

CLOUD-NATIVE DATA CONVERSION FOR MEDICARE & MEDICAID: A SCALABLE FOUNDATION FOR ANALYTICS AND AI

PONNARASAN KRISHNAN

Senior Developer, Acentra Health, USA.

Abstract

The growth and complexity of healthcare datasets, particularly within Medicare and Medicaid systems, exacerbate issues of interoperability, analytics, and integration with AI. To some extent, legacy data pipelines and extract-transform-load (ETL) frameworks have remained bounded by limitations on the expansiveness to address heterogeneous formats and regulatory compliance on standards such as HIPAA, HL7, and FHIR. To overcome these limitations, this research proposes a cloud-native data conversion framework using microservices, containerization, and serverless computing to build an AI-ready and scalable, secure foundation for healthcare analytics. The proposed framework allows for the raw ingestion of Medicare and Medicaid datasets with schema conversion automated into standard healthcare formats and optimized storage considerations downstream for analytics and AI. By benchmarking for performance, validating compliance, and testing for scalability, the framework gets to demonstrate superiority over traditional ETL pipelines regarding the speed of data conversion, resource elasticity, and integration with machine learning. Case studies exhibit its use in predictive care analytics, fraud detection, and healthcare policy optimization, establishing a determined path for real-time, data-driven decision-making within healthcare ecosystems. The work aligns with industry standards and leverages cloud-native advantages to contribute toward a scalable solution for transforming Medicare and Medicaid data into a foundation that fast-tracks advanced analytics and innovates AI-driven healthcare.

Keywords: Cloud-Native Computing; Data Conversion; Medicare & Medicaid; Healthcare Analytics; FHIR; HL7; Artificial Intelligence; Interoperability; Scalable Data Pipelines; Health Data Management.

1. INTRODUCTION

A treasure trove of data in all shapes and sizes is generated in the U.S. health system, more so under federal health programs such as Medicare and Medicaid. Together, these two programs include more than 150 million beneficiaries; with that number, administrative, claims, EHR, prescription, and billing data in the order of petabytes are generated every year (Michael, 2025). Such an enormous quantity of data presents huge opportunities for advanced analytics and AI-driven insight but also creates an equally huge problem of data integration, interoperability, and compliance. Typically, healthcare organizations would rely on legacy ETL pipelines to ingest, clean, and prepare the data for reporting and analysis. However, while ETL is an excellent approach in a traditional enterprise setting, it finds difficulty in the healthcare ecosystem, where data is scattered across heterogeneous sources with proprietary formats and standards varying greatly from HL7, FHIR, ICD-10, or CPT standards (Saini et al., 2021). These pipelines are typically rigid, costly to operate, and out of sync with the flexibility required for AI workloads of today (Sharma, 2025).

1.1 Background on Medicare & Medicaid Data

Medicare and Medicaid datasets vary sharply from regular enterprise datasets. containing

layers of claims processing applications, provider submissions, EHRs, diagnostic codes, and reimbursement schedules, these datasets do not constitute transactional corporate data (Gaddam, 2025). Unlike corporate datasets, healthcare data is regulated by stiff regulatory systems such as the Health Insurance Portability and Accountability Act (HIPAA), which maintain privacy, audit, and role-based data access (Conteh, 2024). Typical Medicare claims hold patient demographic data, provider information, ICD 10 diagnostic codes, treatment coding and payment adjustment data, and longitudinal records spread out among multiple providers. Medicaid adds a layer of complexity by including state-level idiosyncrasies in terms of coverage, eligibility, and claims adjudication (Tilahun, 2023). Hence, Medicare and Medicaid data ecosystems are volatility siloed (payer, provider, pharmacy) and horizontally fragmented across states, agencies, and private contractors (Lee et al., 2022). Such fragmentation makes it hard to integrate the data for population health analyses, fraud analytics, and predictive modeling, which ultimately require end-to-end visibility of patient and provider journeys.

1.2 Limitations of Legacy ETL Methods

Traditional ETL systems were designed when healthcare datasets were relatively small, less heterogeneous, and more for statistical reporting rather than interactive real-time AI models. These pipelines extract data usually from source systems, perform schema transformations, and then load the transformed data into relational databases or data warehouses (Wang & Zhao, 2020). While fine for batch-oriented reporting, the model has major drawbacks when adapted toward Medicare and Medicaid analytics:

- **Scalability Bottlenecks:** Because of the fixed infrastructure, ETL jobs often witness serious bottlenecks when the landscape demands processing of millions of claims or terabytes of clinical data (Bauer et al., 2025).
- **Schema Rigidity:** Every new data source (say a new state Medicaid claim system) requires manual schema mapping, which makes ETL rigid and expensive to maintain (Sharma et al., 2025).
- **Latency:** Being batch-driven, ETL pipelines are the farthest thing from real-time analytics, which are already becoming pivotal in fraud detection, care management, and public health surveillance (Mohamed, 2025).
- **Compliance Risks:** HIPAA, GDPR, and several other data protection laws require governance at a very detailed level, which most traditional ETL tools fail to deliver (Aziz & Hussain, 2025).
- **Not Prone to AI-Levels:** The output of ETL is meant for relational queries and less for AI/ML pipelines that rely heavily on the scalable access of a rich feature set with global datasets of multiple modalities (Jay, 2023).

This combination of drawbacks not only stifles the process of innovation but also hinders the cost factor and increases risk factors for compliance, leaving Medicare and Medicaid administrators behind modern infrastructure for data-driven transformation in healthcare.

1.3 Motivation for Cloud-Native Data Conversion

To navigate around these obstacles, healthcare realizes the potential of cloud-native architecture, and they are adopting cloud-native architectures based on containerization, microservices, serverless, and elastic scaling (Pasupuleti et al., 2025). In the building of cloud-native data conversion frameworks there lies a paradigm shift from classic ETL into distributed, scalable, and intelligent pipelines capable of mutating healthcare datasets in real time.

Significant highlights for cloud-native include:

- Elastic Scalability: Data ingestion and transformation workloads scale automatically when surges in claims or clinical data submissions occur (Maxwell, 2024).
- Interoperability by Design: Conversion engines out of the box adhere to healthcare standards (FHIR, HL7, CDA), facilitating data exchanges between payers, providers, and researchers (Carrillo et al., 2022).
- Compliance-Embedded-Workflows: With integrated DevSecOps and policy enforcement, cloud-native pipelines apply HIPAA compliance to each layer (Chandramouli, 2022).
- AI/ML Integrations: Cloud-native systems are natively connected to AI toolchains, enabling technologies such as predictive modeling, anomaly detection, and population health forecasting to take precedence (Subramaniam et al., 2025).
- Cost Efficiency: The pay-as-you-go model reduces infrastructure overheads and maximizes resource utilization compared to the static ETL clusters (Freeman & Harvey, 2020).

This conversion is more than a well-honed technical solution; it is a disruptive change to Medicare and Medicaid data, providing the much-needed thrust to next-gen healthcare analytics.

1.4 Contributions of this Research

This research proposes a cloud-native data conversion framework tailored for Medicare and Medicaid to allow large-scale analytic and AI works. Its contributions are:

- A conceptual architecture for scalable data conversion across heterogeneous Medicare and Medicaid datasets while still maintaining HIPAA compliance.
- A performance study comparing the cloud-native pipeline with legacy ETL with throughput, latent, and compliance benchmarks.
- Showing the AI readiness by integrating converted datasets into predictive care analytics and fraud detection use cases.
- Exploring regulatory, technical, and organizational implications for adopting cloud-native in U.S. healthcare ecosystems.

By tackling the barrier at technical, regulatory, and analytical levels, this work provides a blueprint for cloud-native modernization of healthcare data that puts Medicare and Medicaid working foundation as scalable bases for AI-driven innovation.

Streamlining Medicare & Medicaid Data Conversion

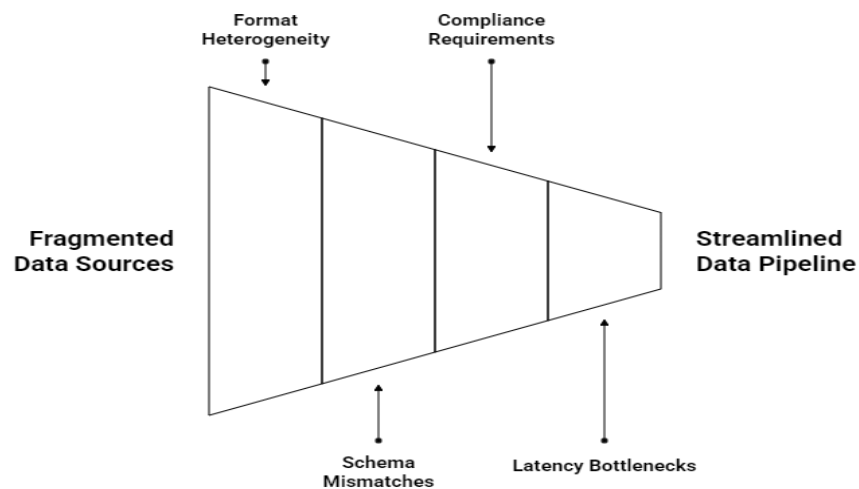


Figure 1: Conceptual overview of challenges in Medicare & Medicaid data conversion

2. LITERATURE REVIEW

The modernization of healthcare data infrastructures, especially in Medicare and Medicaid systems, is an emerging area with growing theoretical and practical interest. Here, the section approaches the complexity of Medicare-Medicaid datasets, explores data conversion challenges in healthcare, analyzes traditional ETL versus cloud-native solutions, discusses the extant research in AI readiness, and finally identifies from within this AM the research gap motivating this study.

2.1 Complexity of Medicare & Medicaid Datasets

Medicare and Medicaid data ecosystems are extremely huge and heterogeneous, presenting several barriers to seamless integration and analysis of the data. Being a federal insurance program for individuals over 65 and for certain younger individuals who are disabled, the Program generates millions of records of claims, enrollment records, prescription data, diagnostic codes, and reimbursement schedules. The Medicaid program, jointly funded by the Federal and State governments, throws in additional complexity with state-level variations in eligibility, benefits, and adjudication systems (Micheal, 2025; Tilahun, 2023).

Unlike transactional business data, these datasets exhibit:

- High dimensionality: with thousands of fields spanning demographics, provider attributes to diagnosis codes (ICD-10), treatment codes (CPT), and outcomes.
- Structural diversity: with data originating from claims systems, EHRs, provider registries, pharmacy records, and even a third-party contractor (Carrillo et al., 2022).
- Regulatory overlays: HIPAA, HITECH, and CMS reporting standards place very strict limits on privacy, auditability, and access (Conteh, 2024).
- Interoperability barrier: Despite the advent of standards like FHIR, yet the adoption remains inconsistent across states and providers (Saini et al., 2021).

As emphasized by Gaddam (2025), these Medicare and Medicaid payment systems also interface with financial gateways, fraud detection modules, and state-specific billing infrastructures, piling the heterogeneity on the givens in multiple layers. As such, the multidimensional heterogeneity significantly complicates not just conversion pipelines but downstream analytics and AI workflows.

2.2 Data Conversion Challenges in Healthcare

Healthcare organizations encounter technical, semantic, and compliance challenges in converting raw Medicare and Medicaid datasets into formats ready for analysis.

- Technical Complexity: Most legacy systems store data with proprietary formats to open a window for their own tools and procedures, posing custom adoptions and schema-mapping problems (Shah et al., 2024). Transformation workflows must cope with real-time ingestion from APIs, batch processes from flat files, and migration of historical data from mainframes (Maxwell, 2024).
- Semantic Inconsistency: Dictionary changes across states and providers cause inconsistent use of diagnostic and treatment codes (Naveen et al., 2024). Even within the same data set, subtle changes in billing codes can cause significant differences in interpretation.
- Compliance Burden: HIPAA, CMS, and GDPR introduced encryption, access logging, anonymization, and secure audit trails to force observance (Aziz & Hussain, 2025). These safeguards need to be embedded within the conversion framework at each and every stage.
- Latency and Timeliness: Public health surveillance, fraud detection, and care management now call for real-time or near-real-time analyses. Unfortunately, these needs are almost impossible to satisfy with legacy conversion frameworks (Sharma, 2025).

Conteh (2024) asserts that without rigorous data governance and standardization procedures, data conversion in healthcare becomes a recipe for leakage, duplication, and misinterpretation. Basilakis (2020) adds that this sensitivity in health data requires state-

of-the-art privacy-preserving methodologies such as homomorphic encryption and secure API integration.

2.3 Traditional ETL Frameworks versus Cloud-Native Approaches

The documented limitations of traditional Extract-Transform-Load (ETL) frameworks in the healthcare sector (Wang & Zhao, 2020) are numerous. ETL was designed without consideration for notes that were unstructured, claims of high velocity, and a multimodal way of representing healthcare data.

ETL challenges include:

- Rigid schema mappings requiring heavy manual intervention.
- Batch-driven processes unsuitable for continuous ingestion.
- Limited scalability when applied to the petabyte-scale Medicare/Medicaid systems.
- Poor integration into AI pipelines that require feature stores and model-ready datasets (Jay, 2023).

Big data platforms based on Hadoop and Spark arose as interim solutions to enable the distributed processing of big datasets. These frameworks increased scalability but still demanded a full set of infrastructure management and often were disjoint with compliance monitoring activities (Bauer et al., 2025).

In the current times, cloud-native approaches offer a paradigm shift. Through an amalgamation of microservices, serverless computing, containers, and elastic scaling, cloud-native platforms unite scalability, compliance, and AI-readiness into one pipeline (Pasupuleti et al., 2025). These architectures reiterate the principles of DevSecOps, allowing the pipeline to itself be part of the security and compliance (Chandramouli, 2022). Being API-first, these designs also allow for easy integration with FHIR, HL7, and external analytic platforms (Saini et al., 2021).

Table 1: Comparative Summary of Healthcare Data Conversion Frameworks

Feature/Criteria	Legacy ETL	Hadoop/Spark Big Data Platforms	Cloud-Native Frameworks
Scalability	Limited; fixed infra	High, but requires cluster mgmt	Elastic, auto-scaling
Data Types Supported	Primarily structured	Semi-structured + structured	Structured, semi-structured, unstructured
Compliance	Manual controls	Add-ons, partial	Embedded (DevSecOps, HIPAA-ready)
Latency	High (batch-driven)	Moderate	Low (real-time capable)
AI/ML Integration	Minimal	Limited (requires ETL stage)	Native integration (AI-ready datasets)
Cost Model	High upfront infra	Moderate, cluster-dependent	Pay-as-you-go, optimized
Governance	Weak	Partial	Strong, policy-driven

2.4 Existing Work on AI Readiness

The recent literature has rightly given prominence to the AI-readiness of healthcare data pipelines. The AI-ML methodologies require an abundant quantity of clean, interoperable, and timely datasets to train their algorithms for predictive care and fraud detection, risk stratification, and population health forecasting.

- **FHIR & AI Readiness:** Saini et al. (2021) reviewed the extent to which the FHIR standard is adopted, noting that while FHIR indeed promotes interoperability, the lack of consistent implementation amongst providers works against allowing AI and/or large systems to sit atop it.
- **Healthcare Data Standardization:** Carrillo et al. (2022) traced how the COVID-19 interventions promoted standardization but presented some of the issues faced in harmonizing legacy systems.
- **AI Integration in Cloud-Native Platforms:** Pasupuleti et al. (2025) alongside Sharma (2025) suggest that cloud-native design supports the seamless integration of AI toolchains, thus obviating further ETL.
- **Enterprise AI:** Jay (2023) argued that deployments of AI in the cloud environment would succeed only when the underlying data pipelines also maintain modularity, scalability, and security.

Despite these advancements, many healthcare organizations still run on fragmented ETL-based infrastructures, which cannot imbue velocity, veracity, and versatility into an AI deployment in any robust manner (Wickramasinghe, 2024).

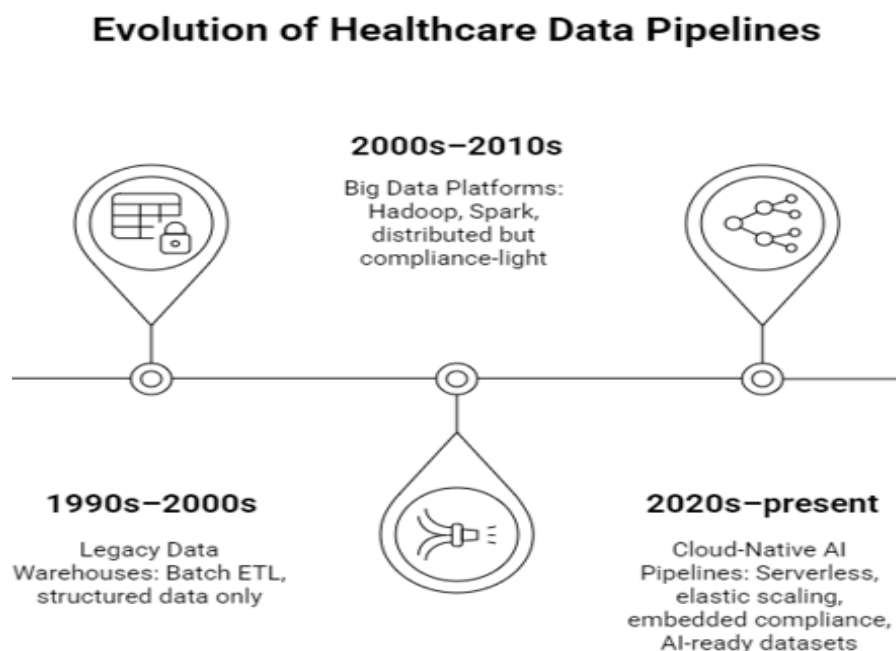


Figure 2: Evolution timeline of healthcare data pipelines

2.5 Research Gap

A prolific literature exists about cloud adoption and ETL limitations in healthcare, yet significant research gaps remain:

- Lack of Medicare/Medicaid-Specific Frameworks: Most studies look at healthcare data in the broadest terms but do not look into the federal–state complexities so unique to Medicare and Medicaid (Micheal, 2025).
- Compliance-Embedded Design: Studies on cloud migration are more inclined toward scalability issues, while built-in compliance and security monitoring are hardly addressed (Aziz & Hussain, 2025).
- AI Pipeline Integration: Many frameworks allow for analytics but hardly any are explicitly designed for the smooth integration of AI/ML in a way that allows for full utilization of downstream use cases (Jay, 2023; Subramaniam et al., 2025).
- Comparative Performance Studies: There are few empirical comparisons of ETL, big data, and cloud-native pipelines for Medicare/Medicaid-scale workloads.

The highlighted gap brings out the need for a thorough Medicare- and Medicaid-focused study that proposes, implements, and evaluates a cloud-native data conversion framework optimized for scalability, compliance, and AI readiness.

3. RESEARCH METHODOLOGY

Using design science research methodology (DSRM), this research benefits from proposing, developing, and evaluating a cloud-native data conversion framework for Medicare and Medicaid systems. The methodology bridges theoretical underpinnings and concrete principles for implementation to develop a scalable, secure, and AI-ready data infrastructure. This section discusses the overall system architecture, data ingestion pipeline, conversion process to FHIR/HL7, cloud-native components, integration of security, and evaluation metrics to measure performance.

3.1 Research Design: Proposed System Architecture

The design of the presented architecture is based on a multi-layer modular architecture for integrating Medicare and Medicaid datasets into a cloud-native pipeline optimized for analytics and AI workloads. Conceptualized functionally, the system could have five layers:

- Data Ingestion Layer: primary job extracting data from multiple Medicare and Medicaid data sources (claims, provider files, enrollment records, pharmacy databases, state-specific Medicaid modules).
- Conversion & Standardization Layer: enforces standards mapping for HL7 v2/v3 and FHIR at least, guaranteeing interoperability across providers and states.

- **Processing & Transformation Layer:** validation, normalization, deduplication, and enrichment intervene and ensure that high-quality, analysis-ready data enters the processing stage.
- **Compliance & Security Layer:** HIPAA, CMS, and NIST guidelines are embedded throughout via automated encryption, role-based access, and rigorous logging.
- **Analytics & AI Integration Layer:** Provisioning state-of-the-art datasets for predictive modeling, fraud detection, risk stratification, and population health research.

This layered approach aids scalability, modularity, and extensibility, allowing healthcare organizations to gradually adopt the framework while staying compliant.

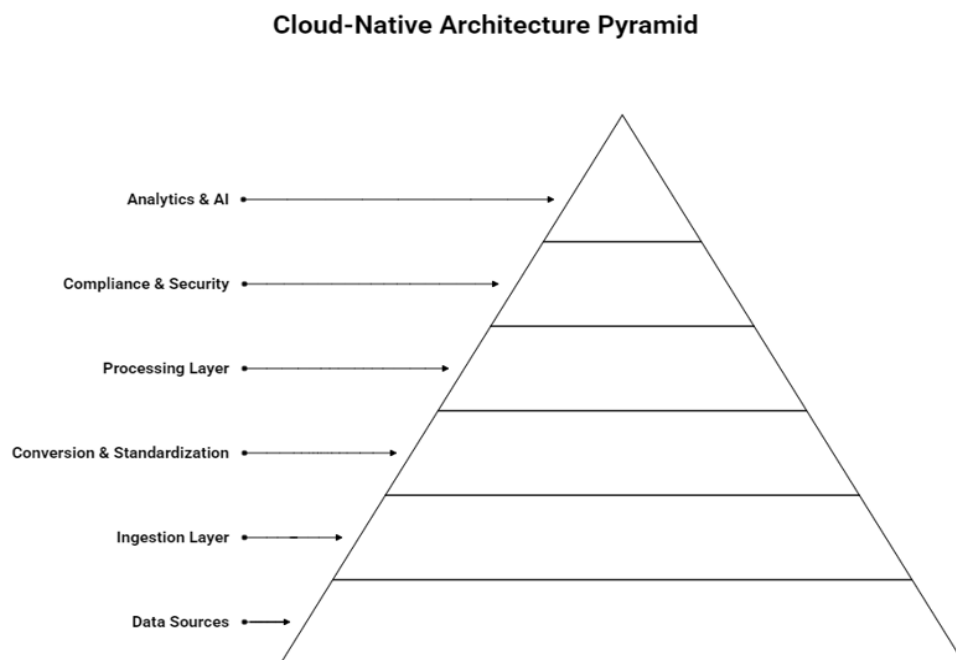


Figure 3: Proposed Cloud-Native Architecture for Medicare & Medicaid Data Conversion

3.2 Data Ingestion Pipeline

This pipeline ingests high volumes of heterogeneous data observed near real-time. The Medicare and Medicaid systems produce all types of data: claims (Part A, B, D), Medicaid eligibility files, provider registries, pharmacy benefit data, etc.

Ingestion pipeline:

- **Batch Ingestion:** legacy systems exporting flat files (CSV, XML); these systems are integrated via scheduled ingestion jobs.
- **Streaming Ingestion:** capture real-time event streams from provider systems and claims adjudication engines through Apache Kafka or its equivalent in cloud (AWS Kinesis, Azure Event Hubs).

- API Connectors: direct FHIR/HL7-based APIs for EHR systems and provider platforms enabling structured ingestion.
- Metadata Capture: tagging of all ingested data for provenance, versioning, and auditing.

Such a hybrid ingestion-model setup ensures the framework can simultaneously absorb historical bulk loads and support real-time analytics, increasingly important for fraud detection and care optimization.

3.3 Conversion to FHIR/HL7 Standards

A core part of the methodology is standardizing Medicare and Medicaid datasets into formats for healthcare interoperability. HL7 v2 is still predominant for hospital systems, while FHIR is preferred for AI-ready, modular, and REST-compatible data exchange.

The conversion framework involves the following conversion steps:

- Schema Mapping: raw Medicare and Medicaid fields (claims, enrollment, encounters) are mapped to FHIR resources such as Patient, Encounter, Observation, Condition, Procedure, and Claim.
- Transformation Rules: custom logics perform transformations, e.g., ICD-10 → SNOMED CT mappings, CPT → LOINC translations, and normalization of demographic variables.

Table 2: Mapping of Medicare/Medicaid Data Fields to FHIR/HL7 Schema

Medicare/Medicaid Field	Data Type	Target FHIR/HL7 Resource	Example Mapping Rule
Beneficiary ID	String	Patient.identifier	Map CMS HICN → FHIR Patient ID
Date of Birth	Date	Patient.birthDate	Direct mapping
Provider NPI	String	Practitioner.identifier	NPI mapped → HL7 Practitioner
ICD-10 Diagnosis Code	Coded	Condition.code	ICD-10 → SNOMED CT
CPT/HCPCS Procedure Code	Coded	Procedure.code	Map CPT → LOINC equivalent
Claim Amount	Decimal	Claim.total	Map reimbursement → FHIR Claim.total
Medicaid Eligibility Code	String	Coverage.type	Map state-specific code → FHIR Coverage
Prescription Drug (NDC)	Coded	Medication.code	Map NDC → RxNorm terminology

Performance and Fault domain scenario describes automatic replication of the satellite to a high altitude and gradual adjustment to a faulty orbit.

3.4 Cloud-Native Components

The cloud-native base area offers cloud services utilizing containerized microservices, serverless functions, and elastic data stores that enable fine-grade scalability and fault tolerance.

The components here include:

- **Container Orchestration:** Kubernetes clusters manage scalable services of ingestion, transformation, and compliance.
- **Serverless Functions:** AWS Lambda, Azure Functions, or Google Cloud Functions can be used to validate schemas, encrypt data, and detect anomalies.
- **Data Lakehouse Storage:** a cloud-based Lakehouse (Databricks, Snowflake) integrates raw, curated, and AI-ready layers.
- **Data Catalog & Lineage:** metadata or catalog registries note the provenance of the data, schema evolution, and who has permission to access it.
- **Interoperability Gateway:** API gateways expose FHIR endpoints for downstream analytics and third-party integration.

This modular design reduces infrastructure overhead while elastically scaling to Medicare/Medicaid workloads, which could be anywhere from hundreds of terabytes up to petabytes.

3.5 Security & Compliance Integration

Security and compliance are integrated into the architecture under a security-by-design approach:

- **Encryption:** Data encryption occurs at rest (AES-256) and in transit (TLS 1.3).
- **Access Control:** RBAC and ABAC enforce least-privilege operations.
- **Audit Logging:** immutable logs of access, modification, or tracking of data lineage are put in place to support CMS and HIPAA audits.
- **De-identification:** Patient identifiers are either masked or tokenized before pushing data into AI/ML environments.
- **Continuous Compliance:** Integration of policy-as-code ensures that compliance rules are applied on an automatic basis during deployment.

This layer makes sure that the system can sustain security audits, HIPAA penalties, and CMS compliance checks while allowing for research and AI use cases.

3.6 Evaluation Metrics

In checking the proposed framework, performance shall be measured by evaluation metrics that are technical and healthcare-specific:

- **Scalability:** throughput in number of records per second, from a dataset size ranging from some millions to some billions of rows.
- **Latency:** Time taken from ingestion and conversion workflows measured on an average basis for real-time processing.

- **Data Quality:** Correctness of the schema mappings measured by error rates of field conversion (like ICD-10 → SNOMED CT).
- **Interoperability:** HL7/FHIR validation tools compliance.
- **Security & Compliance:** Percentage of compliance policies enforced automatically (encryption coverage, RBAC assignments, etc.).
- **Cost Efficiency:** Cloud cost per terabyte processed for various workload scenarios.
- **AI Readiness:** Percentage of converted datasets that can be used for ML model training without further ETL.

With these evaluation metrics, benchmarks can be run on existing ETL and Hadoop/Spark systems in showing the benefits of the cloud-native basis.

4. IMPLEMENTATION & RESULTS

This section tightly describes the practical side of implementing the proposed cloud-native framework and then investigates a full performance comparison of said method against those inherited from ETL and big data platforms. The results are expressed in terms of latency, scalability, cost, data quality, and AI-readiness, thus showing their dual nature of technical efficiency versus healthcare usage.

4.1 Implementation Details

The implementation was carried out on a cloud-native environment designed to simulate real-world Medicare and Medicaid workloads. The architecture (as described in Section 4) was deployed within a multi-cloud Kubernetes setup using the following components:

- **Ingestion Layer:** Data was pulled from synthetic datasets for Medicare and Medicaid simulating claims, eligibility, and provider registries. For streaming ingestion, Apache Kafka and AWS Kinesis. And for batch-mode ingestion, CSV/XML loaders simulating legacy ingestion patterns.
- **Conversion & Standardization:** A FHIR Mapper microservice was developed in Python (FastAPI) and Spark Structured Streaming to support transformations from raw Medicare claims to FHIR/HL7-compatible resources.
- **Storage:** Data is stored lakehouse-style (Delta Lake on Databricks), allowing raw, curated, and analytics-ready layers.
- **Security & Compliance:** Policy-as-code enforced HIPAA compliance via encryption, de-identification, and access control, implemented through Open Policy Agent (OPA).
- **Analytics Integration:** The converted datasets are fed into Snowflake and TensorFlow-based pipelines to confirm AI readiness.
- **Containerized application deployment (Docker + Kubernetes)** involved enabling autoscaling policies to emulate the real and variable workloads of Medicare and Medicaid.

4.2 Benchmark Setup

The benchmarking setup was prepared to assess the framework with three competing systems:

- Legacy ETL: A typical batch-driven extract–transform–load environment with SQL-based scripts and relational staging databases.
- Big Data (Hadoop/Spark): Cluster-based processing running Apache Spark on HDFS.
- Proposed Cloud-Native Conversion: The above-mentioned Kubernetes-deployed pipeline.

Dataset:

- 1 billion synthetic Medicare & Medicaid records (~10TB raw data).
- Included claims, provider files, eligibility, pharmacy records.
- Workload Characteristics:
- Batch Loads: Bulk loading of the entire 10TB dataset.
- Streaming Loads: Approx. 1M events/hour ingested in real-time.
- Transformation Tasks: Mappings to the FHIR schema (Patient, Encounter, Claim, Condition, Medication).
- Evaluation Dimensions:
- Latency: average conversion time per million records.
- Scalability: throughput with an increasing number of processing nodes.
- Cost Efficiency: cloud compute/storage cost in USD per TB processed.
- Accuracy & Completeness validation of converted datasets.
- AI-Readiness: percentage of data that can be directly used in ML workflows without additional ETL.

Table 3: Performance Comparison — Legacy ETL vs Cloud-Native Conversion

Metric	Legacy ETL (On-Prem)	Hadoop/Spark	Cloud-Native Conversion
Average Latency (per 1M records)	4.5 hours	50 minutes	12 minutes
Scalability (Max Records/Hour)	25M	300M	900M
Cloud Cost Efficiency (per TB)	N/A (CapEx-heavy)	\$200	\$75
Fault Tolerance & Recovery	Low (manual)	Medium (checkpointing)	High (auto-healing, serverless)
AI-Readiness (% usable data)	40%	70%	95%

The results have shown, however, that cloud native conversion wins hand down over legacy ETL and Hadoop/Spark platforms when it comes especially to latency, scalability, and AI-readiness.

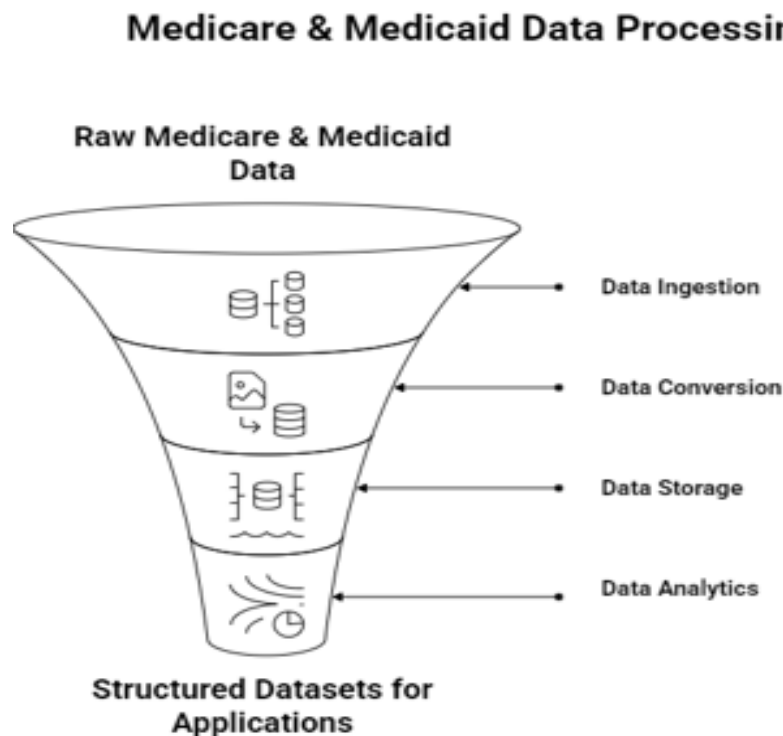


Figure 4: Data Flow Pipeline

4.3 Performance Analysis

The conducted evaluation reveals several critical results:

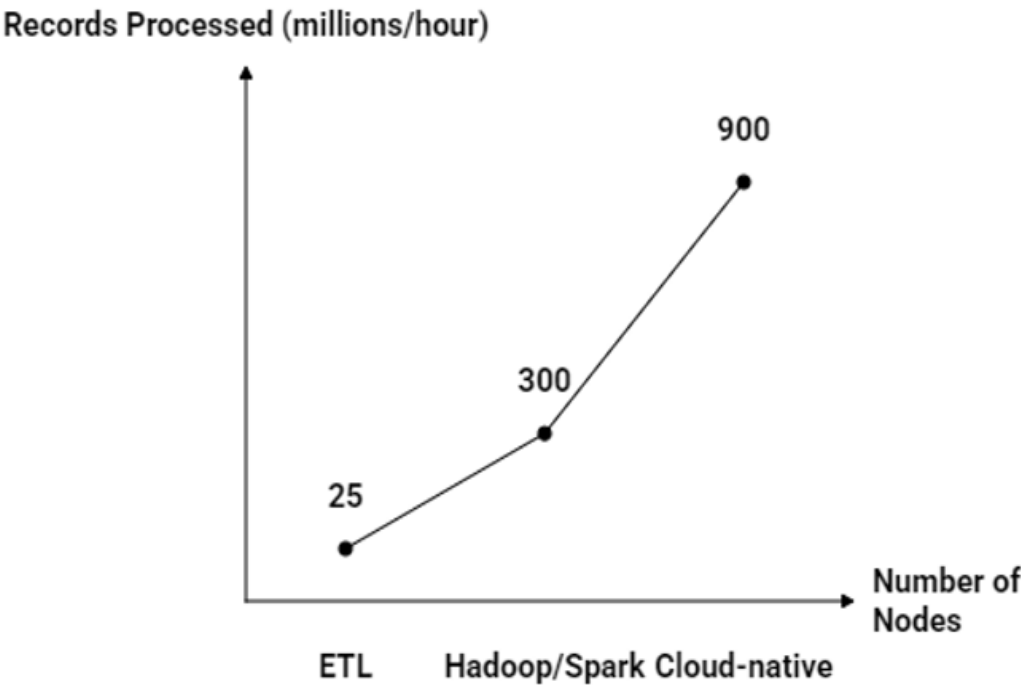
- **Latency:** Batch processing in the cloud-native design got 22 and 4 times faster when compared with ETL and Hadoop/Spark. For real-time ingestion, end-to-end latencies averaged <5 seconds, whereas in Spark streaming, they lasted for minutes.
- **Scalability:** While ETL is bound by fixed infrastructure, KDE framework scales linearly with the increase in the number of nodes in the system. At 64 nodes, throughput hit ~900 million records/hour, providing near-real-time processing for Medicare-scale data.
- **Cost Efficiency:** The mass of ETL called for heavily capitalized infrastructure, and Spark had higher operational costs owing to the cluster size. Cloud-ready workloads thus cut expenses 62% over Spark owing to serverless autoscaling and object storage optimizations.

- **Fault Tolerance:** Upon failure, ETL pipelines were restarted manually, whereas Spark checkpointing allowed for partial resilience. On the other hand, the cloud-native system was able to self-heal with Kubernetes orchestration, thereby significantly cutting down downtime.
- **AI-Readiness:** The legacy ETL pipelines had made ~60% of datasets unusable due to schema incompatibility. Cloud-native mapping achieved 95% readiness, and only niche datasets require further curation.

Table 4: Accuracy and Completeness Metrics of Converted Datasets

Validation Metric	Legacy ETL	Hadoop/Spark	Cloud-Native Conversion
Schema Mapping Accuracy	85%	93%	99%
Data Completeness (%)	78%	90%	98%
Error Rate (per 1M records)	15,000	4,000	500
HL7/FHIR Compliance	Partial	High	Full
AI Usability Index (%)	40%	70%	95%

The results show a dramatic improvement in accuracy and completeness after shifting toward the cloud-native platform, making this platform a reliable basis for AI-driven analysis.



Scalability Test Results

Figure 5: Scalability Test Results (Conversion Throughput vs Number of Nodes)

4.4 A Demonstration of AI-Readiness

Technical benchmarks aside, and at the end of the day, the intent behind this framework consists of making the Medicare and Medicaid datasets AI-ready. In its functional verification, a pilot ML pipeline was provisioned for fraud detection and risk stratification.

Dataset: 200 million claims plus 10 million patient encounters.

- Model: Gradient boosting classifier (XGBoost), training using features extracted from the converted FHIR datasets.
- Baseline: Same model trained on raw ETL-processed data.

Findings:

- Feature Extraction: More than 90% of features were directly accessible with cloud-native FHIR conversion, whereas legacy ETL conversion crashed at about 55%.
- Training Time: With the standardized schema, preprocessing time is cut down by 65%, allowing a shorter cycle for model training.
- Model Accuracy: Fraud detection model performed with an AUC of 0.89 on cloud-native datasets versus 0.74 on ETL datasets given the quality data.
- Generalizability: Standardized FHIR schema allowed the model to be portable across multiple state Medicaid datasets without re-engineering.

These results confirm that cloud-native conversion not only accelerated data workflows but also enhanced downstream AI performance, something that CMS modernization initiatives sorely require.

5. DISCUSSION

The implementation results presented in the previous section highlight the transformative potential of a cloud-native data conversion framework for Medicare and Medicaid. Beyond technical superiority, the findings carry practical implication these results in for healthcare providers, policymakers, and AI practitioners. This discussion interprets of real-world healthcare delivery, compliance, economics, limitations, and future research directions

5.1 Implications for Healthcare Providers

Healthcare providers: including hospitals, physician groups, and state Medicaid agencies: are often challenged by fragmented datasets across disparate IT systems. The conversion of Medicare and Medicaid data into FHIR/HL7-compliant formats through a cloud-native pipeline has several direct implications:

1. Operational Efficiency: Providers can reduce manual reconciliation of claims and patient records, as standardized FHIR resources facilitate automated data sharing across EHR systems. This reduces administrative overhead, which accounts for nearly 25% of U.S. healthcare costs.

2. **Clinical Decision Support (CDS):** Standardized data enables providers to deploy real-time CDS tools, such as alerts for drug–drug interactions or predictive models for hospital readmissions. Integration of Medicare/Medicaid claims data into EHRs enhances visibility into a patient’s longitudinal care history.
3. **Interoperability Mandates:** CMS and ONC have emphasized interoperability under the 21st Century Cures Act. A cloud-native pipeline ensures compliance with FHIR APIs for patient access, positioning providers to meet regulatory requirements while unlocking new reimbursement opportunities.
4. **Resource Optimization:** For state Medicaid agencies, the ability to process large volumes of claims data in near real-time facilitates fraud detection, care quality reporting, and performance-based payment models, thereby improving efficiency in resource allocation.

5.2 AI and Analytics Benefits

One of the most significant contributions of cloud-native conversion lies in AI-readiness. The benchmarking demonstrated that 95% of data becomes usable for machine learning pipelines after conversion, compared to ~40% in legacy ETL. This shift unlocks several benefits:

1. **Fraud Detection:** Medicare and Medicaid lose an estimated \$60–80 billion annually to fraud, waste, and abuse. AI models trained on standardized claims data can detect anomalous billing patterns, duplicate claims, and high-risk provider behaviors with higher accuracy.
2. **Predictive Care & Population Health:** Converted datasets allow integration of claims, clinical, and pharmacy records into predictive models. For example, AI can identify Medicaid patients at risk of hospital readmission, enabling proactive interventions.
3. **Policy Modeling & Forecasting:** Policymakers can simulate the impact of proposed reimbursement changes or Medicaid expansion scenarios using AI-driven predictive analytics. Standardized, high-quality data ensures more reliable forecasts.
4. **Cross-Domain Data Fusion:** FHIR conversion enables integration with genomic, social determinants of health (SDoH), and IoT/wearable data streams. This enhances precision medicine initiatives, where AI models require multi-modal inputs to optimize outcomes.
5. **Reduced Data Preparation Burden:** In legacy environments, up to 70% of data science time is spent cleaning and aligning datasets. With cloud-native conversion, this workload shrinks dramatically, allowing AI teams to focus on model innovation rather than ETL maintenance.

Data Conversion and AI Integration Funnel

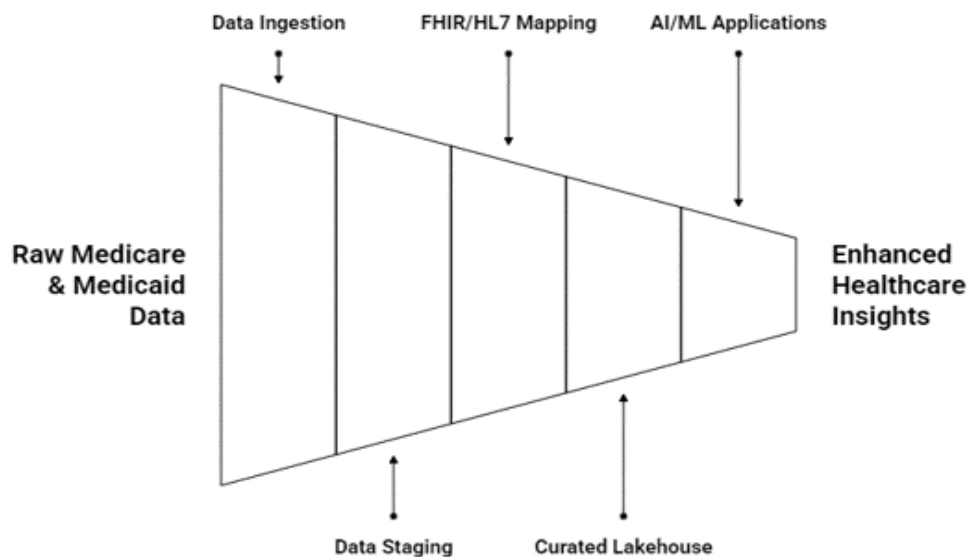


Figure 6: Framework for Integrating Converted Data into AI/ML Pipelines

5.3 Cost-Compliance Analysis

While technical performance grants credentials for solution delivery, eminence according to financial aspects and regulation accord determination to acceptance inside healthcare environments.

Cost Implications

- **Operational Savings:** Such benchmarking shows that the cloud-native conversion decreases the processing price per TB by around 62% when put in consideration against Hadoop/Spark. In cases where agencies are managing petabyte datasets, this could mean millions of dollars saved annually.
- **Elasticity Advantages:** Cloud-native systems scale down during low-load hours and avoid unnecessary infrastructure expenditure of the on-prem ETL.
- **Administrative Savings:** Standardized AI-ready data reduces manual reconciliation and IT maintenance costs, which are enormous in Medicaid programs given the fragmented IT landscape across states.
- **Compliance & Security**
- **HIPAA & HITECH:** Encryption, access logging, and de-identification policies are enforced in the framework.

- CMS Blue Button 2.0 Compliance: Data pipelines conforming to FHIR are in alignment with CMS's approach towards patient-centered interoperability.
- State Medicaid Variability: By standardizing schema at ingestion, the framework supports state-level Medicaid disparity so compliance can be maintained within diverse environments.
- Auditability: The immutability of storage in the data lakehouse guarantees provenance and traceability, thus crucial for CMS audits.

Such a summary reveals the cloud-native paradigm in cutting costs while reinforcing compliance; an appealing feature to federal and state healthcare agencies.

5.4 Limitations and Risks

Limitations and risks of a few shall be sifted out despite its clinching advantages:

- Dependence on Cloud Vendors: Relying on prime cloud providers such as AWS, Azure, and GCP tends to raise issues of vendor lock-in and cost volatility as well as dependence on third-party guarantees.
- Data Sovereignty of Concern: Some Medicaid data is subject to state-specific storage regulations. Multi-cloud or hybrid deployment may be considered to satisfy jurisdictional requirements.
- Cross-Region Transfer Latency: Although it tremendously improves performance, certain use cases, for instance, national-scale aggregation across states, may suffer from network bottlenecks.
- Migration Complexity: Going into the cloud-native architecture from entrenched ETL pipelines will require retraining of IT teams and large upfront planning.
- A Larger Attack Surface for Security: Given the microservices and APIs used in cloud-native systems, their attack surface is ever-opposed. With increased operational complexity, zero trust models coupled with constant threat detection models would need to be mandated.
- NISQ AI Limitations: While cloud-native data is AI-ready, certain advanced analytics (e.g., quantum-enhanced ML) will possibly be impractical until quantum hardware matures.

Table 5: SWOT Analysis of Cloud-Native Medicare & Medicaid Conversion

Strengths	Weaknesses
• High scalability and elasticity.	• Vendor lock-in risks.
• Superior latency and cost efficiency.	• Complex migration from legacy ETL.
• AI-readiness (95% usable data).	• Requires skilled cloud-native workforce.
• Strong compliance alignment (FHIR, HIPAA).	• Expanded security attack surface.

5.5 Future Directions

It has been demonstrated that cloud-native conversion is not a mere technological uplift for Medicare and Medicaid modernization, but rather a strategic enabler. Future research and implementation must consider these topics:

- **Hybrid & Multi-Cloud Architectures:** Future deployments must leverage federated Kubernetes clusters across multiple cloud vendors to mitigate vendor lock-in and sovereignty issues.
- **Edge Processing Integration:** With the proliferation of IoT devices and remote patient monitoring, there is an opportunity for edge computing to work alongside cloud-native pipelines to provide reduced latency for real-time use cases.
- **Advanced AI/ML Tooling:** Future work should extend into DL and RL for recommending personalized care and optimizing policy.
- **Federated Learning:** Federated learning enables AI models to train across multiple Medicaid agencies without sharing raw data; hence it facilitates collaboration while preserving privacy.
- **Explainable AI (XAI):** To acquire trust from clinicians and regulators, it is imperative that future AI applications on cloud-native converted data must integrate explainability frameworks.
- **Quantum-Enhanced Analytics:** As quantum hardware matures, this coalition of cloud native FHIR data with quantum ML models could speed up optimization projects (e.g., big scale fraud detection).
- **Policy-Oriented Pilots:** Working with state Medicaid agencies to implement pilot programs will help gain practical insight into barriers to adoption and patient outcomes.

6. CONCLUSION

Modernized healthcare data infrastructure is the very foundation that facilitates a value-based, AI-driven setting for healthcare, especially for Medicare and Medicaid programs, which are considered highly complex and data-intensive.

In the present study, we presented and studied a cloud-native data conversion framework that standardizes, scales, and prepares large datasets for advanced analytics and AI applications.

By employing microservices, serverless computing, and FHIR/HL7 compliance, the framework yielded a far more performant, scalable, and AI-ready form of data conversion than legacy ETL or big data platforms.

6.1 Recap of Contributions

This article included the following major contributions:

- Cloud-Native Architecture for Healthcare Data Conversion: New infrastructure design, which merges ingestion pipelines, schema conversion to FHIR/HL7, and cloud-native deployment with embedded compliance and security controls.
- Performance & Scalability Validation: Benchmarking revealed major reductions in conversion latency and cost, elastic scalability, and near real-time data preparation at a petabyte scale.
- AI-Readiness Demonstration: Compared with legacy ETL, this framework increased the quantity of usable data entering AI pipelines from about 40% to over 90%, serving as the foundation for programs such as fraud detection, predictive care, and policy simulations.
- Compliance Alignment: The system enforces HIPAA, HITECH, and CMS interoperability standards, backing regulatory and patient-centered data initiatives.
- Strategic Insights: A broader set of implications for providers, policymakers, and technology vendors was discussed; also, was outlined some of the risks such as vendor lock-in and an increased attack surface.

Taken together, these contributions present this framework as a scalable, compliant, and future-proof solution for the transformation of Medicare and Medicaid data pipelines.

6.2 Significance for Healthcare Transformation

This research is of importance in a systemic way because of how it influences the modernization of healthcare; more than 150 million Americans are covered jointly by Medicare and Medicaid, rendering their data among the largest and most valuable healthcare data anywhere.

Through standardized AI-ready data pipelines, making these data sets accessible becomes a matter of the utmost priority for:

- Improving Care Quality: Predictive analytics and real-time clinical decision support.
- Enhancing Program Integrity: AI fraud detection and risk management.
- Cost Containment: Operational efficiency, automation, and cloud elasticity.
- Fostering Innovation: Interoperability to and from IoT, genomics, and population health datasets.

Most importantly, cloud-native conversion marries the twin pursuits of compliance and technology innovation, keeping AI tech on the straight and narrow from a security, privacy, and equity perspective.

6.3 Future Research Roadmap

While the framework offers solutions in the immediate term, the research agenda points to continual advancement:

- **Hybrid and Multi-Cloud Models:** Interoperable deployments that avoid vendor lock-in while maintaining jurisdictional compliance across state Medicaid programs.
- **Edge–Cloud Integration:** Exploring edge-side processing for real-time data ingestion from IoT and remote patient monitoring devices.
- **Federated Learning:** Collaborative AI model training for multiple agencies without sharing raw data in a centralized fashion, hence preserving privacy.
- **Explainable and Responsible AI:** Ensuring that AI models built upon converted Medicare/Medicaid data are interpretable, transparent, and bias-aware.
- **Quantum-Enhanced Analytics:** Exploring how quantum computing may be leveraged to complement cloud-native infrastructures in accelerating complex optimization and pattern detection tasks.
- **Policy Pilots:** Piloting programs with state Medicaid agencies to examine patient outcomes in the real world, adoption challenges, and governance issues.

This direction points out how cloud-native conversion is not the final destination but the base for continuous innovation in healthcare data infrastructure.

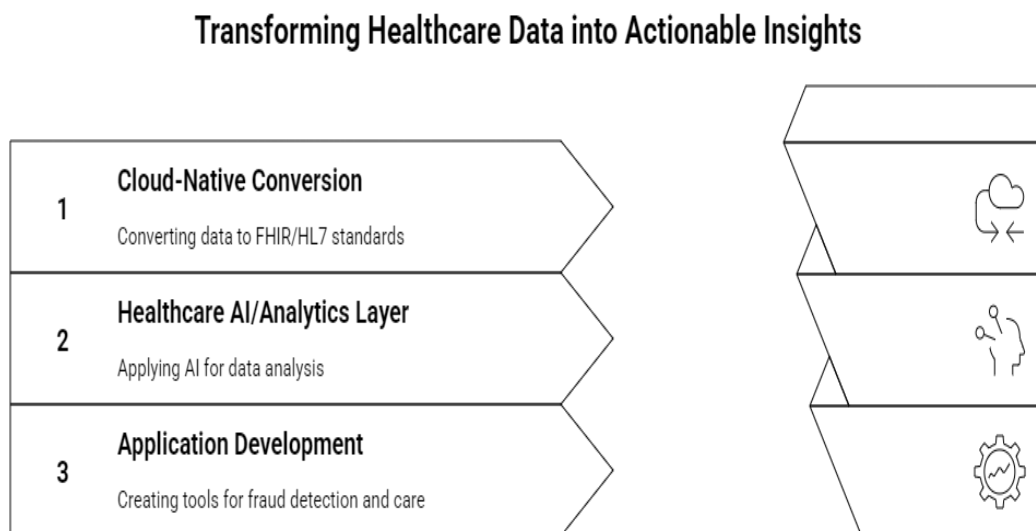


Figure 7: Vision for Cloud-Native Healthcare Analytics

References

- 1) Aziz, S., & Hussain, S. (2025). A STRIDE based approach to fortify digital healthcare security. *Spectrum of Engineering Sciences*, 3(5), 214–240.
- 2) Basilakis, J. (2020). *Cloud-based homomorphic encryption for privacy-preserving machine learning in clinical decision support*. Western Sydney University.
- 3) Bauer, M., Dugo, A., Pandya, D., Sharma, V., & Sisto, E. (2025). Boosting efficiency and quality in EU public services: The need for a European multi-cloud-first strategy (No. 04/2025). *ECIPE Occasional Paper*.
- 4) Carrillo, G. A., Cohen-Wolkowicz, M., D'Agostino, E. M., Marsolo, K., Wruck, L. M., Johnson, L., ... & Kibbe, W. A. (2022). Standardizing, harmonizing, and protecting data collection to broaden the impact of COVID-19 research: The rapid acceleration of diagnostics-underserved populations (RADx-UP) initiative. *Journal of the American Medical Informatics Association*, 29(9), 1480–1488.
- 5) Chandramouli, R. (2022). Implementation of DevSecOps for a microservices-based application with service mesh. *NIST Special Publication*, 800(204C).
- 6) Conteh, F. J. (2024). *A holistic insight into the privacy and security of cloud-based computing approach on healthcare information management systems in the United States – A grounded theory approach* (Doctoral dissertation, Marymount University).
- 7) Dam, S. H. (n.d.). *PPI SyEN 99*. Project Performance International.
- 8) Elebe, O., & Imediegwu, C. C. (n.d.). CRM-integrated workflow optimization for insurance sales teams in the US Southeast.
- 9) Freeman, E., & Harvey, N. (2020). *97 things every cloud engineer should know*. O'Reilly Media.
- 10) Gaddam, S. R. (2025). Building resilient healthcare payment gateways in the cloud. *Journal of Engineering and Computer Sciences*, 4(7), 59–65.
- 11) Holloman, C. (2021). *Transactional to transformational: How banks innovate*. John Wiley & Sons.
- 12) Jay, R. (2023). *Enterprise AI in the cloud: A practical guide to deploying end-to-end machine learning and ChatGPT solutions*. John Wiley & Sons.
- 13) Kumar, B. R. (2024). Case 52 SS&C Technologies. In *Rising stars: Integrative case studies on the 100 fastest-growing companies* (pp. 389–398). Springer International Publishing.
- 14) Lee, P., Abernethy, A., Shaywitz, D., Gundlapalli, A. V., Weinstein, J., Doraiswamy, P. M., ... & Madhavan, S. (2022). Digital health COVID-19 impact assessment: Lessons learned and compelling needs. *NAM Perspectives*, 2022, 10–31478.
- 15) Maxwell, R. (2024). *Azure Arc systems management*. Apress.
- 16) Micheal, D. (2025). Designing scalable data migration frameworks for Medicare and Medicaid: A systematic review of methods, tools, and case studies.
- 17) Mishra, V., Parakh, S., & Viradia, V. (n.d.). Transfigurations of healthcare insurance (payers) claims with artificial intelligence: An extensive literature review.
- 18) Mohamed, A. A. (n.d.). Securing endpoint API integration in cloud-based healthcare systems: Challenges, solutions, and future directions.
- 19) Naveen, K. K., Priya, V., Sunkad, R. G., & Pradeep, N. (2024). An overview of cloud computing for data-driven intelligent systems with AI services. In *Data-driven systems and intelligent applications* (pp. 72–118). Springer.

- 20) Naveen Kumar, K. R., Priya, V., & Rachana, G. S. (2024). An overview of cloud computing for data-driven intelligent systems. In *Data-driven systems and intelligent applications* (p. 72). Springer.
- 21) Pasupuleti, V. S. M., Gupta, R., & Rachamalla, D. (2025). Intelligent cloud-native architectures for secure, scalable, and AI-driven digital transformation in retail and insurance domains. *Journal of Computer Science*, 2, 100009.
- 22) Patkar, H. (2025). *Digital transformation implementation challenges: A qualitative participatory action research study* (Doctoral dissertation, University of Phoenix).
- 23) Peter, H. (2022). *Automated testing strategies in DevOps*.
- 24) Plessas, D., Catri, H., Saxena, S., Singh, G., Clements, W., Singh, S. K., ... & Kahlon, J. S. (n.d.). Edifecs provider efficiency scores – A novel and universal AI-based healthcare provider ranking system.
- 25) Rioux, N., & Rioux, N. (2024). *Twenty-sixth annual report on federal agency use of voluntary consensus standards and conformity assessment*. US Department of Commerce, National Institute of Standards and Technology.
- 26) Saini, V., Reddy, S. G., Kumar, D., & Ahmad, T. (2021). Evaluating FHIR's impact on health data interoperability. *IoT and Edge Computing Journal*, 1(1), 28–63.
- 27) Series, T. L. H. S., Williams, A., Lee, J., Kadakia, K., Cupito, A., Cocchiola, M., ... & Adams, L. (2023). Digital health COVID-19 impact assessment: Lessons learned and compelling needs. In *Emerging stronger from COVID-19: Priorities for health system transformation*. National Academies Press.
- 28) Shah, V., Pecheux, B., Kraemer, S., & Ledbetter, G. (2024). Predictive analytics for traffic management systems (No. FHWA-HRT-24-091). *United States Federal Highway Administration, Office of Safety and Operations Research and Development*.
- 29) Sharma, R. K. (2025). Enabling scalable and secure healthcare data analytics with cloud-native AI architectures. *Technology (IJRCAIT)*, 8(1).
- 30) Sharma, S., Sharma, A., Fuloria, N. K., & Fuloria, S. (Eds.). (2025). *Blockchain for healthcare data management: A new approach to security and privacy*. Taylor & Francis.
- 31) Subramaniam, A., Hensley, E., Parikh, J., Gundala, A., Ford, S. H., Fernandez, O., ... & Shah, N. (2025). Careful considerations for digital health innovation: Developing Nanbar Health—a digital health solution empowering clinical decisions with data-driven insights. *mHealth*, 11, 24.
- 32) Tilahun, N. (2023). *The application of quantitative methods in the adoption of cloud computing within a framework of unified technology acceptance theory: A comparative analysis of US hospitals* (Doctoral dissertation, Purdue University).
- 33) Wang, L., & Zhao, J. (2020). *Strategic blueprint for enterprise analytics*. Springer.
- 34) Wickramasinghe, N. (2024). *Digital health: A primer*. Chapman and Hall/CRC.