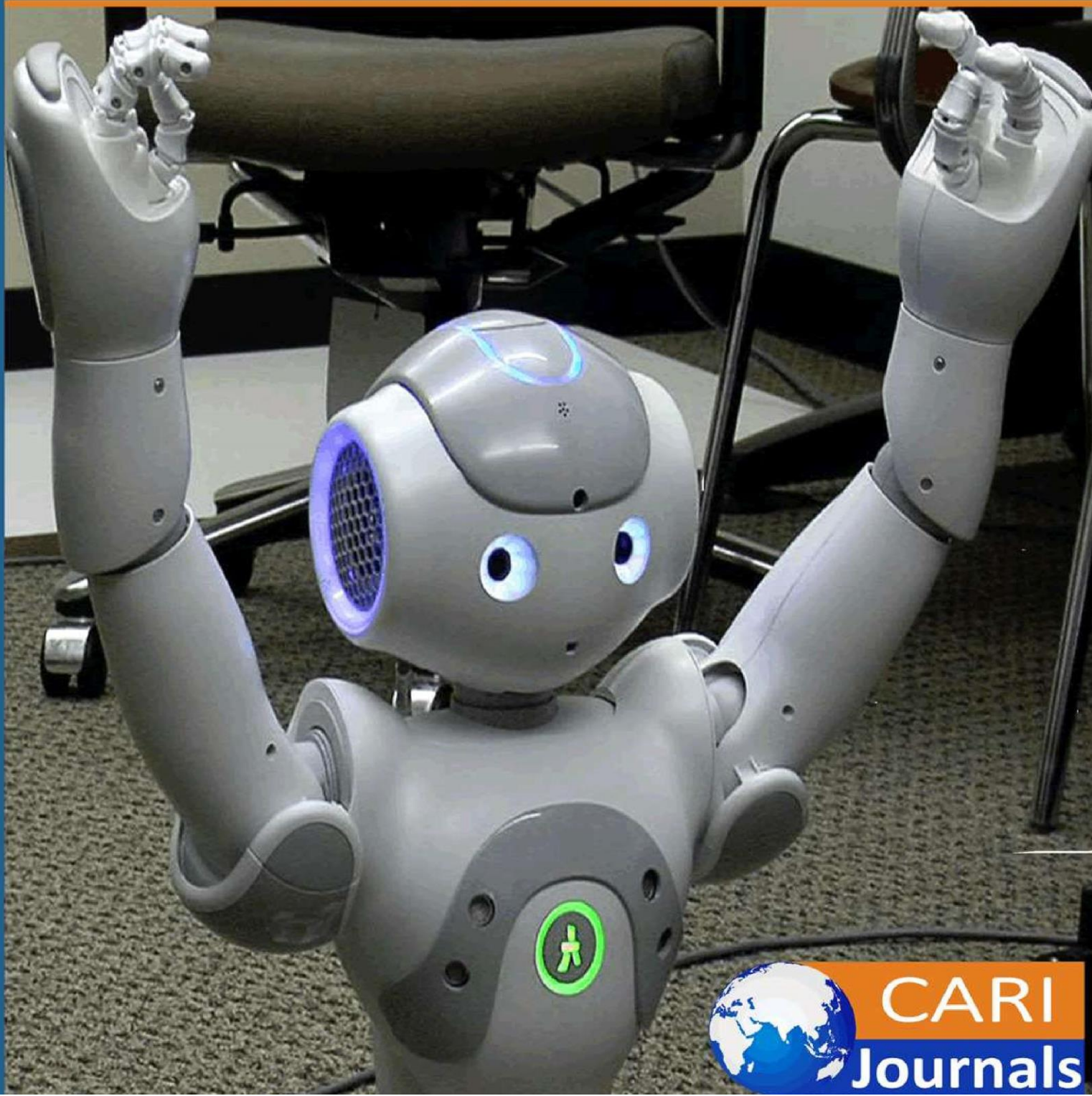


International Journal of Computing and Engineering

(IJCE) **Ensuring Data Security and Compliance in ETL Processes for
Healthcare and Financial Services**



**CARI
Journals**

Ensuring Data Security and Compliance in ETL Processes for Healthcare and Financial Services



Santosh Kumar, Singu

Senior Solution Specialist, Deloitte Consulting LLP

<https://orcid.org/0009-0001-8954-6776>

Accepted: 30th Oct, 2024, Received in Revised Form: 14th Nov, 2024, Published: 23rd Nov, 2024

Abstract

Purpose: This study explores the critical role of ETL (Extract, Transform, Load) processes in managing data within highly regulated industries such as healthcare and finance. It highlights the challenges posed by stringent legal frameworks like HIPAA in the health sector and GDPR in finance while emphasizing the importance of ETL in ensuring compliance and improving data utility.

Methodology: The research examines the application of ETL processes in consolidating and transforming data for organizational use, with specific examples from healthcare (e.g., Electronic Health Records) and financial sectors (e.g., Basel III reporting). It also reviews current best practices for addressing data-related challenges, including governance, encryption, validation, and containerization.

Findings: Key challenges in ETL processes include data privacy, regulatory compliance, data quality issues, and technical limitations. However, implementing best practices such as robust data governance, advanced encryption methods, intelligent validation mechanisms, and containerized workflows significantly mitigates these risks. These practices ensure secure data handling and enhance organizational compliance with regulatory standards.

Unique Contribution to Theory, Policy and Practice: The study contributes to the theoretical understanding of ETL processes as a linchpin for data management in regulated environments. It offers policy insights into how organizations can meet compliance requirements effectively. Practically, it provides actionable recommendations for organizations to adopt ETL best practices, ensuring secure, efficient, and legally compliant data operations. These advancements strengthen client trust, reduce legal risks, and empower organizations to leverage data for strategic advantage.

Keywords: *ETL Process, Data Privacy, Healthcare Regulation and Financial Regulation, Data Control.*

Introduction

Some of the compliance measures were used to ensure the security of the data, especially in industries that have maximum regulation, such as the health sector and financial sectors. ETL, or extract, transform, and load, is an essential factor in data management concerning such sectors, as it compiles data from numerous sources and formats into another usable, accessible reporting, analysis, and storage style. However, the advanced ETL processes open a set of new and rather delicate security and compliance issues, such as data protection, compliance with regulations and laws, and technical problems [3]. To this end, this essay responds to the following research questions: This essay discusses a variety of stakeholders differentiated by the unique security and compliance requirements in the ETL frameworks of healthcare and financial services, identifies the significant challenges and offers suggestions on how to enhance the ETL procedures to ensure compliance with strict data security and regulation standards in the respective industries.

The Importance of ETL in Healthcare and Financial Services

From a business point of view, ETL procedures are most valuable for processes in the healthcare and the financial industry. In healthcare, ETL is required to convert different EHR systems and clinical information into CDM for analysis, data presentation and the enhancement of patient care. Financial institutions apply ETL approaches to process data to meet various regulatory reports to ETL to provide data for various compliance purposes, such as Compliance with Basel III Liquidity to Coverage ratio, which consumes daily processed data [4].

How ETL Process Works

The ETL's full form is Extract, Transform, and Load, a common strategy used in data management and in data warehousing and analytics. This process helps in the management of the data obtained from various sources in a format that can easily be analyzed. All ETL steps are essential to provide data accuracy, coherency, and availability so that decision-making is based on actual information within an organization. The ETL process is typically divided into three sequential phases: Extraction, Transformation, and loading.

The First step **Extraction** involves gathering information from multiple sources, such as relational databases, APIs, or flat files. Data in today's digital environment can be in one format or another, from one system or several systems within or even from outside an organization. The extraction process brings together data in this form: It may need to be revised, questionable as to accuracy, or not in the appropriate format for analysis required for the data analysis objectives [1]. This phase is the most important because it establishes the kind of data quality that should be expected. Sometimes, because of inaccuracy or lack of data at this stage, the ETL process affects and delivers flawed insights at the later stage. Indeed, the Extraction is usually followed by data validation stages to have the least number of errors in the first data set.

After Extraction, the data moves **to Transformation**, where data is processed and made ready for use in a particular form. The transformation phase is the most challenging since it includes tasks such as data cleansing, selection, and condensing. Cleaning would be defining,

identifying, invalidating and eradicating redundancies, lacuna and inaccuracies in data to improve its quality [2]. Filtering eliminates extraneous data from the discovered sources to apply an interest analysis. At the same time, aggregation transforms data from different sources into a similar form. This phase also ensures that data required within a specific target data repository has been formatted appropriately to be compatible and functional. Transformation prepares the data to make it equally valuable per a given analytical purpose and business objective.

The final step in the ETL process is **Loading**. In this phase, the transformed data gets to its right destination, such as a data warehouse, data lake or cloud storage. It is a data storage and management centre where the data is consolidated and can then be processed and utilized for reports. Depending on the organization's needs, loading can be done in large volumes at some point or in gradual batches constantly [5]. The load step is compassionate because it helps those who make decisions get fresh and accurate information required for analysis, reporting, and other strategic activities.

See the diagram below:

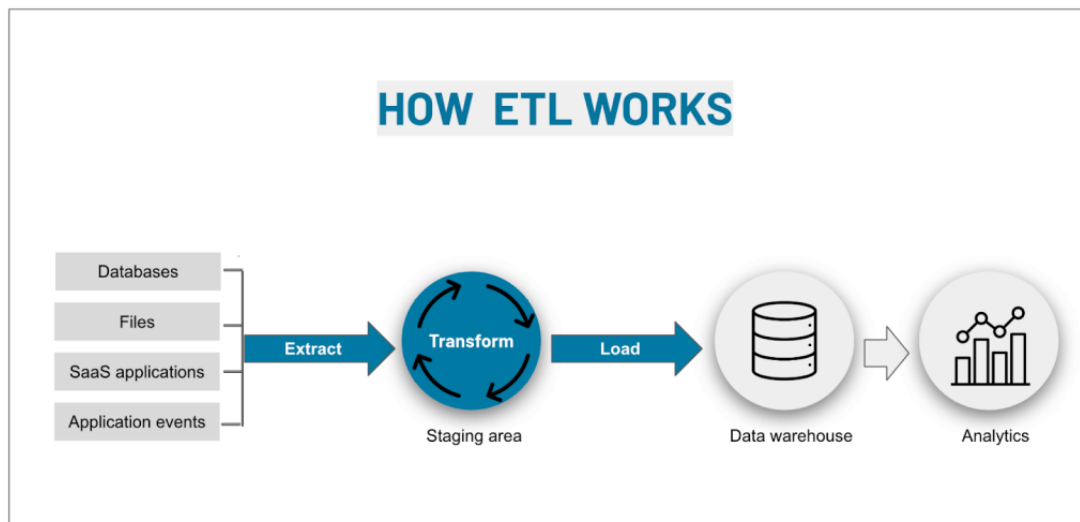


Figure 1

ETL Integration

In healthcare, ETL integration is essential for gathering data from different sources for analysis, reporting, and use of the resulting data. Based on the provided image, we can break down the ETL integration process into three main areas. The three key categories in health informatics include Data Acquisition, Data Storage, and Information Delivery.

Data Acquisition Area: Ultimately, this phase involves internal data (HE rea, operational databases, hospitals, etc.) and external data (public health databases, health insurance claims, third parties, etc.). In health care, this source can be numerous and include clinical data, patient data, lab data, and bills, among others. Data acquisition means all health data must reach the point of integration to create the proper framework [7].

Data Staging and ETL Process: The ETL process occurs in the Data Staging subarea. Here, data is taken from its source and then passes through the process of extraction transformation and Transformation, which is staged. In the Transform phase of European Healthcare Data Lake, data such as Health care data is groomed and structured correctly to make them homogeneous in various forms and places. For example, separate hospitals may employ somewhat distinguishing coding systems for diagnoses or treatment so that data transformation will harmonize those discrepancies into a consistent standard, for instance, ICD codes or HL7 [6]. This process is critical in healthcare settings for attaining data interoperability, where data analysis is done consistently across facilities or units.

Data Storage Area: Transformed data is stored in a Data Warehouse or a set of better-sized repositories called Data Marts. The Data Warehouse is an aggregate where firms can holistically store and consolidate a colossus of healthcare information that would later be used for retrieving and analyzing data. [metadata] is also saved here; it contains the data's structural, historical and geometric characteristics. On the other hand, a data mart is a small part of the data warehouse that focuses more on information and provides solutions to a specific need, whether clinical data results or patient satisfaction ratings [8]. Such a structure allows users, particularly extended individuals, healthcare providers, and researchers, to quickly obtain the needed specifics for analysis.

Information Delivery Area: The last tier, known as an OLAP (Online Analytical Processing) Server, forwards the processed data to the end consumers in the form of dashboards, reports or analytical tools. These tools help healthcare administrators, clinicians, and researchers to develop awareness, predict, and audit patient results, therapeutic effectiveness, and organizational execution. For instance, an expert in healthcare delivery might employ OLAP-based reports to analyze cases of patient readmission or outbreaks of some epidemics.

In conclusion, ETL integration in healthcare guarantees that health data obtained from several sources is efficiently matched, processed, and prepared for analysis in the shortest time possible [9]. They help to recognize and promote a full range of patient care, enhance overall operational efficiency in the delivery of services, and facilitate public health surveillance and early detection of health threats. See the figure below:

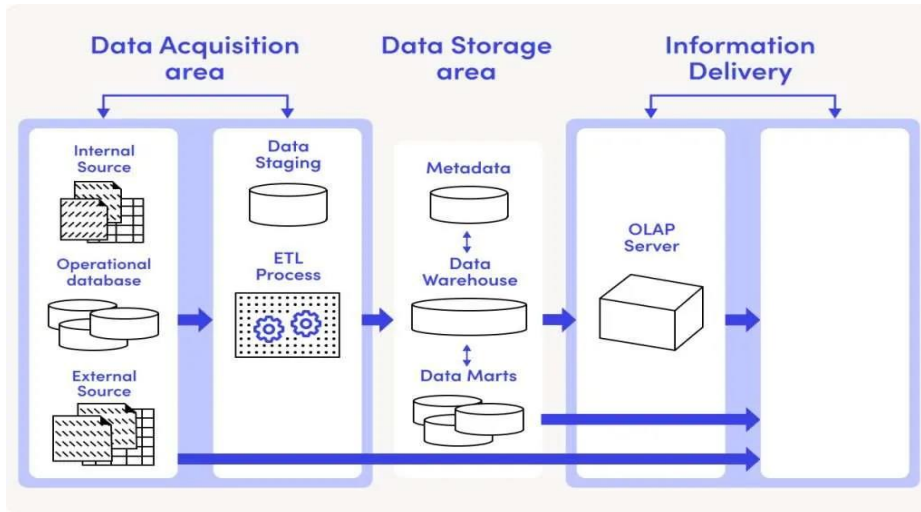


Figure 2

ETL vs. ELT

ETL (Extract, Transform, and Load) and ELT (Extract, Load, then Transform) are data integration techniques. In ETL, data is converted and cleansed before being loaded into the target system, while in ELT, data is loaded first and then converted into the target system. ELT aligns with today's big data and cloud computing systems, apart from utilizing the storage capability to perform transformations on big data.

ELT vs. ETL: What is the difference?

Aspect	ETL (Extract, Transform, Load)	ELT (Extract, Load, Transform)
Transformation Location	Performed in a separate phase after extraction and before loading into the target system	Directly performed in the target data storage system after the data is loaded
Suitable for	Smaller datasets where transformations can be efficiently performed before loading	Large datasets and big data applications, especially in cloud-based data warehouses
Processing Efficiency	May require additional computational resources for transforming the data before loading	Leverages the processing power of modern data storage systems for efficient transformations
System Requirements	Less intensive system requirements as transformations are performed before loading	Requires robust data storage systems capable of handling transformations after the data is loaded
Flexibility	Offers more flexibility in terms of data manipulation and transformation options before loading	Provides flexibility in scaling and handling large datasets efficiently

Figure 3

Compliance and Security Challenges in ETL

Since dealing with healthcare and financial data, ETL processes must be safe, adhere to all the required regulations, and be fast. Some primary challenges include:

1. **Data Privacy and Confidentiality:** In healthcare, laws such as HIPAA in the United States set very specific guidelines on the management of patient data. In the same way, financial institutions must keep the customer's information private to meet the requirements of GDPR or the Gramm-Leach-Bliley Act. Experience also shows that implementing these regulations affects all phases of the ETL process and requires anonymization and even encryption [10].
2. **Regulatory Compliance:** Every industry has its compliance challenges. When data is exchanged between different institutions in the healthcare sector, it has to be in compliance with CDM standards and regional standards [11]. Financial services have to adhere to the new Basel III regulation that demands the provision of liquidity coverage ratios or LCR on a cyclic basis to reduce systemic risk.
3. **Technical Constraints and Resource Limitations:** Most healthcare and financial organizations need to be better positioned to invest in specific instruments like software or have the competence to address ETL management adequately. This is attended by a limitation that is capable of causing errors, delays, and eventual low use of resources, which in turn affects the integrity of the data and the endorsement of the regulations.
4. **Data Integrity and Quality:** The quality of each piece of data is crucial in both industries since the wrong information leads to wrong decisions and regulatory violations. The problem is made worse by the ability of data to be incompatible between different systems, the level of data model sophistication, and the technical nature of data quality checks, which require specialist knowledge.

Approaches for Ensuring that ETL is Safe and Lawful

Thus, to overcome issues connected with data safety and meet the requirements of regulations, healthcare and financial services companies must include the best security practices of ETL processes. These strategies do more than improve data security; they also meet compliance with regulatory standards, which is essential when dealing with sensitive data. Below are recommended approaches:

1. Implement Robust Data Governance Frameworks:

Implementing an overarching data governance framework is the primary recommendation for managing data in compliance with the organizational policies, regulatory guidelines, and standards. Such frameworks set up proper standard practices for managing data related to ETL works, including but not limited to access rights, data sorting ancestry records, and logging, all of which bring about data recurrent nature and credibility [12]. For example, governance tools can be used at the healthcare organization to check that patient data is sufficiently guarded at each ETL phase. At the same time, at the financial institution, it is possible to ensure that data management complies with laws such as GDPR so that the company will not get entangled in legal violations.

2. Utilize Advanced Encryption Techniques:

Securing files during Extraction, Transformation, and loading helps safeguard the files from unauthorized personnel access and adds to the safety issues in the ETL process. For instance, the encryption mechanism must meet HIPAA standards in providing healthcare services, ensuring the protection of patient information at every stage. Like any other company, financial institutions require cryptography techniques to address GDPR and PCI DSS regulations where personal and transactional information is processed [13]. This approach goes a long way toward improving data security and compliance with regulatory industry standards, and we reduce the organization's vulnerability to security breaches and fines that come with them.

3. Automate Data Quality Checks:

Pre-storage data audit prevents augmentations such as missing, twofold, and conflicting data from penetrating the storage arrangement. This is especially important in a healthcare institution and the financial industry, where reliable data determines decisions and financial reporting. For instance, in financial ETL procedures, data quality checks for automation can enhance liquidity ratio computations necessary for keeping with the guidelines of financial reporting standards. Similarly, such controls are also used in health care to verify the accuracy of a patient and research record, which leads to improved results about the patient and record and improved research results.

4. Adopt Containerization and Distributed Computing:

Tools like Docker and PySpark enhance the general effectiveness of an ETL process when working with large volumes of data, as well as security and scalability. Containerization using services such as Docker entails the construction of closed systems that safeguard the data processing configuration to external interferences and guarantee a standard execution platform. Computer programs such as PySpark enhance the use of resources and time so that institutions can effectively work on large volumes of information [14]. They also help meet processing timelines required of regulated industries and boost data security and compliance. For better understanding and illustration, see the figure below.

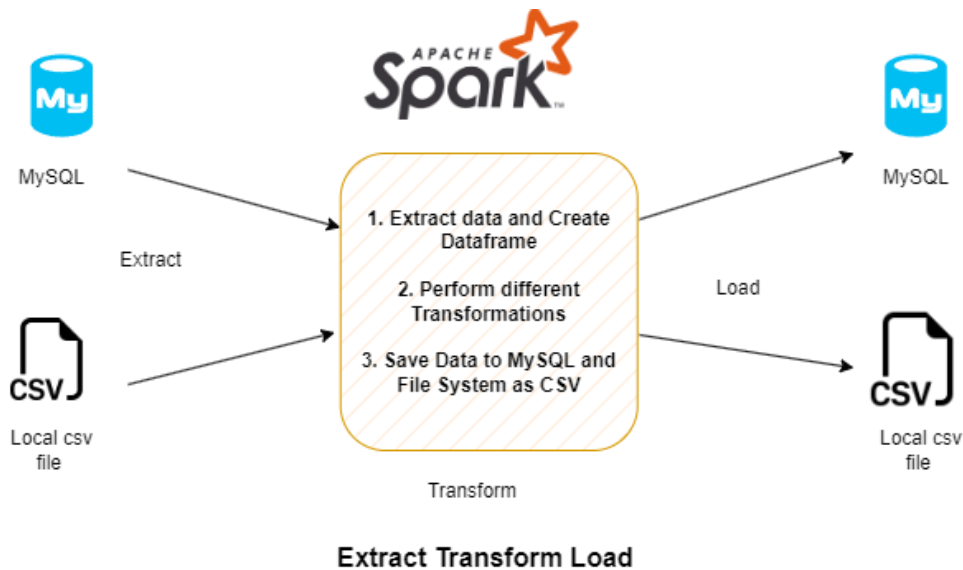


Figure 4

The following is an integration scenario that takes the form of an Apache Spark, an unyielding big data processing engine, in the middle of different MySQL databases. In this integration, PySTAPR (Python-Spark-Transformation-Aggregation-Pipeline) is adopted to improve the standard distributed application framework to handle TransfTransformationther processes in Spark. PySTAPR connects to MySQL to pull data from source tables, and the operation performs data operations in Spark clusters and writes the resulting data back to MySQL or another database. It resembles the typical ETL (Extract, Transform, Load) setup. This is because PySTAPR minimizes the complexity and time required for the transformation steps and optimizes data handling with Spark.

Case Studies and Industry Applications

Several case studies illustrate the successful application of these ETL strategies in healthcare and financial services:

Healthcare Data Transformation:

In the MIMIC-III project, we prove the applicability of containerized ETL integration to the processing of vast healthcare data. Employing PySpark to develop scalable ML models, this project adopted containerization to improve scalability while meeting high-security standards for data storage. Through the implementation of containerization, the project was able to enhance data management and security and, at the same time, create a yardstick for ETL in healthcare [15].

Bank Liquidity Coverage Ratio Reporting:

In financial services, adopting ETL frameworks to support Basel III reporting has been dramatically helpful in enhancing the LCR reporting quality and efficiency. Banks have reduced operational risks and complied with the reporting requirements in line with industry requirements

through automation. This case outlines how using an automated version of ETL aided by software tools assists financial institutions in transforming data and addressing the need for regulatory compliance.

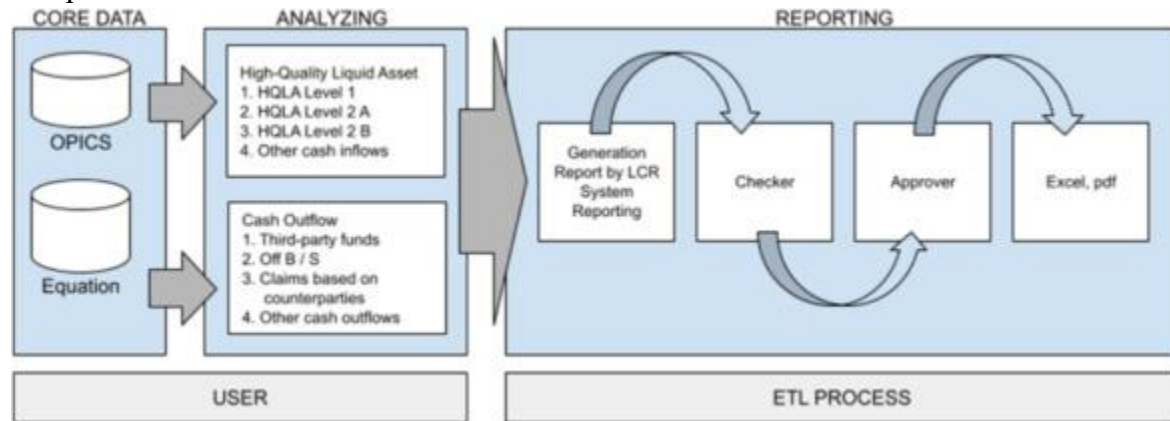


Figure 5

The following diagram illustrates the overall ETL process for Basel III for banks, mainly covering LCR requirements. Information from the OPICS and Equation, tree core banking data, disaggregates and classifies assets and cash flows, including HQLA and cash expenditures. The LCR document and computations continue to the following phase, which is the reporting phase of a firm's LCR. The workflow includes a checker who verifies the data, the data sent to the Approver, and the ETL's generation of the Excel and PDF reports through the framework, ensuring compliance and accurate LCR reporting [16]. This arrangement is helpful for banks because it helps manage compliance with the Basel III standards as it organizes data processing and report creation.

Conclusion

There are many reasons why the security and compliance of ETL processes are essential for healthcare and finance businesses. Data governance standards, high-level security measures, automated quality controls, and the like, in addition to emerging technologies such as containerization and distributed processing, help organizations improve their ETL process to handle data privacy or regulatory constraints efficiently. With the advancement of these techniques in use, healthcare and financial institutions, which are the significant customers using these techniques, will enhance data management, moving to greater heights towards improving data security and ensuring compliance with new set standards by regulatory bodies. Apart from this, it is also proactive in protecting organizations and creating confidence with the stakeholders who provide their sensitive information to the organizations.

Bibliography

[1]

A. Krylov, “Data Security in Healthcare: Tips for Cybersecurity,” Mar. 14, 2023. <https://kodjin.com/blog/why-healthcare-data-security-solutions-are-important/>

[2]

K. Hoffmann *et al.*, “Data integration between clinical research and patient care: A framework for context-depending data sharing and in silico predictions,” *PLOS Digital Health*, vol. 2, no. 5, p. e0000140, May 2023, doi: <https://doi.org/10.1371/journal.pdig.0000140>.

[3]

P. Shojaei, E. V. Gjorgievska, and Y.-W. Chow, “Security and Privacy of Technologies in Health Information Systems: A Systematic Literature Review,” *Computers*, vol. 13, no. 2, p. 41, Feb. 2024, doi: <https://doi.org/10.3390/computers13020041>.

[4]

T. Ong, R. Pradhananga, E. Holve, Iii, and M. Kahn, “A Framework for Classification of Electronic Health Data Extraction-Transformation-Loading Challenges in Data Network Participation,” 2019. Accessed: Nov. 07, 2024. [Online]. Available: <https://pmc.ncbi.nlm.nih.gov/articles/PMC5994935/pdf/egems-5-1-222.pdf>

[5]

Ehsan Soltanmohammadi and Neset Hikmet, “Optimizing Healthcare Big Data Processing with Containerized PySpark and Parallel Computing: A Study on ETL Pipeline Efficiency,” *Journal of Data Analysis and Information Processing*, vol. 12, no. 04, pp. 544–565, Jan. 2024, doi: <https://doi.org/10.4236/jdaip.2024.124029>.

[6]

A. Itsekson, “The Importance of ETL in Healthcare: All You Need To Know,” *Jelvix*, 2023. <https://jelvix.com/blog/etl-process-in-healthcare-benefits-challenges-and-best-practices>

[7]

S. Khanra, A. Dhir, A. K. M. N. Islam, and M. Mäntymäki, “Big data analytics in healthcare: a systematic literature review,” *Enterprise Information Systems*, vol. 14, no. 7, pp. 878–912, Aug. 2020, doi: <https://doi.org/10.1080/17517575.2020.1812005>.

[8]

V. Ehrenstein, H. Kharrazi, H. Lehmann, and C. O. Taylor, *Obtaining Data From Electronic Health Records*. Agency for Healthcare Research and Quality (US), 2020. Available: <https://www.ncbi.nlm.nih.gov/books/NBK551878/>

[9]

W. Raghupathi and V. Raghupathi, “Big data analytics in healthcare: promise and potential,” *Health Information Science and Systems*, vol. 2, no. 1, pp. 1–10, Feb. 2019, doi: <https://doi.org/10.1186/2047-2501-2-3>.

[10]

N. Berros, F. El Mendili, Y. Filaly, and Y. El Bouzekri El Idrissi, “Enhancing Digital Health Services with Big Data Analytics,” *Big Data and Cognitive Computing*, vol. 7, no. 2, p. 64, Mar. 2023, doi: <https://doi.org/10.3390/bdcc7020064>.

[11]

C. Peng, P. Goswami, and G. Bai, “A literature review of current technologies on health data integration for patient-centered health management,” *Health Informatics Journal*, vol. 26, no. 3, p. 146045821989238, Dec. 2019, doi: <https://doi.org/10.1177/1460458219892387>.

[12]

B. Ozaydin, F. Zengul, N. Oner, and S. S. Feldman, “Healthcare Research and Analytics Data Infrastructure Solution: A Data Warehouse for Health Services Research,” *Journal of Medical Internet Research*, vol. 22, no. 6, p. e18579, Jun. 2020, doi: <https://doi.org/10.2196/18579>.

[13]

V. Manickam and M. Rajasekaran Indra, “Dynamic multi-variant relational scheme-based intelligent ETL framework for healthcare management,” *Soft Computing*, Mar. 2022, doi: <https://doi.org/10.1007/s00500-022-06938-8>.

[14]

R. Raja, I. Mukherjee, and B. K. Sarkar, “A Systematic Review of Healthcare Big Data,” *Scientific Programming*, vol. 2020, no. 1, pp. 1–15, Jul. 2020, doi: <https://doi.org/10.1155/2020/5471849>.

[15]

A. Almalawi, A. I. Khan, F. Alsolami, Y. B. Abushark, and A. S. Alfakeeh, “Managing Security of Healthcare Data for a Modern Healthcare System,” *Sensors*, vol. 23, no. 7, p. 3612, Jan. 2023, doi: <https://doi.org/10.3390/s23073612>.

[16]

F. Prasser, H. Spengler, R. Bild, J. Eicher, and K. A. Kuhn, “Privacy-enhancing ETL-processes for biomedical data,” *International Journal of Medical Informatics*, vol. 126, pp. 72–81, Jun. 2019, doi: <https://doi.org/10.1016/j.ijmedinf.2019.03.006>.



©2024 by the Authors. This Article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>)