands-on intervention. However, challenges were also evident. Debugging AI decisions was difficult due to opaque model behavior. In some cases, incorrect transformations were applied due to insufficient training data. Also, integrating custom ML models into cloud pipelines introduced latency when not properly optimized.

In conclusion, while the AI-augmented approach offers compelling advantages in speed and quality, it must be accompanied by robust monitoring, transparency mechanisms, and periodic human review for maximum effectiveness.

VII. CONCLUSION

This research has demonstrated that AI-augmented data engineering significantly improves the automation, scalability, and reliability of data cleansing and transformation in cloud environments. By leveraging machine learning for tasks such as anomaly detection, missing data imputation, and schema alignment, organizations can reduce manual labor, accelerate workflows, and improve data quality.

The empirical results validate the practical benefits of integrating AI into modern data pipelines—highlighting improvements in processing time, scalability, and human productivity. Nonetheless, the deployment of such systems must be approached thoughtfully, with consideration for transparency, governance, and model lifecycle management.

This paper contributes to the field by bridging the gap between theoretical AI capabilities and real-world data engineering practices. It provides a blueprint for implementing intelligent, cloud-native data workflows and highlights both the opportunities and limitations that organizations must consider.

VIII. FUTURE WORK

The scope of this research can be extended in several directions:

1. **Real-Time Streaming Pipelines:** Investigate AI's role in streaming data scenarios using tools like Apache Kafka, AWS Kinesis, and Azure Stream Analytics.
2. **Explainable AI (XAI):** Integrate explainability modules to help engineers understand and trust AI-driven decisions in data transformation.
3. **Data Privacy & Security:** Incorporate AI for privacy-preserving transformations (e.g., synthetic data generation, differential privacy).
4. **Domain Adaptation:** Expand ML models to support industry-specific transformation rules (e.g., finance, healthcare).
5. **Generative AI:** Explore the use of large language models to automate transformation rule generation and semantic data enrichment.

REFERENCES

1. Rahm, E., & Do, H. H. (2000). *Data cleaning: Problems and current approaches*. IEEE Data Eng. Bull.
2. Chu, X., Ilyas, I. F., Krishnan, S., & Wang, J. (2016). *Data Civilizer System*. CIDR.
3. Rekatsinas, T., Chu, X., Ilyas, I. F., & Ré, C. (2017). *HoloClean: Holistic Data Repairs with Probabilistic Inference*. VLDB.
4. AWS Glue Documentation. <https://docs.aws.amazon.com/glue/>
5. Google Cloud DataPrep. <https://cloud.google.com/dataprep>
6. Azure Data Factory. <https://azure.microsoft.com/en-us/services/data-factory/>
7. PyCaret Documentation. <https://pycaret.org/>
8. Zhang, Y., Chen, L., & Wang, W. (2020). *Automatic Data Cleaning Techniques: A Survey*. ACM Computing Surveys.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



International Journal of Multidisciplinary and Scientific Emerging Research (IJMSERH)

Impact Factor: 9.274

✉ ijmserh@gmail.com

🌐 www.ijmserh.com