

ETL Strategies for Large-Scale Retail Data Warehouses

Ravi Kiran Koppichetti

koppichettiravikiran@gmail.com

Abstract

Large-scale retail data warehouses are critical for storing and analyzing vast amounts of transactional, operational, and customer data. Effective ETL (Extract, Transform, Load) strategies are essential for ensuring that data is accurately extracted from diverse sources, transformed into a usable format, and loaded into the data warehouse for analysis. This paper explores the challenges of implementing ETL processes in large-scale retail data warehouses and provides strategies for optimizing ETL workflows. Key topics include data integration, scalability, performance optimization, and the use of modern ETL tools and technologies. The paper concludes with recommendations for designing robust ETL pipelines that meet the demands of the retail industry.

Keywords: ETL, Extract Transform Load, Data Warehouse Architectures, Data Warehouse Data Mining, Machine Learning, Artificial Intelligence, Data Science, Data Preprocessing, Data Pre-Processing, Data Preparation

I. Introduction

The retail industry is undergoing a significant transformation due to the rapid adoption of digital technologies and a focus on data-driven strategies. Retailers operate in a competitive environment where understanding customer behavior, optimizing supply chains, and responding swiftly to market changes are crucial. To meet these goals, they utilize data warehouses that collect and organize information from sources such as POS systems, e-commerce platforms, and CRM systems. These warehouses support analytics, providing insights for business growth.

However, establishing and maintaining these data warehouses is complex. Retailers must aggregate data from various sources, convert it into a usable format, and load it for analysis—a process known as ETL (Extract, Transform, Load). Effective ETL strategies ensure data accuracy and accessibility, enabling informed decisions and quick adjustments to market dynamics.

Nonetheless, executing ETL processes in large retail data warehouses presents challenges. Retailers handle vast volumes of structured and unstructured data, making consolidation into a uniform format intricate. The need for real-time data processing adds complexity, as traditional batch ETL methods may not meet the demands of scenarios like personalized marketing and fraud detection.

Data quality is another key concern in ETL. Poor data quality can lead to misguided insights and faulty decisions, necessitating robust validation, cleansing, and profiling techniques. As data volumes grow, scalable and efficient ETL processes are essential to prevent bottlenecks that hinder decision-making.

To tackle these issues, retailers are increasingly adopting modern ETL tools. Cloud ETL platforms like AWS Glue, Google Dataflow, and Azure Data Factory offer scalable solutions, while streaming tools like Apache Kafka enable real-time processing for immediate insights. The shift to an ELT model, where data is loaded first before transformation, optimizes complex transformations in modern systems like Snowflake.

This paper discusses the challenges of implementing ETL in retail data warehouses and offers strategies to improve workflows, including unified data models and parallel processing. A case study of a retail chain showcases the application of these strategies, demonstrating their positive impact on efficiency and customer experience. Best practices for constructing robust ETL pipelines that meet retail demands are also presented.

By utilizing effective ETL strategies, retailers can maximize their data warehouses, enabling real-time analytics, boosting efficiency, and enhancing customer experiences. As the retail landscape evolves, ETL processes will be critical in fostering innovation and growth, helping retailers build resilient and scalable pipelines that maintain their competitiveness.

II. Challenges in ETL for Large-Scale Retail Data Warehouses

Implementing ETL (Extract, Transform, Load) processes within large-scale retail data warehouses is a complex and multifaceted endeavor. Retailers must navigate various challenges to guarantee that data is accurately extracted, transformed, and subsequently loaded into the warehouse for analysis. The following discussion will delve into these challenges in detail, emphasizing their implications for ETL workflows and their broader impact on retail operations.

A. Data Volume and Variety: Retailers accumulate a substantial volume of data from various sources, including point-of-sale (POS) systems, e-commerce platforms, customer relationship management (CRM) systems, Internet of Things (IoT) devices, and social media. This information encompasses a diverse range of types, comprising structured data such as transactional records, semi-structured data like JSON logs from web applications, and unstructured data, including customer reviews and images. However, the management of this extensive array of data presents significant challenges for Extract, Transform, Load (ETL) processes. Conventional ETL tools and systems frequently encounter difficulties in accommodating the scale and complexity inherent in retail data, which can result in performance issues and delays in accessing critical information. For instance, a considerable retail chain with multiple locations and a robust e-commerce platform may generate terabytes of data daily, incorporating transaction logs, customer interactions, and updates from the supply chain. Therefore, to process this data efficiently and in a timely manner, the implementation of scalable and adaptable ETL solutions is indispensable. Ineffective management of data volume and variety may result in incomplete or inaccurate data within the data warehouse, thereby adversely affecting the reliability of analytical and decision-making processes [1, 2, 3, 4, 5].

B. Data Integration: Retailers are undertaking the significant endeavor of integrating data from a diverse array of sources, each characterized by its distinct format, structure, and semantics. For

instance, while data sourced from Point of Sale (POS) systems may be systematically organized within relational databases, data extracted from e-commerce platforms could be presented in JSON or XML formats. The consolidation of this data into a cohesive and uniform format presents a considerable challenge. A primary obstacle involves ensuring data consistency and compatibility across these varied sources during the Extract, Transform, Load (ETL) processes. To address this issue, retailers can develop robust data integration strategies that not only harmonize the data but also rectify any inconsistencies that may arise. For example, a retailer may find itself integrating data from in-store POS systems, online transactions, and mobile applications. Given that each of these sources may utilize different identifiers for customers and products, this necessitates careful mapping and transformation logic to ensure seamless integration during the ETL process. The significance of effective data integration cannot be overstated; inadequate integration can result in data silos, wherein information becomes fragmented and difficult to access. This fragmentation can severely impede a retailer's capacity to conduct comprehensive analyses and derive valuable insights [1, 6, 7, 8].

- C. Performance and Scalability:** As the volume of data continues to expand, it is imperative for ETL processes to adapt and manage the increasing workloads efficiently without compromising performance. Retailers frequently encounter performance challenges during the extraction, transformation, and loading phases, particularly when dealing with large datasets or complex transformations. The primary challenge is to design ETL processes that are both efficient and scalable. Retailers are advised to optimize their ETL workflows to minimize processing time and resource consumption while ensuring that their systems can accommodate the continuously growing data volumes. For example, during peak shopping seasons such as Black Friday or holiday sales, retailers witness a considerable surge in transaction data. It is essential for ETL processes to effectively manage this influx, ensuring that there are no delays or failures. When performance bottlenecks arise, they can postpone the availability of data, which may impede real-time analytics and decision-making. Furthermore, issues related to scalability may result in elevated infrastructure costs and operational inefficiencies, circumstances that no retailer wishes to confront [1, 9, 10].
- D. Data Quality:** Ensuring the accuracy, completeness, and consistency of our data is of paramount importance for reliable analytics. When data quality declines, it may result from errors in source systems, discrepancies among sources, or complications during the ETL process. A significant challenge lies in maintaining high data quality throughout this process. To address this, retailers can employ comprehensive data validation, cleansing, and profiling techniques to identify and rectify any issues. For example, a retailer may encounter problems such as duplicate customer records, incomplete product information, or inconsistent pricing data. It is essential to resolve these issues during the ETL process to ensure that the data warehouse is populated with precise and trustworthy information. Should data quality deteriorate, it could result in misleading insights and suboptimal business decisions, thereby undermining the value of the data warehouse [6, 11, 12].

- E. Real-Time Processing:** Retailers are increasingly recognizing the necessity for real-time or near-real-time data processing to enhance personalized marketing, dynamic pricing, and fraud detection. Traditional batch-based ETL (Extract, Transform, Load) processes, which process data at scheduled intervals, frequently do not suffice in addressing these evolving demands. The primary challenge resides in the implementation of real-time ETL processes, which necessitate extremely low-latency data processing and high availability. To facilitate seamless real-time data integration and analysis, retailers should consider adopting streaming ETL tools and methodologies. For example, consider a retailer that needs to identify and respond to fraudulent transactions immediately; this requires the prompt processing of transaction data. Likewise, personalized marketing campaigns benefit significantly from real-time analysis of customer behavior, thereby enabling the effective delivery of tailored promotions. If retailers are unable to support real-time processing capabilities, they may encounter difficulties in swiftly responding to fluctuating market conditions and customer needs. This inadequacy could result in a decline in competitiveness and overall customer satisfaction, which is certainly an undesirable outcome [13, 14].
- F. Legacy Systems and Technical Debt:** Numerous retailers encounter challenges associated with legacy systems that were not designed to accommodate the requirements of contemporary ETL processes. These systems frequently exhibit deficiencies in flexibility, scalability, and integration—qualities essential for the effective management of extensive data warehousing. A significant obstacle lies in the integration of these older systems with modern ETL tools and technologies. To bridge this gap, retailers often resort to middleware, APIs, or custom connectors. For example, a retailer may operate an antiquated point-of-sale (POS) system that maintains data in a proprietary format, necessitating the application of custom ETL logic to extract and transform that data for the data warehouse. The reality is that legacy systems and the associated burden of technical debt can complicate ETL processes and increase costs, thereby hindering the implementation of efficient data warehousing solutions [15, 16].

III. ETL Strategies for Large-Scale Retail Data Warehouses

To tackle the challenges of deploying Extract, Transform, Load (ETL) processes in large retail data warehouses, it is vital for retailers to adopt various strategies focused on improving data integration, scalability, performance, and quality. These approaches are essential for ensuring accurate extraction, transformation, and loading of data into the warehouse, thus enabling retailers to gain actionable insights and support business growth. The following discussion delves into these strategies in detail, emphasizing their benefits and practical applications. Effectively implementing ETL strategies is crucial for creating robust and scalable data warehouses in the retail sector. By utilizing methods like unified data models, parallel processing, streaming ETL, and modern ETL tools, retailers can effectively navigate the challenges of high data volumes, integration complexities, and real-time processing needs. These techniques help retailers maintain data accuracy, boost performance, and adapt to dynamic use cases, thereby fostering innovation and encouraging business expansion. The next section of this paper shares a case study of a leading retail chain that successfully embraced these strategies, providing valuable insights and lessons learned.

A. Data Integration Strategies: A crucial aspect of ETL is integrating data from various sources into a unified format. Retailers manage data from POS systems, e-commerce platforms, CRM systems, and IoT devices, each having unique structures. A unified data model simplifies this process by offering a common framework for integration. For instance, retailers may develop a dimensional model, such as a star schema, where transactional data (sales) resides in fact tables and descriptive data (products, customers) is placed in dimension tables, ensuring consistency and improving query performance.

API-based integration effectively enables real-time data extraction from systems such as e-commerce platforms and cloud applications. APIs empower retailers to instantly extract data from source systems and update the warehouse. Retailers can utilize APIs to collect real-time sales data from online stores for immediate analysis, reducing latency and supporting dynamic scenarios like personalized marketing and fraud detection.

For retailers managing large datasets, Change Data Capture (CDC) techniques minimize the data processed during ETL. CDC only captures alterations in source data, such as new or updated records, preventing the need to reprocess the entire dataset. Retailers can utilize CDC to update customer records from their CRM, ensuring the data warehouse remains current without overloading the ETL process, thereby enhancing performance, reducing resource consumption, and preserving data accuracy [1, 6, 7, 8].

B. Scalability and Performance Optimization: As data volume grows, ETL (Extract, Transform, Load) processes must scale effectively without losing performance. A key strategy is parallel processing, which distributes workloads across multiple nodes, allowing retailers to process large datasets efficiently. For example, retailers with multiple locations can use Apache Spark to process sales data in parallel, significantly speeding up data loading into warehouses and managing increasing data volumes.

Another optimization method is incremental loading, which involves adding only new or updated data to the warehouse. Instead of reprocessing the entire dataset, retailers can extract changes using timestamps or Change Data Capture (CDC). For instance, loading only current-day sales transactions minimizes processing time and resource use while keeping the data warehouse up-to-date.

Data partitioning also enhances performance and scalability. By dividing large datasets into smaller segments based on key attributes (like date or region), query performance improves, and processing time decreases. For example, a retailer might partition sales data by month, resulting in faster queries and better data processing, especially for those handling vast datasets [1, 9, 10].

C. Data Quality Management: Data quality is crucial in the ETL process, as poor data can lead to faulty insights and poor decisions. Data validation is key for ensuring accuracy, with retailers using rules to detect errors like missing values. For example, validating sales data checks that

transactions match valid product and customer IDs. Automated tools like Great Expectations help transfer high-quality data to the warehouse.

Data cleansing is also vital. It involves error identification and correction, such as removing duplicates and fixing inconsistencies. Retailers might use tools like Trifacta or OpenRefine to clean customer data, remove duplicates, and standardize addresses, enhancing analytics reliability.

Before loading data into the warehouse, retailers should profile the data to evaluate its structure and quality. Tools like Talend or Informatica can identify issues such as missing attributes or inconsistent pricing, steering ETL workflows, and resolving data quality concerns before they impact the data warehouse [6, 11, 12].

D. Real-Time and Near-Real-Time ETL: In today's fast-paced retail landscape characterized by rapid dynamism, the capability for real-time or near-real-time data processing is imperative for facilitating adaptive applications such as personalized marketing, dynamic pricing strategies, and fraud detection mechanisms. Streaming Extract, Transform, Load (ETL) constitutes a robust methodology for accomplishing this objective. Technologies such as Apache Kafka, Apache Flink, and AWS Kinesis empower retailers to process data as it is generated, thus furnishing immediate insights and enabling prompt responses. For instance, a retailer may leverage Kafka to analyze real-time transaction data, thereby facilitating instantaneous fraud detection and adjustments to dynamic pricing.

For retailers striving to equilibrate real-time processing with system efficacy, micro-batching presents a viable alternative. This methodology entails the processing of data in small, recurrent batches, thereby mitigating latency while preventing system overload. For example, a retailer might execute sales data processing in batches of five minutes using Apache Spark Streaming or Google Dataflow, thereby providing near-real-time insights while preserving system performance. This approach is especially advantageous for retailers managing substantial data volumes or operating under constrained infrastructure capabilities [13, 14].

E. Modern ETL Tools and Technologies: The adoption of contemporary ETL (Extract, Transform, Load) tools and technologies is imperative for the development of scalable and efficient data pipelines. Cloud-based ETL platforms, including AWS Glue, Google Dataflow, and Azure Data Factory, present flexible and cost-effective solutions for data processing. These platforms utilize cloud infrastructure to deliver scalability, enhanced performance, and user-friendly interfaces. For instance, a retailer may employ AWS Glue to process and load data from multiple sources into a cloud-based data warehouse, thereby reducing infrastructure costs while enhancing scalability.

Another notable trend is the ELT methodology, which entails loading data into the data warehouse prior to executing transformations. This approach capitalizes on the capabilities of advanced data warehouses such as Snowflake and BigQuery, which are adept at efficiently

handling complex transformations through SQL or built-in functions. For example, a retailer could load raw sales data into Snowflake and execute transformations via SQL queries, thereby streamlining the ETL workflow and enhancing overall performance.

Retailers may also leverage data orchestration tools like Apache Airflow and Prefect to automate and oversee intricate ETL workflows. These tools permit retailers to define workflows as directed acyclic graphs (DAGs), ensuring that task dependencies are met and providing clarity on the ETL process. For instance, a retailer may utilize Apache Airflow to automate the ETL process for daily sales data, thereby ensuring that tasks are executed in the appropriate sequence and enabling real-time monitoring of the workflow [15, 16, 17, 18, 19].

IV. Case Study: Implementing ETL in a Large Retail Chain

To illustrate the practical application of ETL strategies, let's examine the case of a large retail chain with thousands of stores and a significant e-commerce presence. The retail chain faced several challenges in managing its data, including integrating data from disparate sources, ensuring data quality, and supporting real-time analytics. By adopting a range of ETL strategies, the Retail chain successfully built a robust data pipeline that transformed its operations and enabled data-driven decision-making. The retail chain data ecosystem was highly fragmented, with data coming from a variety of sources, including in-store POS systems, online transactions, mobile apps, and supply chain operations. Each source used different formats, structures, and identifiers, making it difficult to integrate data into a unified data warehouse. Additionally, the sheer volume of data generated by thousands of stores and millions of customers posed significant scalability and performance challenges. The retail chain also needs to support real-time analytics for dynamic use cases, such as personalized marketing, dynamic pricing, and fraud detection. Finally, ensuring data quality was a top priority, as poor-quality data could lead to incorrect insights and flawed business decisions.

- A. **Unified Data Model:** The retail chain established a unified data model founded on a star schema to standardize data from various sources. The model encompassed fact tables for transactional data (e.g., sales, returns) and dimension tables for descriptive data (e.g., products, customers, stores). This methodology ensured consistency throughout the data warehouse and facilitated data integration. For instance, sales data from in-store Point of Sale (POS) systems and online transactions were mapped to a common schema, thereby enabling seamless analysis across different channels.
- B. **Streaming ETL:** To facilitate real-time analytics, the retail chain has implemented streaming ETL utilizing Apache Kafka. The real-time transaction data obtained from point-of-sale (POS) systems and e-commerce platforms is processed through Kafka, thereby enabling immediate insights and responses. For instance, the retail chain employs streaming ETL to detect and prevent fraudulent transactions in real time, which significantly reduces financial losses and enhances customer trust.
- C. **Cloud-Based ETL:** The retail chain utilized cloud-based ETL platforms, including AWS Glue, to design and manage its ETL workflows. Additionally, the organization embraced an ELT

(Extract, Load, Transform) methodology, which involved loading raw data into the data warehouse prior to executing transformations. This method capitalized on the capabilities of Snowflake, the retail chain's cloud-based data warehouse, to efficiently manage complex transformations utilizing SQL.

- D. **Data Quality Management:** The retail chain emphasized the importance of data Quality. During the Extract, Transform, Load (ETL) process, it established strong data validation rules to detect errors, including but not limited to missing values and inconsistent formats. Automated tools were used to clean customer data, remove duplicates, and standardize address formats. Additionally, the chain performed data profiling to examine the structure and content of the source data, effectively identifying potential issues before they affected the data warehouse.
- E. **Parallel Processing and Incremental Loading:** To manage the substantial volume of data produced by its retail locations, the retail chain adopted parallel processing using Apache Spark. Sales data from thousands of stores was processed simultaneously, significantly reducing the time needed for data to load into the data warehouse. Additionally, the retail chain implemented incremental loading to update the data warehouse only with new or modified records. For instance, daily sales transactions were loaded incrementally, minimizing both processing time and resource consumption.

The successful implementation of Extract, Transform, Load (ETL) strategies by the retail chain illustrates the significant impact of effective data integration and processing within the retail sector. By adopting a unified data model, utilizing modern tools, and prioritizing data quality, the retail chain established a robust data pipeline that enabled real-time analytics, increased operational efficiency, and enhanced customer experiences. This case study provides valuable insights and best practices for other retailers looking to optimize their ETL processes and fully harness the potential of their data.

V. Optimal Strategies for ETL in Retail Data Warehousing

Implementing effective ETL (Extract, Transform, Load) processes in retail data warehouses requires careful planning, robust strategies, and adherence to best practices. Retailers must ensure that their ETL workflows are scalable, performant, and capable of handling the industry's unique challenges. Below, we explore the best practices for designing and managing ETL processes in retail data warehouses, providing actionable insights and practical examples [20, 21].

- A. **Scalability in Design:** Retailers generate substantial volumes of data from point-of-sale systems, e-commerce platforms, and Internet of Things devices. To manage the increasing data loads effectively, it is advisable to employ distributed processing frameworks such as Apache Spark to parallelize Extract, Transform, and Load (ETL) tasks. For instance, sales data from thousands of stores can be processed in parallel to minimize load times.
- B. **Prioritizing Data Quality:** Implementing validation, cleansing, and profiling techniques is crucial to ensuring data accuracy and consistency. Use tools like Great Expectations or Trifacta

to automate data quality verification processes. For instance, it is essential to validate sales data for missing product IDs or cleanse customer data to eliminate duplicates.

- C. Utilize Contemporary ETL Tools:** To achieve scalable and cost-efficient data processing, embrace cloud-based ETL solutions such as AWS Glue, Google Dataflow, or Azure Data Factory. These tools facilitate streamlined ETL workflows and integrate effortlessly with cloud data warehouses, including Snowflake and BigQuery.
- D. Implement an ELT Approach:** Ingest raw data into the data warehouse prior to executing transformations (ELT). This methodology capitalizes on the advanced processing capabilities of contemporary data warehouses for intricate transformations. For instance, employ SQL queries within Snowflake to transform raw sales data with enhanced efficiency.
- E. Implement Real-Time Data Processing:** Employ streaming Extract, Transform, Load (ETL) tools such as Apache Kafka or AWS Kinesis for the purpose of real-time data processing. This approach facilitates dynamic applications, including fraud detection and personalized marketing strategies. For situations requiring near-real-time processing, it is advisable to consider micro-batching methodologies utilizing tools such as Apache Spark Streaming.
- F. Automate Extract, Transform, Load (ETL) Workflows:**
To effectively automate and manage ETL workflows, employ orchestration tools such as Apache Airflow or Prefect. Define workflows as Directed Acyclic Graphs (DAGs) to ensure that tasks are executed in the appropriate sequence while facilitating real-time performance monitoring.
- G. Ensure Data Security and Compliance:** Protect sensitive customer information by implementing strong encryption protocols, strict access controls, and comprehensive audit logs. Maintain compliance with regulatory frameworks such as the General Data Protection Regulation (GDPR) and the California Consumer Privacy Act (CCPA). For example, encrypt customer data during both transmission and storage to reduce the risk of unauthorized access.
- H. Monitor and Optimize Performance:** Consistently oversee the performance of Extract, Transform, and Load (ETL) processes by utilizing monitoring tools such as Datadog or Prometheus. Assess key metrics, including processing time and error rates, and refine workflows to minimize latency while enhancing overall efficiency.
- I. Strategic Plan for Future Growth:** Develop Extract, Transform, Load (ETL) processes that are adaptable to incorporate emerging data sources and changing business requirements. Implement modular architectures conducive to seamless updates, facilitating the integration of the Internet of Things (IoT) and social media data without necessitating a comprehensive system overhaul.
- J. Encourage Collaborative Efforts:** Promote collaboration among data engineers, data scientists, and business stakeholders to ensure the alignment of ETL processes with organizational

objectives. Utilize platforms such as Confluence or Slack to enhance communication and guarantee the usability of data.

VI. Conclusion

ETL processes are the backbone of large-scale retail data warehouses, enabling retailers to consolidate and analyze data from diverse sources. However, implementing ETL in the retail industry presents unique challenges, including data volume, integration complexity, and the need for real-time processing. By adopting strategies such as unified data models, parallel processing, streaming ETL, and modern ETL tools, retailers can overcome these challenges and build robust data pipelines.

The case study of a large retail chain demonstrates the practical application of these strategies, highlighting the importance of scalability, data quality, and real-time processing. As the retail industry continues to evolve, ETL processes will play an increasingly critical role in enabling data-driven decision-making and driving business growth. By following best practices and leveraging emerging technologies, retailers can ensure that their ETL workflows meet the demands of a dynamic and competitive market.

References

- [1] P. S. Diouf, A. Boly, and S. Ndiaye, "Variety of data in the ETL processes in the cloud: State of the art," in *2018 IEEE International Conference on Innovative Research and Development (ICIRD)*, May 2018, pp. 1–5.
- [2] N. Berkani, L. Bellatreche, and L. Guittet, "ETL processes in the era of variety," *Transactions on Large-Scale Data-and Knowledge-Centered Systems XXXIX: Special Issue on Database-and Expert-Systems Applications*, pp. 98–129, 2018.
- [3] S. M. F. Ali, "Next-generation ETL framework to address the challenges posed by big data," in *DOLAP*, Mar. 2018.
- [4] S. K. Bansal, "Towards a semantic extract-transform-load (ETL) framework for big data integration," in *2014 IEEE International Congress on Big Data*, Jun. 2014, pp. 522–529.
- [5] S. Souissi and M. BenAyed, "Genus: An ETL tool treating the big data variety," in *2016 IEEE/ACS 13th International Conference of Computer Systems and Applications (AICCSA)*, Nov. 2016, pp. 1–8.
- [6] V. Köppen, B. Brüggemann, and B. Berendt, "Designing data integration: The ETL pattern approach," *UPGRADE: The European Journal for the Informatics Professional*, no. 3, pp. 49–55, 2011.
- [7] J. Sreemathy, S. Nisha, and G. P. RM, "Data integration in ETL using TALEND," in *2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS)*, Mar. 2020, pp. 1444–1448.

- [8] J. Sreemathy et al., “Overview of ETL tools and Talend-data integration,” in *2021 7th International Conference on Advanced Computing and Communication Systems (ICACCS)*, vol. 1, Mar. 2021, pp. 1650–1654.
- [9] P. Badgujar, “Optimizing ETL processes for large-scale data warehouses,” *Journal of Technological Innovations*, vol. 2, no. 4, 2021.
- [10] M. Bodziony, S. Roszyk, and R. Wrembel, “On evaluating performance of balanced optimization of ETL processes for streaming data sources,” in *DOLAP*, 2020, pp. 74–78.
- [11] S. B. Dakrory, T. M. Mahmoud, and A. A. Ali, “Automated ETL testing on the data quality of a data warehouse,” *International Journal of Computer Applications*, vol. 131, no. 16, pp. 9–16, 2015.
- [12] R. Singh and K. Singh, “A descriptive classification of causes of data quality problems in data warehousing,” *International Journal of Computer Science Issues (IJCSI)*, vol. 7, no. 3, pp. 41, 2010.
- [13] X. Li and Y. Mao, “Real-time data ETL framework for big real-time data analysis,” in *2015 IEEE International Conference on Information and Automation*, Aug. 2015, pp. 1289–1294.
- [14] H. Zhou, D. Yang, and Y. Xu, “An ETL strategy for real-time data warehouse,” in *Practical Applications of Intelligent Systems: Proceedings of the Sixth International Conference on Intelligent Systems and Knowledge Engineering (ISKE2011)*, Dec. 2011, pp. 329–336.
- [15] S. Jha, *A big data architecture for integration of legacy systems and data*, Doctoral dissertation, CQUniversity, 2021.
- [16] M. Guerrero, M. Segura, and J. Lucio, “Proposal of a framework for information migration from legacy applications in solidarity financial sector entities,” in *International Conference on Systems and Information Sciences*, Jul. 2020, pp. 309–320.
- [17] M. Patel and D. B. Patel, “Progressive growth of ETL tools: A literature review of past to equip future,” in *Rising Threats in Expert Applications and Solutions: Proceedings of FICR-TEAS 2020*, pp. 389–398.
- [18] J. Goldfedder and J. Goldfedder, “Choosing an ETL tool,” in *Building a Data Integration Team: Skills, Requirements, and Solutions for Designing Integrations*, pp. 75–101, 2020.
- [19] A. Karagiannis, P. Vassiliadis, and A. Simitsis, “Scheduling strategies for efficient ETL execution,” *Information Systems*, vol. 38, no. 6, pp. 927–945, 2013.
- [20] H. Gadde, “AI-enhanced data warehousing: Optimizing ETL processes for real-time analytics,” *Revista de Inteligencia Artificial en Medicina*, vol. 11, no. 1, pp. 300–327, 2020.



[21] A. Simitsis, P. Vassiliadis, and T. Sellis, “Optimizing ETL processes in data warehouses,” in *21st International Conference on Data Engineering (ICDE'05)*, Apr. 2005, pp. 564–575.