# A Technical Review of Dynamic and Mixed Approach for Health Data Extraction, Transformation and Loading Process

## Adya Mishra

Independent Researcher, Virginia, USA
adyamishra29@gmail.com

**Abstract**

Healthcare data originates from a diverse array of sources—including electronic health records (EHRs), laboratory systems, wearable devices, and unstructured clinical text—making its integration a complex endeavor. Traditional extraction, transformation, and loading (ETL) pipelines, though foundational, often struggle to keep pace with evolving data schemas, regulatory obligations, and the need for real-time insights. This paper provides a technical review of dynamic and mixed ETL strategies tailored specifically for health data. Dynamic approaches emphasize adaptive schema discovery, rule-based transformations, and metadata-driven designs that automatically adjust to new and updated data sources, reducing manual reconfiguration. Mixed ETL models integrate both real-time streaming and batch processing, enabling healthcare organizations to process time-critical clinical data immediately while performing more complex transformations on larger datasets at scheduled intervals. Key challenges—including data quality assurance, regulatory compliance, and the requirement for robust security—are addressed alongside recommended best practices for metadata management, rule engines, and orchestration tools. The review further highlights implementation considerations and future trends, such as AI-driven data integration, serverless architectures, and data mesh paradigms. By adopting these flexible, scalable ETL approaches, healthcare institutions can enhance the accuracy, timeliness, and security of patient data analytics—ultimately improving clinical decision-making and patient outcomes.

**Keywords:** Electronic health records, Extraction, Transformation and loading

## 1. INTRODUCTION

Health data stands at the forefront of modern analytics, with healthcare providers, researchers, and policy makers increasingly relying on data-driven insights to improve patient outcomes, reduce operational costs, and streamline care delivery. Such data may be generated from electronic health records (EHRs), medical devices, laboratory results, imaging systems, and even personal health-tracking apps. Each data source has distinct characteristics—ranging from structured patient demographics to unstructured clinical notes—making data integration a fundamentally complex and challenging task.

Extraction, Transformation, and Loading (ETL) processes in healthcare must not only consolidate disparate data formats, but also ensure that the data meets stringent quality, security, and compliance standards (e.g., HIPAA, GDPR). Traditional, static ETL pipelines often assume that input data structures are stable and that transformations remain constant over time. However, healthcare environments are inherently dynamic:

data standards evolve, regulatory requirements shift, and new data sources emerge regularly. As a result, ETL approaches
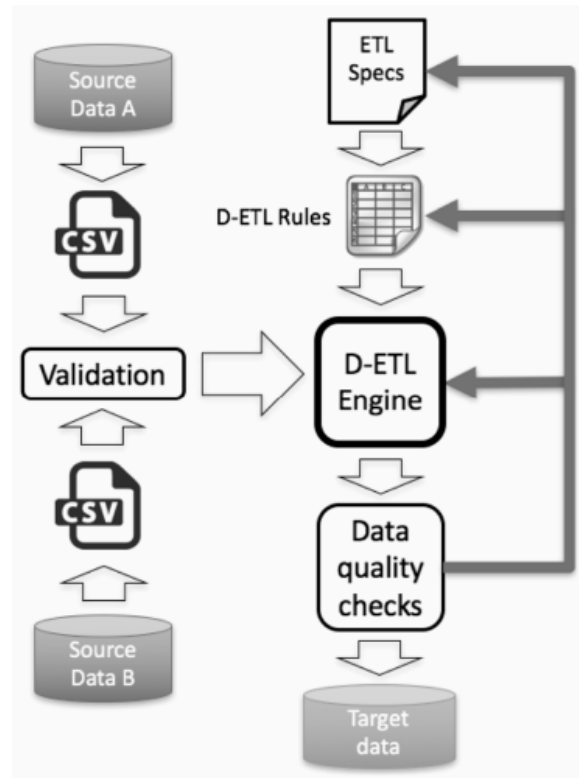


**Fig. 1. ETL Workflow approach [2].**

that are static or monolithic are increasingly ill-suited to meet the pace and variability of today's health data ecosystem.

This paper provides a comprehensive technical review of dynamic and mixed approaches for health data ETL, discussing how these paradigms offer greater flexibility, scalability, and robustness. We begin by examining the complexities inherent to health data, followed by an overview of traditional ETL fundamentals. We then delve into the principles behind dynamic ETL frameworks, highlighting their adaptability to evolving data structures. Afterward, we address mixed approaches, which combine both real-time and batch processing for a balanced data integration strategy. Key challenges, implementation considerations, and future outlooks are presented to give readers a 360-degree view of where health data ETL is heading in the coming years.

## 2. ETL SYSTEM OVERVIEW

ETL System is a broad presentation of data movement and transactional processes from extraction of multiple application sources, transforming data into dimensional fact table and conformed format to feed data into data warehouse environment such as data marts, data stores or other target systems. The process is widely applied in data integration, migration, staging and master data management efforts [1].

Healthcare ETL inherits this paradigm but faces additional challenges such as ensuring compliance, maintaining auditable logs, and merging clinical metadata across varied systems. Moreover, the velocity and volume of data in modern healthcare settings can be significant—particularly when real-time monitoring and telemedicine are involved—necessitating new design considerations [3].

## 3. THE COMPLEXITY OF HEALTH DATA

### A. Heterogeneity of Data Sources

One of the most significant challenges in healthcare data integration is its heterogeneity. Data sources include, but are not limited to:

1. Electronic Health Records (EHRs): Structured and semi-structured data capturing patient demographics, diagnoses, medications, and billing codes.

2. Diagnostic and Imaging Systems: Large file formats like DICOM for radiology and cardiology images, requiring specialized handling.

3. Laboratory Information Management Systems (LIMS): Often produce standardized results (e.g., HL7 messages) but can vary in format depending on the vendor.

4. Real-Time Monitoring Devices (IoT/Medical Devices): Devices that track patient vitals, generating continuous streams of time-series data.

5. Unstructured Data: Clinical notes, physician narratives, and transcribed reports that need natural language processing (NLP) for integration.

Effectively combining these diverse formats demands an ETL process that can adapt to evolving schemas and a wide array of communication protocols (HL7, FHIR, DICOM, etc.).

### B. Regulatory and Privacy Constraints

Healthcare data is among the most sensitive types of information, subject to strict rules around privacy and consent. In the United States, the Health Insurance Portability and Accountability Act (HIPAA) outlines rigorous requirements for data handling, encryption, and access control. In Europe, the General Data Protection Regulation (GDPR) imposes additional constraints for data storage and transfer. These constraints significantly influence how the ETL pipeline is designed, deployed, and maintained, adding layers of security, auditing, and authorization [5].

### C. Data Quality and Clinical Relevance

A critical aspect of health data is the emphasis on quality and clinical relevance. Inaccurate or incomplete data can lead to flawed diagnoses, dangerous prescriptions, or incorrect research outcomes. Data transformations must ensure proper labeling, validation, and sometimes normalization (e.g., converting diagnostic codes from ICD-9 to ICD-10 or SNOMED CT). Traditional ETL might rely on fixed data
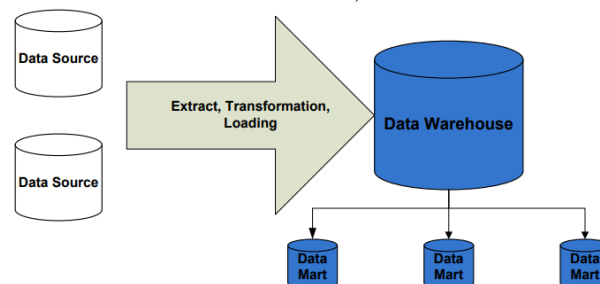


**Fig. 2. ETL Architecture [4].**

quality checks, but the evolving nature of healthcare data standards and continuous updates to code sets demand more flexible approaches.

### D. Dynamic Approaches to Health Data ETL

A "dynamic" ETL approach is one that automatically adapts to changes in source schemas, data volume, and quality requirements with minimal manual intervention. Rather than relying on static, hard-coded mappings, dynamic ETL pipelines monitor source systems for schema evolution and can update mappings,

transformations, or metadata configurations on-the-fly. This flexibility helps address the ever-changing nature of healthcare data standards and ensures a stable integration pipeline even as new data sources come online [4].

## E. Key Principles of a Dynamic ETL Framework

1. Schema Discovery and Evolution: Dynamic ETL pipelines continuously discover or detect changes in source schemas, adjusting data mappings automatically. This might involve using machine learning techniques to infer schema from raw data, or scanning metadata repositories for changes.

2. Rule-Based Transformation Logic: Instead of explicitly coding transformations, dynamic pipelines rely on rule engines and parameterized transformations. For instance, a rule might state, "All medication names in Source A must map to RxNorm codes," which can be updated as new medication codes emerge.

3. Metadata-Driven Design: Dynamic frameworks place heavy emphasis on metadata, storing information about schemas, transformations, quality rules, and lineage in a centralized repository. Automated processes then use these metadata entries to execute transformations without manual reconfiguration.

4. Event-Driven Orchestration: With dynamic ETL, updates to data sources or new data arrival can trigger reprocessing, partial reloading, or schema evolution procedures. This event-driven style ensures that the pipeline remains responsive to real-time changes.

## F. Mixed Approaches: Combining Real-Time and Batch

While dynamic ETL focuses on adaptability, a mixed approach addresses the varied latency requirements of health data. Some information (e.g., critical vitals from an Intensive Care Unit) must be processed in near real-time for immediate clinical decision-making. Other data, such as administrative records or historical lab results, can be processed in batch to optimize resource usage.

Mixed approaches typically involve a hybrid architecture:

1. Streaming or Real-Time Layer: Uses technologies like Apache Kafka, Apache Flink, or AWS Kinesis to handle immediate data feeds. Data is quickly validated, transformed (at least partially), and made available for time-sensitive analytics.

2. Batch Layer: Periodically processes large volumes of historical or less time-critical data, often in a data warehouse or data lake environment. The batch transformations can be more complex, employing machine learning, advanced NLP, or computationally expensive matching algorithms.

## G. Technical Components of a Mixed ETL Pipeline

1. Data Ingestion Layer**:** A combination of message queues (e.g., Kafka topics) and file-based ingestion (e.g., HL7 messages via SFTP).

2. Real-Time Processing Engine: A streaming platform that applies lightweight transformations—data validation, patient ID resolution, immediate alert generation—and routes data to monitoring dashboards or operational data stores.

3. Batch Processing Engine: A big data framework (Spark, Hadoop) or a traditional ETL tool (Informatica, Talend) that handles deep transformations, complex code conversions (ICD-9 to ICD-10), and aggregates historical data for longitudinal analytics.

4. Data Storage and Serving Layer: Depending on the design, this may include a data lake (e.g., AWS S3, Azure Data Lake), a structured data warehouse (Snowflake, SQL Server), and specialized analytics platforms for large-scale queries [6].

## 4. TOOLS AND TECHNOLOGIES SUPPORTING DYNAMIC AND MIXED ETL

### A. Metadata Management Systems

Metadata repositories such as Apache Atlas or commercial solutions like Informatica Metadata Manager can maintain comprehensive dictionaries of data schemas, lineage, and transformation rules. In dynamic ETL frameworks, these repositories act as the "brain," guiding the pipeline whenever a new field or updated code appears in a source system.

### B. Workflow Orchestration

Tools like Apache Airflow, Luigi, or Azure Data Factory are frequently used for pipeline orchestration. These platforms enable event-driven triggers, parameterization, and parallel execution—essential for dynamic or hybrid scenarios. In a healthcare setting, a newly deployed EHR interface can trigger an automated job to update transformations for the entire pipeline.

### C. Rule Engines and Machine Learning

Rule engines (Drools, Camunda) and ML libraries (TensorFlow, PyTorch) can support dynamic schema detection or anomaly-based transformations. For instance, an ML model might automatically identify new medication codes from textual data and map them to standard RxNorm codes if a matching pattern is detected [8].

### D. Real-Time Stream Processing

Apache Kafka, Apache Flink, AWS Kinesis, and Google Cloud Pub/Sub are widely used for streaming ingestion in mixed architectures. They allow real-time data ingestion from IoT medical devices, EHR event logs, or HL7 messages. Flink, Spark Streaming, or Storm can further process the streams in memory, applying minimal transformation before distributing data to various sinks (databases, dashboards, or alerting systems).

### E. Traditional ETL Platforms

Despite the rise of big data tools, platforms like Informatica PowerCenter, Talend, SSIS, and Pentaho remain staples in healthcare. These tools offer user-friendly interfaces and extensive libraries for HL7 processing, code set mappings, and connectivity to major EHR vendors. Within a dynamic or mixed framework, they often serve as the batch processing engine [7].
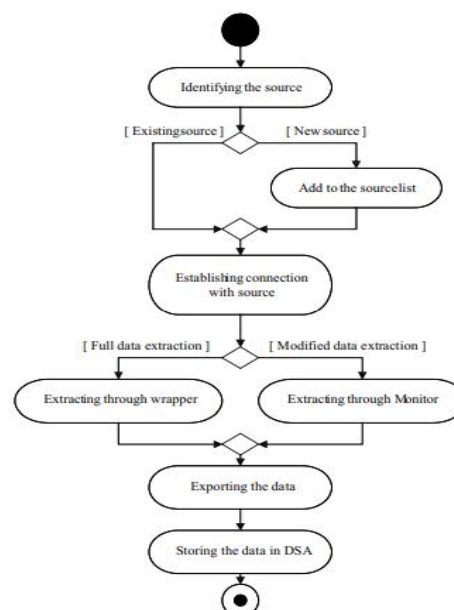


**Fig. 3. Data Extraction Diagram [5].**

## 5. CONSIDERATIONS

Performance considerations for health data ETL revolve around several key factors, beginning with the latency versus accuracy trade-off. In real-time processing scenarios, concerns about computational overhead and complexity often arise, as streaming frameworks require simplified transformation logic to maintain low latency. Ensuring that these minimal transformations still produce clinically relevant insights may involve iterative design and close collaboration with healthcare professionals. Conversely, batch processes can accommodate more complex data manipulations but inherently introduce time delays, making them less suitable for acute care settings where immediacy is crucial. Another factor involves scaling in the cloud: many healthcare organizations leverage platforms such as AWS, Azure, or GCP, which provide autoscaling features, serverless computing (e.g., AWS Lambda, Azure Functions), and managed big data services (e.g., EMR, Databricks) to respond dynamically to fluctuating data volumes. However, this approach necessitates strict adherence to HIPAA/HITECH compliance, requiring HIPAA-eligible services and robust network security configurations—such as virtual private clouds (VPCs) and private endpoints—to safeguard protected health information (PHI). Finally, caching and data storage optimization play a significant role in performance. Caching strategies, such as storing reference tables for repetitive lookups, reduce overhead, while the use of columnar storage formats (e.g., Parquet or ORC) in batch environments can dramatically decrease read times for large-scale queries common in population health analyses.

## 6. CONCLUSION

Healthcare data represents one of the most intricate domains for ETL, given its diverse sources, regulatory constraints, and clinical significance. Traditional ETL pipelines—while foundational—often struggle to keep pace with rapid changes in data structure, vocabulary, and volume. Dynamic ETL approaches address this challenge by automatically adapting to evolving schemas and leveraging metadata-driven, rule-based transformations. This adaptability ensures better data quality, regulatory compliance, and reduced maintenance overhead.

Meanwhile, mixed ETL architectures reconcile the need for real-time patient care and long-term analytics. By coupling streaming frameworks with batch processing, organizations can deliver immediate insights to clinicians while still performing deep, computationally intensive transformations on historical data. Although this hybrid model requires careful orchestration and data governance, it provides a flexible path forward for healthcare institutions seeking to maximize the value of their data.

Implementing dynamic or mixed ETL strategies is not without challenges. Organizations must invest in metadata management, rule engines, streaming platforms, and robust security measures. Effective data governance is essential to maintain consistency and ensure that data remains clinically valid. Moreover, staff training and change management become critical to ensure adoption of new tools and processes.

Looking ahead, advancements in serverless computing, AI-driven integration, and emerging paradigms like data mesh are likely to further transform the healthcare ETL landscape. The continuous evolution of security technologies will also play a decisive role in allowing healthcare organizations to safely harness patient data for research, population health, and precision medicine. By embracing dynamic and mixed approaches, healthcare providers and researchers can keep pace with the rapidly changing demands of modern medicine, driving improved patient outcomes and system-wide efficiencies.

In summary, dynamic and mixed ETL approaches stand as pivotal strategies for healthcare organizations seeking agility, scalability, and robust data integration. As medical data ecosystems continue to expand and

incorporate new data types, these adaptive ETL paradigms will be key enablers of timely, accurate, and secure healthcare analytics—ultimately contributing to better patient care and innovation in the field.

### REFERENCES

1. Sabtu, A., Azmi, N. F. M., Sjarif, N. N. A., Ismail, S. A., Yusop, O. M., Sarkan, H., & Chuprat, S. (2017, July). The challenges of Extract, Transform and Loading (ETL) system implementation for near real-time environment. In 2017 International Conference on Research and Innovation in Information Systems (ICRIIS) (pp. 1-5). IEEE.

2. Ong, T. C., Kahn, M. G., Kwan, B. M., Yamashita, T., Brandt, E., Hosokawa, P., ... & Schilling, L. M. (2017). Dynamic-ETL: a hybrid approach for health data extraction, transformation and loading. BMC medical informatics and decision making, 17, 1-12.

3. Ong T, Pradhananga R, Holve E, Kahn MG. A Framework for Classification of Electronic Health Data Extraction-Transformation-Loading Challenges in Data Network Participation. EGEMS (Wash DC). 2017 Jun 13;5(1):10. doi: 10.5334/egems.222. PMID: 29930958; PMCID: PMC5994935.

4. Simitsis, A., Vassiliadis, P., & Sellis, T. (2005). Extraction-transformation-loading processes. In Encyclopedia of Database Technologies and Applications (pp. 240-245). IGI Global.

5. Runtuwene, J. P. A., Tangkawarow, I. R., Manoppo, C. T. M., & Salaki, R. J. (2018, February). A comparative analysis of extract, transformation and loading (ETL) process. In IOP Conference Series: Materials Science and Engineering (Vol. 306, No. 1, p. 012066). IOP Publishing.

6. Mrunalini, M., Kumar, T. S., & Kanth, K. R. (2009, November). Simulating secure data extraction in extraction transformation loading (ETL) processes. In 2009 Third UKSim European Symposium on Computer Modeling and Simulation (pp. 142-147). IEEE.

7. Kumar, S., & Nadeem, M. (2008). Extraction, Transformation, Loading (ETL) and Data Cleaning Problems. Journal of Independent Studies and Research on Computing, 6(1).

8. Dhanda, P., & Sharma, N. (2016). Extract Transform Load Data with ETL Tools. International journal of advanced research in computer science, 7(3)