

Original Article

# End-to-End Data Pipeline Design for Medicaid Claims and Encounter Reporting

\*Rangopal Baddam

Independent Researcher, USA.

## Abstract:

Designing a reliable end-to-end data pipeline for Medicaid claims and encounter reporting has become increasingly critical as healthcare systems transition toward value-based care, interoperability, and real-time analytics. This study presents a scalable and standards-compliant pipeline architecture that supports the ingestion, transformation, validation, and reporting of Medicaid claims and encounter data across heterogeneous sources. The proposed framework integrates batch and streaming data processing using modern data engineering tools while ensuring compliance with regulatory requirements such as HIPAA and CMS reporting standards. The pipeline begins with multi-source data ingestion, incorporating Electronic Health Records (EHRs), managed care organization (MCO) submissions, and third-party billing systems. Data is processed through a layered architecture consisting of raw, curated, and analytics-ready zones, enabling efficient schema standardization and data harmonization using HL7 FHIR and X12 transaction formats. Advanced data validation mechanisms, including rule-based and machine learning-assisted anomaly detection, are implemented to improve data quality, reduce claim denials, and ensure accurate encounter reporting. A key contribution of this work is the integration of metadata-driven orchestration and automated quality checks, which significantly enhance pipeline transparency and traceability. The architecture also incorporates cloud-native technologies such as distributed storage, containerized workflows, and serverless processing to support scalability and cost optimization. Additionally, the pipeline supports near real-time reporting capabilities, allowing state Medicaid agencies and stakeholders to monitor utilization patterns, detect fraud, and improve decision-making. Evaluation of the proposed design demonstrates improvements in data latency, reporting accuracy, and operational efficiency compared to traditional legacy systems. The findings suggest that adopting a modern, interoperable data pipeline can significantly enhance Medicaid program oversight and healthcare delivery outcomes.

## Keywords:

Medicaid Claims Reporting, Encounter Data Management, Healthcare Data Pipelines, HL7 FHIR Integration, X12 Transaction Standards, CMS Compliance, Healthcare Interoperability, Cloud-Native Data Architecture, Real-Time Healthcare Analytics.

## Article History:

**Received:** 24.03.2024

**Revised:** 28.04.2024

**Accepted:** 09.05.2024

**Published:** 17.05.2024



## 1. Introduction

### 1.1. Background of Medicaid Data Ecosystems (Claims vs. Encounter Data)

The Medicaid data ecosystem is a complex, multi-stakeholder environment that captures healthcare service utilization, reimbursement activities, and patient outcomes across diverse delivery models. At its core, the ecosystem relies on two primary data types: claims data and encounter data. Claims data typically originate from fee-for-service (FFS) models, where providers submit structured billing records for reimbursement. In contrast, encounter data are generated within managed care arrangements, where providers report services delivered to beneficiaries under capitated payment models.

While both data types serve similar analytical and regulatory purposes, encounter data are often less standardized and more prone to inconsistencies due to variations in managed care organization (MCO) reporting practices. The increasing shift toward managed care in Medicaid programs has amplified the importance of encounter data, making its quality and completeness critical for policy evaluation, risk adjustment, and federal reporting requirements (CMS, 2023).

### 1.2. Importance of Data Pipelines in Healthcare Analytics and Reimbursement

Modern healthcare systems depend heavily on robust data pipelines to process large volumes of structured and semi-structured data efficiently. In the context of Medicaid, data pipelines enable the seamless flow of information from source systems such as Electronic Health Records (EHRs), billing platforms, and third-party administrators into centralized repositories for validation, transformation, and analysis.

Efficient pipeline design directly impacts reimbursement accuracy, fraud detection, population health analytics, and policy decision-making. Poorly designed pipelines can lead to delayed payments, inaccurate reporting, and compliance risks. Conversely, scalable and automated pipelines support real-time analytics, improve data quality, and facilitate value-based care initiatives by ensuring that high-quality data is available for timely decision-making (Zaharia et al., 2016; Rajkomar et al., 2019).

### 1.3. Role of the Centers for Medicare & Medicaid Services (CMS)

The Centers for Medicare & Medicaid Services plays a central role in regulating, standardizing, and monitoring Medicaid data reporting across states. CMS mandates the submission of Medicaid data through systems such as the Transformed Medicaid Statistical Information System (T-MSIS), which standardizes data elements and improves national-level visibility into Medicaid programs.

CMS also provides frameworks, toolkits, and compliance guidelines to ensure data accuracy, completeness, and timeliness. These include encounter data reporting requirements, data quality metrics, and interoperability standards aligned with initiatives such as HL7 FHIR and the Trusted Exchange Framework and Common Agreement (TEFCA). Through these efforts, CMS aims to enhance transparency, support program integrity, and enable data-driven healthcare policy (CMS, 2021; ONC, 2022).

### 1.4. Problem Statement: Fragmentation, Data Quality Issues, and Reporting Delays

Despite ongoing modernization efforts, Medicaid data systems continue to face significant challenges. One of the primary issues is data fragmentation, where information is distributed across multiple, often incompatible systems operated by different stakeholders. This fragmentation leads to inconsistencies in data formats, incomplete records, and limited interoperability.

Data quality issues further compound the problem. Encounter data, in particular, often suffer from missing fields, coding errors, and delayed submissions, which undermine their reliability for analytics and reporting. Additionally, legacy systems and manual processes contribute to reporting delays, preventing timely insights and hindering effective program management. These challenges not only impact operational efficiency but also pose risks to compliance with federal reporting standards and reimbursement accuracy (Kahn et al., 2016; CMS, 2023).

### 1.5. Research Objectives and Scope

This study aims to address the aforementioned challenges by proposing a comprehensive design for an end-to-end data pipeline tailored to Medicaid claims and encounter reporting. The key objectives include:

- Designing a unified architecture that supports both claims and encounter data processing
- Enhancing data quality through automated validation and anomaly detection mechanisms
- Reducing latency in data processing and reporting

- Ensuring compliance with CMS and HIPAA regulations
- Supporting scalable and flexible data integration across heterogeneous systems

The scope of this research focuses on architectural design principles, implementation strategies, and evaluation metrics relevant to modern healthcare data pipelines. It does not aim to replace existing CMS frameworks but rather to complement them with improved technical solutions.

### 1.6. Focus on Pipeline Design, Compliance, Scalability, and Interoperability

The proposed approach emphasizes four critical dimensions: pipeline design, compliance, scalability, and interoperability. A well-structured pipeline design ensures efficient data ingestion, transformation, and storage. Compliance is addressed through adherence to regulatory standards such as HIPAA, CMS reporting requirements, and data governance policies.

Scalability is achieved through cloud-native technologies, distributed processing frameworks, and modular architectures that can handle increasing data volumes and complexity. Interoperability is facilitated through standardized data formats such as HL7 FHIR and X12, enabling seamless data exchange across systems and stakeholders.

By integrating these dimensions, the research seeks to provide a practical and future-ready solution for improving Medicaid data management, ultimately supporting better healthcare outcomes and more efficient program administration.

## 2. Overview of Medicaid Claims and Encounter Data

### 2.1. Definition and Differences: Claims vs. Encounter Data

Medicaid data reporting is fundamentally built on two distinct yet complementary data types: claims data and encounter data, each reflecting different reimbursement and care delivery models.

Claims data are generated within fee-for-service (FFS) arrangements, where healthcare providers submit detailed billing records to receive payment for services rendered. These records typically include procedure codes, diagnosis codes, service dates, provider identifiers, and cost information. Because reimbursement is directly tied to each submitted claim, claims data are generally well-structured, standardized, and subject to strict validation processes. This makes them relatively reliable for financial reconciliation, utilization tracking, and audit purposes.

In contrast, encounter data are produced under managed care models, where providers are reimbursed through capitated payments rather than per-service billing. Managed Care Organizations (MCOs) are responsible for collecting and submitting records of services delivered to beneficiaries. Although encounter data mirror claims in structure, they are not directly tied to immediate reimbursement, which often results in variability in completeness, timeliness, and accuracy. Despite these challenges, encounter data are essential for risk adjustment, quality measurement, rate setting, and federal oversight, particularly as Medicaid programs increasingly shift toward managed care delivery models (CMS, 2023).

### 2.2. Standard Formats Used in Medicaid Data Exchange

To ensure interoperability and consistency across diverse healthcare systems, Medicaid claims and encounter data are transmitted using standardized electronic transaction formats. The most widely used formats are based on the X12 and NCPDP standards:

- 837I (Institutional Claims): Used for billing institutional services such as hospital stays, skilled nursing facilities, and inpatient care.
- 837P (Professional Claims): Designed for outpatient and professional services, including physician visits, laboratory tests, and outpatient procedures.
- NCPDP (Pharmacy Claims): A specialized standard for pharmacy transactions, capturing prescription drug dispensing and related information.
- Dental Formats (e.g., 837D): Used for dental services, including preventive, diagnostic, and restorative procedures.

These standardized formats ensure that data elements such as diagnosis codes (ICD), procedure codes (CPT/HCPCS), and provider identifiers are consistently represented, enabling automated processing and integration across systems. However, despite the use of

common standards, implementation variations across states and organizations can still introduce inconsistencies that complicate downstream processing and analytics.

### 2.3. Submission Requirements and Reporting Timelines

Medicaid programs enforce strict submission requirements to ensure timely and accurate reporting of healthcare services. Both claims and encounter data must be submitted within defined timeframes, typically within 180 days ( $\leq 180$  days) after the date of service, although specific timelines may vary by state and program.

For managed care programs, encounter data are usually submitted on a monthly basis, with MCOs required to aggregate and transmit service records to state Medicaid agencies. These submissions are then further processed and reported to federal systems such as the Transformed Medicaid Statistical Information System (T-MSIS).

Timely submission is critical for several reasons. It ensures accurate rate setting, supports real-time monitoring of healthcare utilization, and enables compliance with federal reporting mandates. Delays in submission can lead to incomplete datasets, inaccuracies in actuarial calculations, and potential financial penalties for non-compliance.

### 2.4. Regulatory Framework Governing Encounter Data

The reporting of Medicaid encounter data is governed by federal regulations, particularly 42 CFR § 438, which establishes requirements for managed care programs. These regulations mandate that states collect, validate, and submit complete and accurate encounter data to support program oversight and payment integrity.

Key provisions under 42 CFR § 438 include:

- Ensuring that encounter data are accurate, complete, and timely
- Implementing validation processes to verify data reliability
- Using encounter data for rate setting, quality measurement, and program evaluation
- Maintaining audit trails and documentation to support federal reviews

The Centers for Medicare & Medicaid Services enforces these requirements and provides guidance to states on data quality standards, reporting methodologies, and compliance expectations. Non-compliance can result in corrective action plans, financial penalties, or adjustments to federal funding.

### 2.5. Challenges: Inconsistency across Data Sources and Reporting Systems

Despite the presence of standardized formats and regulatory frameworks, Medicaid data systems continue to face significant challenges, particularly related to inconsistency and fragmentation. One major issue is the heterogeneity of data sources, which include EHR systems, billing platforms, third-party administrators, and MCO data warehouses. Each of these systems may use different configurations, coding practices, and data validation rules, leading to discrepancies in reported data.

Additionally, variations in state-level implementations and reporting requirements further exacerbate inconsistencies. Encounter data are especially vulnerable to issues such as missing fields, incorrect coding, duplicate records, and delayed submissions. These problems reduce data reliability and limit the effectiveness of analytics, policy evaluation, and reimbursement processes.

Another critical challenge lies in the lack of seamless interoperability between systems. Even when standardized formats are used, differences in interpretation and implementation can hinder data integration. As a result, significant manual effort is often required to reconcile discrepancies, increasing operational costs and delaying reporting timelines.

Addressing these challenges requires not only improved technical solutions such as robust data pipelines and automated validation but also stronger governance frameworks and alignment across stakeholders.

Overview of Medicaid Claims and Encounter Data

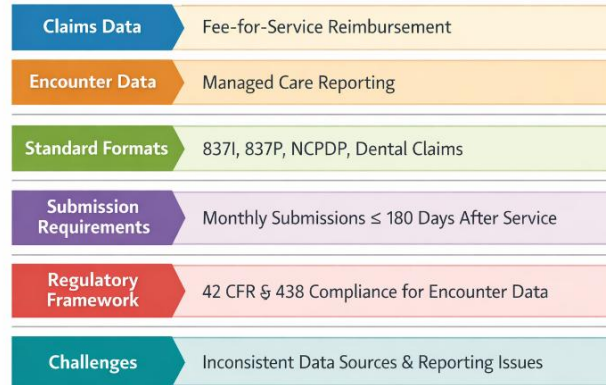


Figure 1. Overview of Medicaid and Encounter Data

Table 1. Overview of Medicaid Claims and Encounter Data

Category	Claims Data (Fee-for-Service)	Encounter Data (Managed Care)
Definition	Billing records submitted by providers for each service rendered and reimbursed individually	Records of services delivered under capitated payment models, reported by Managed Care Organizations (MCOs)
Primary Purpose	Direct reimbursement and financial transactions	Monitoring service utilization, quality measurement, and regulatory reporting
Payment Model	Fee-for-Service (FFS)	Managed Care (Capitated payments)
Data Reliability	Generally high due to direct linkage with payment validation	Often variable due to indirect linkage with reimbursement
Standard Formats	837I (Institutional), 837P (Professional), 837D (Dental), NCPDP (Pharmacy)	Same formats as claims but with additional reporting layers and aggregation by MCOs
Submission Frequency	Typically continuous or periodic submissions tied to billing cycles	Monthly submissions aggregated by MCOs
Submission Timeline	Usually ≤180 days after date of service (state-dependent)	Typically ≤180 days after service, with stricter compliance monitoring
Regulatory Requirement	CMS billing and auditing guidelines	Governed under 42 CFR § 438 for completeness, accuracy, and timeliness
Use Cases	Reimbursement processing, auditing, utilization tracking	Risk adjustment, rate setting, program evaluation, policy analysis
Data Challenges	Occasional coding errors or claim denials	Incomplete records, delays, inconsistencies across MCOs and systems
Interoperability Issues	Moderate (standardized but varies by implementation)	High (due to multiple stakeholders and reporting variations)

3. Regulatory and Interoperability Landscape

The regulatory and interoperability landscape surrounding Medicaid data reporting has evolved significantly over the past decade, driven by the need for transparency, accountability, and seamless data exchange across healthcare systems. Federal oversight primarily led by the Centers for Medicare & Medicaid Services has introduced structured guidelines and interoperability mandates that directly influence how Medicaid claims and encounter data are collected, validated, and shared.

These regulatory frameworks are not merely compliance checkpoints; they fundamentally shape the design of modern healthcare data pipelines. As healthcare systems transition toward real-time analytics and patient-centered care, interoperability standards such as HL7 FHIR and API-driven architectures have become essential components of Medicaid data infrastructure.

### 3.1. CMS Guidelines for Encounter Data Validation and Reporting

CMS has established comprehensive guidelines to ensure that Medicaid encounter data meet strict standards for accuracy, completeness, and timeliness. These guidelines require state Medicaid agencies and Managed Care Organizations (MCOs) to implement systematic validation processes before submitting data to federal systems such as T-MSIS.

Validation typically occurs across multiple layers, including:

- Syntactic validation, ensuring that data conform to required formats (e.g., X12 standards)
- Semantic validation, verifying that coded values (ICD, CPT, HCPCS) are clinically and logically consistent
- Referential integrity checks, ensuring consistency across related data elements (e.g., patient IDs, provider identifiers)

CMS also mandates ongoing monitoring of data quality through predefined metrics, such as error rates, missing data percentages, and timeliness indicators. These measures are critical for maintaining the integrity of Medicaid reporting and supporting downstream use cases such as rate setting and program evaluation (CMS, 2023).

### 3.2. Medicaid Managed Care Encounter Data Toolkit Overview

The Medicaid Managed Care Encounter Data Toolkit, developed by CMS, serves as a practical guide for states to improve encounter data collection, validation, and submission processes. The toolkit outlines best practices for:

- Establishing standardized data submission workflows
- Implementing automated validation and error correction mechanisms
- Enhancing collaboration between states and MCOs
- Monitoring data quality through dashboards and reporting tools

A key feature of the toolkit is its emphasis on data governance and accountability, encouraging states to define clear roles and responsibilities for data submission and validation. It also provides recommendations for aligning encounter data processes with federal reporting requirements, thereby reducing inconsistencies and improving overall data reliability.

From a pipeline design perspective, the toolkit reinforces the importance of embedding validation and quality checks directly into data processing workflows rather than treating them as post-processing activities.

### 3.3. Data Validation, Audit Requirements, and Completeness Standards

Regulatory compliance extends beyond initial data submission to include rigorous audit and validation requirements. CMS requires states to conduct regular audits of encounter data to ensure that submitted records accurately reflect services delivered. These audits may involve:

- Cross-referencing encounter data with source systems such as EHRs and claims records
- Verifying provider credentials and service eligibility
- Assessing data completeness across critical fields (e.g., diagnosis codes, service dates, beneficiary identifiers)

Completeness standards are particularly important, as missing or incomplete data can compromise the reliability of analytics and financial calculations. CMS defines thresholds for acceptable levels of completeness and may impose corrective actions if states fail to meet these standards.

These requirements necessitate the integration of continuous data quality monitoring and audit trails within data pipelines, ensuring that every transformation and validation step is traceable and compliant with federal guidelines (Kahn et al., 2016; CMS, 2021).

### 3.4. 2024 CMS Interoperability Rule (CMS-0057-F)

The CMS Interoperability and Prior Authorization Final Rule (CMS-0057-F), finalized in 2024, represents a major advancement in healthcare data exchange. This rule mandates the adoption of FHIR-based APIs to facilitate standardized, secure, and real-time data sharing across healthcare stakeholders.

Key API requirements include:

- Patient Access API: Enables beneficiaries to access their health data through third-party applications
- Provider Access API: Allows providers to retrieve patient data for care coordination
- Payer-to-Payer API: Supports seamless data exchange between health plans, particularly during coverage transitions
- Prior Authorization API: Streamlines authorization workflows, reducing administrative burden and improving turnaround times

These APIs are built on the HL7 FHIR standard, which provides a flexible and interoperable framework for representing healthcare data. The rule also emphasizes security, privacy, and standardized implementation to ensure consistent adoption across states and organizations (CMS, 2024).

### 3.5. Implications for Pipeline Design

#### 3.5.1. Real-Time Data Exchange

The shift toward FHIR-based APIs introduces a paradigm change from batch-oriented processing to near real-time data exchange. Traditional pipelines designed for periodic submissions must evolve to support continuous data ingestion and streaming architectures. This requires the integration of event-driven systems, message queues, and low-latency processing frameworks to handle real-time data flows efficiently.

#### 3.5.2. API-First Architecture

The interoperability rule strongly encourages an API-first approach, where data pipelines are designed around standardized interfaces rather than internal data structures. This approach enhances modularity, scalability, and ease of integration with external systems. APIs act as the primary mechanism for data access and exchange, enabling seamless communication between providers, payers, and patients.

#### 3.5.3. Longitudinal Data Sharing

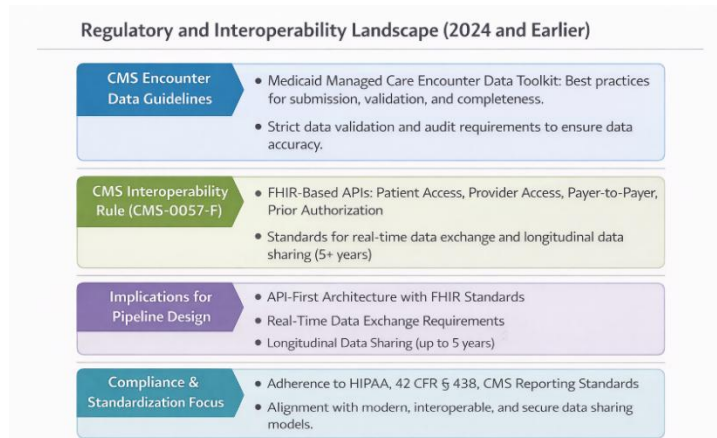
One of the most significant requirements of CMS-0057-F is the support for longitudinal data sharing, allowing up to five years of patient data to be exchanged between payers. This necessitates robust data storage and retrieval mechanisms capable of handling large volumes of historical data while maintaining data integrity and accessibility.

From a pipeline perspective, this requirement introduces challenges related to:

- Data retention and archival strategies
- Efficient indexing and querying of historical records
- Ensuring consistency across longitudinal datasets

### 3.6. Concluding Perspective

Overall, the regulatory and interoperability landscape up to 2024 underscores a clear shift toward standardization, transparency, and real-time data accessibility. Compliance with CMS guidelines and interoperability rules is no longer optional but a foundational requirement for modern Medicaid data systems. For pipeline designers, this means moving beyond traditional batch processing models and adopting architectures that are scalable, API-driven, and deeply integrated with validation and governance mechanisms. By aligning technical design with regulatory expectations, organizations can not only ensure compliance but also unlock the full potential of Medicaid data for analytics, policy development, and improved healthcare outcomes.



**Figure 2. Regulatory and Interoperability Landscape (2024 and Earlier)**

**Table 2. Regulatory and Interoperability Landscape for Medicaid**

Category	Key Components	Description	Implications for Data Pipeline Design
CMS Encounter Data Guidelines	Accuracy, Completeness, Timeliness	CMS mandates that encounter data must meet strict quality standards before submission to systems like T-MSIS	Requires built-in validation layers, data quality checks, and automated error handling within pipelines
Encounter Data Toolkit (CMS)	Submission workflows, validation rules, governance practices	Provides best practices for states and MCOs to standardize data collection, validation, and reporting	Encourages metadata-driven pipelines and standardized ingestion/processing frameworks
Data Validation & Audit Requirements	Syntactic, semantic, and referential validation; audit trails	Ensures data integrity through multi-level validation and periodic audits	Necessitates audit logging, traceability, and continuous monitoring mechanisms in pipeline architecture
Completeness Standards	Mandatory fields, error thresholds, reporting metrics	Defines acceptable levels of missing or incomplete data and enforces corrective actions	Requires completeness checks, data profiling, and automated reconciliation processes
Regulatory Framework (42 CFR § 438)	Managed care reporting requirements	Governs encounter data submission, validation, and usage for rate setting and program oversight	Drives compliance-focused pipeline design with enforceable validation rules and reporting layers
CMS Interoperability Rule (CMS-0057-F, 2024)	FHIR-based API mandates	Introduces standardized APIs for secure and real-time data exchange across stakeholders	Requires API integration layers and FHIR-compatible data models within pipelines
FHIR-Based APIs	Patient Access, Provider Access, Payer-to-Payer, Prior Authorization	Enables standardized, interoperable data exchange using HL7 FHIR resources	Promotes API-first architecture and modular pipeline components for data access and sharing
Real-Time Data Exchange	Event-driven data sharing	Moves from batch processing to near real-time data availability	Requires streaming pipelines, message queues, and low-latency processing systems

Longitudinal Data Sharing	Up to 5 years of patient data exchange	Supports continuity of care and historical data analysis across payers	Requires scalable storage, indexing, and efficient retrieval mechanisms for historical datasets
Compliance & Security Standards	HIPAA, CMS reporting policies	Ensures privacy, security, and regulatory compliance in data handling	Requires encryption, access controls, and secure data transmission protocols

## 4. End-to-End Data Pipeline Architecture

An end-to-end data pipeline for Medicaid claims and encounter reporting represents the backbone of modern healthcare data infrastructure. It enables the seamless movement of data from diverse source systems to regulatory reporting platforms while ensuring compliance, scalability, and high data quality. Unlike traditional fragmented systems, a well-designed pipeline integrates ingestion, transformation, validation, storage, and reporting into a unified architecture.

Given the complexity of Medicaid ecosystems characterized by heterogeneous data sources, strict compliance requirements, and increasing interoperability mandates pipeline design must balance performance, traceability, and regulatory alignment. This section presents a layered architecture that supports both batch and real-time processing while aligning with CMS standards and interoperability frameworks such as HL7 FHIR.

### 4.1. Data Sources

The pipeline begins with diverse upstream data sources, each contributing different types of healthcare and administrative data. These sources vary in structure, frequency, and quality, making integration a non-trivial challenge.

- **Provider EHR Systems:** Electronic Health Records capture clinical data such as diagnoses, procedures, lab results, and patient demographics. These systems often follow HL7 or FHIR standards but may differ in implementation across providers.
- **Billing Systems:** These systems generate structured claims data (e.g., 837 transactions) used for reimbursement. They are typically well-standardized but may still contain coding inconsistencies.
- **Pharmacy Systems:** Pharmacy data include prescription fills, drug codes (NDC), and dispensing information. These systems often use NCPDP standards and are critical for medication tracking and cost analysis.
- **Managed Care Organization (MCO) Feeds:** MCOs aggregate encounter data from multiple providers and submit them to state Medicaid agencies. These feeds are often heterogeneous and may require extensive preprocessing.

The diversity of these sources necessitates a flexible ingestion framework capable of handling multiple formats and transmission protocols.

### 4.2. Data Ingestion Layer

The ingestion layer is responsible for capturing data from upstream systems and delivering it into the pipeline in a consistent and reliable manner.

- **Batch Ingestion:** Traditional Medicaid workflows rely heavily on batch ingestion of EDI transactions such as 837I, 837P, and NCPDP files. These files are typically transferred via secure protocols (e.g., SFTP) and processed in scheduled intervals.
- **Streaming Ingestion:** With the adoption of FHIR APIs and interoperability mandates, streaming ingestion has become increasingly important. Event-driven systems enable near real-time data capture, supporting faster analytics and decision-making.
- **Validation at Ingestion Stage:** Early-stage validation is critical to prevent downstream errors. This includes format validation, schema checks, and basic business rule enforcement (e.g., required fields, valid codes).

A hybrid ingestion model combining batch and streaming ensures compatibility with legacy systems while enabling modern real-time capabilities.

### 4.3 Data Transformation and Standardization

Once ingested, data must be transformed into a standardized format suitable for integration and analysis.

- **Code Mapping:** Healthcare data rely on multiple coding systems, including ICD (diagnoses), CPT/HCPCS (procedures), and NDC (drug codes). These codes often need to be mapped to standardized vocabularies such as RxNorm or ATC to support interoperability and analytics.
- **Data Normalization and Schema Alignment:** Data from different sources are normalized into a unified schema, ensuring consistency in field names, data types, and relationships.
- **Example – Drug Classification Pipelines:** In Medicaid claims processing, drug-related data are often transformed through classification pipelines that map NDC codes to therapeutic classes (e.g., ATC levels). This enables cost analysis, utilization tracking, and formulary management.

This layer plays a critical role in ensuring that downstream processes operate on consistent and interoperable datasets.

#### 4.4. Data Storage Layer

The storage layer provides the foundation for managing large volumes of healthcare data while supporting both raw data retention and structured analytics.

- **Data Lake (Raw Ingestion):** A data lake stores raw, unprocessed data in its original format. This allows for flexibility in reprocessing and auditing.
- **Data Warehouse (Structured Reporting):** Processed and standardized data are stored in a data warehouse optimized for querying and reporting. This layer supports regulatory submissions and analytics.
- **Master Data Management (MDM):** MDM ensures consistency of key entities such as patients, providers, and organizations across the pipeline. It resolves duplicates and maintains a single source of truth.

Together, these storage components enable both operational efficiency and regulatory compliance.

#### 4.5. Data Processing and Enrichment

After storage, data undergo further processing to enhance its value and usability.

- **Deduplication and Record Linkage:** Duplicate records are identified and merged using deterministic or probabilistic matching techniques. Record linkage ensures that data from different sources are correctly associated with the same entity.
- **Risk Adjustment Enrichment:** Data are enriched with risk scores based on patient conditions and utilization patterns. These scores are critical for payment models and policy analysis.
- **Quality Scoring and Validation:** Additional metrics are computed to assess data quality, including completeness, consistency, and anomaly detection.

This layer transforms raw data into actionable insights while improving its reliability.

#### 4.6. Data Validation and Quality Assurance

Validation is a continuous process throughout the pipeline but is particularly critical before data submission.

- **Completeness Checks:** Ensure that all required fields are populated.
- **Accuracy Checks:** Validate that data values are correct and consistent with clinical and business rules.
- **Timeliness Checks:** Verify that data are submitted within required timeframes (e.g.,  $\leq 180$  days).
- **Audit Pipelines:** Maintain detailed logs and traceability for each data element, aligning with CMS validation requirements and audit standards.

Embedding validation into the pipeline reduces errors, improves compliance, and enhances trust in the data.

#### 4.7 Data Submission Layer

The submission layer prepares and transmits data to regulatory systems.

- **Formatting into CMS-Compliant EDI Transactions:** Data are converted into standardized formats such as 837I, 837P, or NCPDP for submission.
- **Submission to Encounter Data Systems (EDS):** Processed data are transmitted to state and federal systems, including T-MSIS.

- Error Handling and Resubmission Workflows: Errors identified during submission are logged, corrected, and resubmitted through automated workflows.

This layer ensures that data meet all regulatory requirements and are successfully delivered to CMS systems.

#### 4.8. Reporting and Analytics Layer

The final layer focuses on deriving insights and supporting decision-making.

- Regulatory Reporting Dashboards: Provide visibility into compliance metrics, submission status, and data quality indicators.
- Operational Analytics: Support day-to-day operations, including utilization monitoring and provider performance analysis.
- Financial Reconciliation and Risk Scoring: Enable accurate reimbursement calculations and risk-based payment models.

This layer transforms processed data into meaningful insights that drive policy, operational improvements, and financial outcomes.

#### 4.9. Concluding Perspective

The end-to-end pipeline architecture described above reflects a shift from siloed, batch-oriented systems to integrated, scalable, and interoperable data platforms. By aligning each layer with regulatory requirements and interoperability standards, organizations can build pipelines that not only ensure compliance but also unlock the full analytical potential of Medicaid data.



Figure 3. End-to-End Data Pipeline Architecture for Medicaid Claims and Encounter Reporting

Table 3. End-to-End Medicaid Data Pipeline Architecture

Pipeline Layer	Key Components	Core Functions	Technologies / Standards	Design Considerations
Data Sources	EHR systems, Billing systems, Pharmacy systems, MCO feeds	Capture clinical, financial, and administrative data	HL7, FHIR, X12, NCPDP	High heterogeneity, varying data quality, multiple formats

		from multiple upstream systems		
Data Ingestion Layer	Batch ingestion (EDI), Streaming ingestion (APIs/events), Initial validation	Ingest data from multiple sources and perform early-stage validation	837I, 837P, Kafka, FHIR APIs, SFTP	Hybrid ingestion model (batch + streaming), low-latency ingestion
Data Transformation & Standardization	Code mapping, Schema alignment, Normalization pipelines	Convert raw data into standardized formats and harmonize schemas	ICD, CPT, HCPCS, NDC → RxNorm, ATC	Ensuring interoperability, consistent schema design
Data Storage Layer	Data lake, Data warehouse, Master Data Management (MDM)	Store raw and processed data, maintain master records	Cloud storage (S3, ADLS), SQL warehouses	Scalability, data partitioning, entity resolution
Data Processing & Enrichment	Deduplication, Record linkage, Risk adjustment, Quality scoring	Enhance data with derived metrics and resolve duplicates	ML models, probabilistic matching, risk scoring algorithms	Accuracy vs. performance trade-offs, enrichment consistency
Data Validation & Quality Assurance	Completeness checks, Accuracy validation, Timeliness monitoring, Audit pipelines	Ensure compliance with CMS data quality standards	CMS validation rules, data profiling tools	Continuous monitoring, traceability, audit readiness
Data Submission Layer	EDI formatting, Submission to EDS/T-MSIS, Error handling workflows	Prepare and submit data to regulatory systems	X12 (837), NCPDP, CMS EDS	Automated error correction, resubmission pipelines
Reporting & Analytics Layer	Dashboards, Operational analytics, Financial reconciliation, Risk scoring	Generate insights for compliance, operations, and decision-making	BI tools (Power BI, Tableau), SQL, ML analytics	Real-time reporting, scalability, user accessibility

### 5. Technical Design Considerations

Designing a robust end-to-end data pipeline for Medicaid claims and encounter reporting requires careful attention to several technical dimensions that directly influence performance, reliability, and compliance. Given the scale and complexity of Medicaid ecosystems where millions of transactions are processed across multiple stakeholders pipeline architectures must be engineered to handle high data volumes, ensure fault tolerance, maintain traceability, and support evolving interoperability standards.

This section outlines the critical technical considerations that underpin a resilient and future-ready pipeline, emphasizing scalability, reliability, governance, and seamless integration between legacy and modern systems.

#### 5.1. Scalability (Handling Millions of Claims/Encounters)

One of the foremost challenges in Medicaid data processing is the sheer volume of data generated ежедневно across providers, payers, and managed care organizations. A modern pipeline must be capable of handling millions of claims and encounter records daily, with the ability to scale horizontally as data volumes grow.

Scalability is typically achieved through distributed computing frameworks and cloud-based infrastructure that allow workloads to be dynamically scaled based on demand. Key strategies include:

- Horizontal scaling: Adding more compute nodes to distribute processing workloads
- Partitioning and sharding: Splitting large datasets into manageable segments for parallel processing
- Elastic resource allocation: Automatically scaling resources up or down based on workload intensity

Without proper scalability, pipelines can experience bottlenecks, increased latency, and system failures, ultimately affecting reporting timelines and compliance.

## 5.2. Fault Tolerance and Retry Mechanisms

Given the critical nature of healthcare data, pipelines must be designed to handle failures gracefully without data loss. Fault tolerance ensures that the system can continue operating even in the presence of hardware, network, or application-level failures.

Key mechanisms include:

- Checkpointing: Periodically saving the state of data processing to allow recovery from failures
- Retry logic: Automatically reprocessing failed transactions with configurable retry policies
- Dead-letter queues: Isolating problematic records for later analysis without disrupting the pipeline
- Redundancy: Replicating data and services across multiple nodes or regions

These mechanisms are essential for maintaining data integrity and ensuring that no critical records are lost during processing or submission.

## 5.3. Data Lineage and Traceability

In a highly regulated environment like Medicaid, it is crucial to maintain a clear record of how data flows through the pipeline. Data lineage refers to the ability to trace data from its origin to its final destination, including all transformations and validations applied along the way.

Traceability supports:

- Regulatory compliance: Demonstrating adherence to CMS and audit requirements
- Error resolution: Identifying the source of discrepancies or anomalies
- Transparency: Providing stakeholders with visibility into data processing workflows

Implementing lineage requires capturing metadata at each stage of the pipeline, including timestamps, transformation logic, and data provenance. This ensures that every data element can be audited and verified.

## 5.4. Metadata Management

Metadata management is a foundational component of modern data pipelines, enabling efficient organization, governance, and discovery of data assets. Metadata includes information about data structure, source, transformations, quality metrics, and usage.

Effective metadata management supports:

- Schema management: Tracking changes in data structures over time
- Data cataloging: Providing searchable inventories of datasets
- Governance and compliance: Ensuring adherence to regulatory standards
- Automation: Driving pipeline orchestration and validation processes

A metadata-driven approach allows pipelines to be more flexible and adaptable, reducing manual intervention and improving overall efficiency.

## 5.5. Cloud-Native Architectures (AWS, Azure, GCP)

The adoption of cloud-native architectures has become a cornerstone of scalable and resilient healthcare data pipelines. Cloud platforms such as Amazon Web Services, Microsoft Azure, and Google Cloud Platform provide a wide range of services for data storage, processing, and analytics.

Key advantages of cloud-native architectures include:

- Elastic scalability: Dynamically adjust resources based on workload
- Managed services: Reduce operational overhead with fully managed databases, streaming services, and analytics tools
- High availability: Built-in redundancy and disaster recovery capabilities
- Cost optimization: Pay-as-you-go pricing models

Cloud-native designs often leverage microservices, containerization, and serverless computing to create modular and flexible pipelines that can evolve with changing requirements.

### 5.6. Integration of FHIR APIs into Legacy EDI Pipelines

One of the most significant technical challenges in modern Medicaid systems is integrating FHIR-based APIs with existing legacy EDI pipelines. While EDI formats such as 837 transactions remain the standard for claims submission, interoperability mandates require the adoption of FHIR APIs for real-time data exchange.

This integration involves:

- Data translation layers: Converting EDI transactions into FHIR resources and vice versa
- API gateways: Managing access, authentication, and routing of API requests
- Hybrid architectures: Supporting both batch (EDI) and real-time (FHIR) processing within the same pipeline
- Backward compatibility: Ensuring that legacy systems continue to function while new capabilities are introduced

Successfully integrating these two paradigms enables organizations to modernize their data infrastructure without disrupting existing workflows. It also positions the pipeline to support future interoperability requirements and real-time analytics.

### 5.7. Concluding Perspective

Technical design considerations are central to building a resilient and future-proof Medicaid data pipeline. By addressing scalability, fault tolerance, data governance, and interoperability, organizations can create systems that not only meet current regulatory requirements but also adapt to emerging trends in healthcare data exchange.

A well-engineered pipeline is not just a technical asset it is a strategic enabler of efficient healthcare delivery, regulatory compliance, and data-driven decision-making.

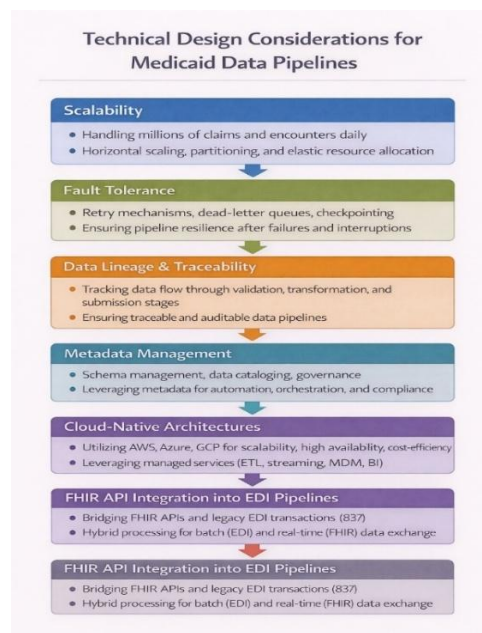


Figure 4. Technical Architecture and Design Considerations for Medicaid Data Pipelines

**Table 4. Technical Design Considerations for Medicaid Data Pipelines**

Design Consideration	Key Components	Technical Approach	Pipeline Impact
<b>Scalability</b>	High-volume claims/encounter processing	Distributed computing, horizontal scaling, data partitioning, elastic cloud resources	Enables processing of millions of records efficiently with minimal latency
<b>Fault Tolerance &amp; Retry Mechanisms</b>	Checkpointing, retries, dead-letter queues, redundancy	Automated retry policies, failure isolation, distributed system resilience	Prevents data loss, ensures continuous pipeline operation under failures
<b>Data Lineage &amp; Traceability</b>	Data provenance tracking, audit logs, transformation history	Metadata capture at each stage, lineage tracking tools, audit pipelines	Supports regulatory compliance, debugging, and transparency
<b>Metadata Management</b>	Schema management, data cataloging, governance policies	Centralized metadata repositories, automated schema versioning, data catalogs	Improves data discoverability, governance, and pipeline automation
<b>Cloud-Native Architecture</b>	AWS, Azure, GCP services, microservices, serverless computing	Use of managed services (ETL, storage, streaming), containerization, orchestration	Provides scalability, high availability, and cost-efficient infrastructure
<b>FHIR API Integration with EDI</b>	API gateways, data translation layers, hybrid ingestion	Mapping EDI (837) ↔ FHIR resources, API-first design, event-driven architecture	Enables real-time interoperability while maintaining legacy compatibility
<b>Security &amp; Compliance (Cross-Cutting)</b>	HIPAA compliance, encryption, access control	Data encryption (at rest/in transit), IAM policies, secure APIs	Ensures data privacy, regulatory compliance, and secure data exchange
<b>Performance Optimization</b>	Low-latency processing, efficient querying	Indexing, caching, parallel processing, optimized data models	Enhances pipeline speed and responsiveness for analytics and reporting

## 6. Data Quality, Governance, and Compliance

Data quality, governance, and compliance form the foundation of trust in Medicaid claims and encounter reporting systems. Without robust governance structures and rigorous validation mechanisms, even the most advanced data pipelines can produce unreliable or non-compliant outputs. Given the regulatory oversight imposed by agencies such as the Centers for Medicare & Medicaid Services, ensuring high data quality is not optional it is a critical requirement with direct financial and operational implications. This section explores how governance frameworks, validation strategies, and security controls collectively ensure that Medicaid data pipelines produce accurate, consistent, and compliant data outputs.

### 6.1. Data Governance Frameworks

A well-defined data governance framework establishes the policies, standards, and responsibilities required to manage data effectively across its lifecycle. In the Medicaid context, governance frameworks ensure alignment between state agencies, Managed Care Organizations (MCOs), providers, and federal regulators.

Key components of data governance include:

- Data ownership and stewardship: Clearly defining roles for data creation, validation, and maintenance
- Standardization policies: Enforcing consistent data definitions, formats, and coding systems across stakeholders
- Data lifecycle management: Governing how data are collected, stored, archived, and deleted
- Governance committees and oversight bodies: Providing accountability and decision-making authority

Effective governance frameworks enable organizations to reduce inconsistencies, improve data reliability, and ensure compliance with regulatory standards. They also support interoperability by aligning data practices across disparate systems.

## 6.2. Data Validation Rules and Audit Trails

Data validation is a critical mechanism for ensuring that Medicaid data meet required standards for accuracy, completeness, and consistency. Validation rules are typically implemented at multiple stages of the pipeline, including ingestion, transformation, and pre-submission.

Common validation rules include:

- Format validation: Ensuring data conform to required standards (e.g., X12, FHIR schemas)
- Business rule validation: Verifying logical consistency (e.g., valid diagnosis-procedure combinations)
- Referential integrity checks: Ensuring relationships between entities (e.g., patient-provider linkage)

In addition to validation, **audit trails** play a crucial role in compliance. Audit trails capture detailed logs of data transformations, validations, and access events, enabling:

- Traceability of data from source to submission
- Identification and resolution of errors
- Demonstration of compliance during regulatory audits

Together, validation rules and audit trails ensure that data pipelines produce trustworthy outputs while maintaining transparency and accountability.

## 6.3. Handling Inconsistencies Across Reporting Systems

One of the persistent challenges in Medicaid data management is the presence of inconsistencies across multiple reporting systems. These inconsistencies arise from differences in data formats, coding practices, system configurations, and reporting timelines among providers, MCOs, and state agencies.

Addressing these challenges requires a combination of technical and governance strategies:

- Standardized data models: Aligning data structures across systems to ensure consistency
- Data reconciliation processes: Comparing and resolving discrepancies between datasets
- Master Data Management (MDM): Maintaining a unified view of key entities such as patients and providers
- Automated anomaly detection: Identifying outliers and inconsistencies using rule-based or machine learning techniques

By systematically addressing inconsistencies, organizations can improve data quality and ensure that reporting outputs are reliable and compliant.

## 6.4. Privacy and Security

### 6.4.1. HIPAA Compliance

The Health Insurance Portability and Accountability Act (HIPAA) establishes national standards for protecting sensitive patient health information. Medicaid data pipelines must comply with HIPAA requirements to ensure the confidentiality, integrity, and availability of Protected Health Information (PHI).

Key HIPAA considerations include:

- Data encryption: Protecting data both at rest and in transit
- Secure data transmission: Using protocols such as HTTPS and secure file transfer mechanisms
- Audit controls: Monitoring access and usage of sensitive data
- Breach notification procedures: Ensuring timely reporting of security incidents

Failure to comply with HIPAA can result in significant legal and financial penalties, as well as reputational damage.

### 6.5. Role-Based Access Control (RBAC)

Role-Based Access Control (RBAC) is a fundamental security mechanism that restricts access to data based on user roles and responsibilities. In Medicaid data pipelines, RBAC ensures that only authorized personnel can access sensitive information.

Key aspects of RBAC include:

- Least privilege principle: Granting users only the access necessary to perform their **وظائف**
- Role definition and management: Assigning roles based on job functions (e.g., analyst, auditor, administrator)
- Access monitoring and auditing: Tracking user activity to detect unauthorized access

RBAC not only enhances security but also supports compliance with regulatory requirements by ensuring controlled access to sensitive data.

### 6.6. CMS Reporting Accuracy Requirements and Financial Implications

CMS imposes strict requirements on the accuracy and completeness of Medicaid data submissions. These requirements are critical for ensuring the integrity of program oversight, rate setting, and financial reconciliation.

Inaccurate or incomplete data can have significant financial implications, including:

- Incorrect capitation rate calculations: Leading to overpayment or underpayment to MCOs
- Penalties and corrective action plans: Imposed for non-compliance with reporting standards
- Reduced federal funding: **نتيجة** failure to meet data quality benchmarks
- Increased audit scrutiny: Resulting in additional operational costs and administrative burden

To mitigate these risks, organizations must implement robust validation, governance, and monitoring mechanisms throughout the pipeline.

### 6.7. Concluding Perspective

Data quality, governance, and compliance are not isolated components but integrated pillars of a successful Medicaid data pipeline. By establishing strong governance frameworks, enforcing rigorous validation rules, and ensuring robust security controls, organizations can build systems that are both reliable and compliant. Ultimately, high-quality data enable better decision-making, improved healthcare outcomes, and more efficient program management making governance and compliance essential investments rather than regulatory obligations.

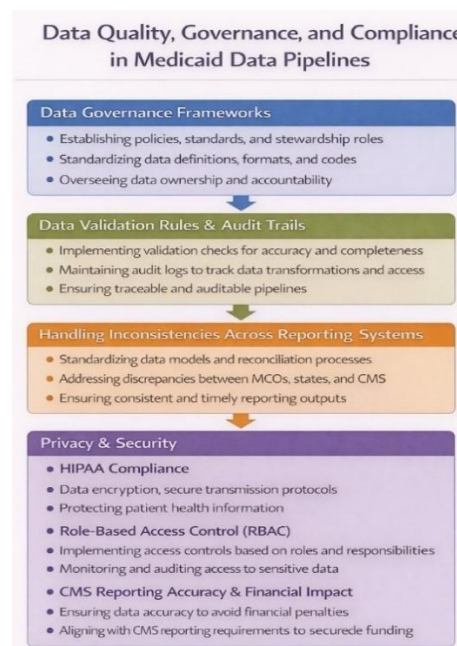


Figure 5. Data Governance, Quality, and Compliance Strategies for Medicaid Data Pipelines

**Table 5: Data Quality, Governance, and Compliance Framework for Medicaid Data Pipelines**

Domain	Key Components	Implementation Approach	Impact on Pipeline & Compliance
Data Governance Frameworks	Data ownership, stewardship, policies, standards	Establish governance committees, define roles, enforce standardized data definitions and lifecycle policies	Ensures consistency, accountability, and alignment across stakeholders (providers, MCOs, states)
Data Validation Rules	Format checks, business rules, referential integrity	Automated validation engines at ingestion and transformation stages, rule-based and ML-assisted validation	Improves data accuracy, reduces errors, ensures compliance with CMS quality standards
Audit Trails & Traceability	Data lineage, transformation logs, access logs	Implement logging frameworks, metadata tracking, audit pipelines aligned with CMS requirements	Enables traceability, supports audits, and ensures regulatory transparency
Handling Data Inconsistencies	Data reconciliation, standardization, MDM	Cross-system validation, schema alignment, entity resolution (patient/provider matching)	Reduces discrepancies across systems, improves data reliability and reporting accuracy
HIPAA Compliance	Data privacy, encryption, secure transmission	Encryption (at rest/in transit), secure APIs, access monitoring, breach detection systems	Protects patient data, ensures legal compliance, avoids penalties
Role-Based Access Control (RBAC)	User roles, access permissions, least privilege	Identity and Access Management (IAM), role definition, activity monitoring and auditing	Prevents unauthorized access, enhances data security and governance
CMS Reporting Accuracy Requirements	Accuracy, completeness, timeliness standards	Continuous data quality monitoring, validation thresholds, automated reporting checks	Ensures compliance with CMS mandates and reduces risk of penalties
Financial Implications of Data Quality	Rate setting accuracy, reimbursement integrity	Data quality scoring, reconciliation workflows, audit readiness mechanisms	Prevents over/underpayments, protects funding, improves financial outcomes
Compliance Monitoring & Governance Automation	KPIs, dashboards, alerts	Real-time monitoring tools, compliance dashboards, automated alerts for anomalies	Enables proactive compliance management and operational efficiency

## 7. Emerging Trends and Innovations

The Medicaid data ecosystem is undergoing a rapid transformation driven by advances in interoperability standards, cloud computing, and artificial intelligence. Traditional batch-oriented systems largely dependent on legacy EDI transactions are being re-engineered into real-time, intelligent, and interoperable data platforms. These innovations are not only improving operational efficiency but also enabling more proactive healthcare delivery, enhanced fraud detection, and data-driven policy decisions. This section explores the key emerging trends shaping the future of Medicaid data pipelines, highlighting how modern technologies are redefining data processing, validation, and exchange.

### 7.1. Transition from Batch EDI to Real-Time FHIR-Based Pipelines

Historically, Medicaid data pipelines have relied heavily on batch processing of EDI transactions such as 837 claims, which are submitted periodically and processed in scheduled intervals. While this approach has been effective for compliance and reimbursement, it introduces latency that limits real-time decision-making.

The emergence of FHIR-based APIs is driving a shift toward real-time or near real-time data pipelines. These pipelines enable continuous data exchange between providers, payers, and patients, supporting use cases such as care coordination, prior authorization, and patient access to health records.

From a technical perspective, this transition requires:

- Event-driven architectures capable of handling streaming data

- API gateways and microservices for real-time data access
- Hybrid models that integrate batch (EDI) and streaming (FHIR) workflows

This shift represents a fundamental change in how healthcare data are processed, moving from delayed reporting to immediate data availability and responsiveness.

## 7.2. AI-Driven Anomaly Detection in Claims Data

Artificial intelligence (AI) and machine learning are increasingly being integrated into Medicaid data pipelines to enhance data quality and detect anomalies. Claims and encounter datasets are highly complex and often contain hidden patterns that are difficult to identify باستخدام traditional rule-based methods.

AI-driven anomaly detection systems can:

- Identify unusual billing patterns that may indicate fraud or abuse
- Detect outliers in service utilization or cost trends
- Flag inconsistencies in coding or reporting across providers

These systems leverage techniques such as supervised learning, unsupervised clustering, and deep learning to analyze large-scale datasets in real time. By integrating AI into the pipeline, organizations can move from reactive error detection to proactive risk identification and mitigation.

## 7.3. Automated Data Validation and Reconciliation

Manual data validation processes are increasingly being replaced by automated validation and reconciliation frameworks. These systems use predefined rules, machine learning models, and workflow automation to ensure data accuracy and consistency across the pipeline.

Key capabilities include:

- Automated rule enforcement: Validating data against predefined business and regulatory rules
- Cross-system reconciliation: Comparing datasets from different sources to identify discrepancies
- Self-healing pipelines: Automatically correcting common errors and resubmitting data

Automation reduces operational overhead, improves data quality, and accelerates reporting timelines. It also enhances compliance by ensuring that validation processes are consistently applied across all data flows.

## 7.4. Interoperability-Driven Architectures (API Ecosystems)

Modern Medicaid systems are increasingly adopting interoperability-driven architectures, where APIs serve as the primary mechanism for data exchange. Instead of relying solely on file-based transfers, organizations are building API ecosystems that enable seamless communication between systems.

These architectures are characterized by:

- API-first design: Prioritizing standardized interfaces for data access and exchange
- Microservices-based architecture: Breaking down complex systems into modular, reusable components
- Ecosystem integration: Enabling data sharing across providers, payers, third-party applications, and patients

Interoperability-driven architectures not only improve data accessibility but also foster innovation by allowing external developers to build applications on top of healthcare data platforms.

## 7.5. Use of Data Lakes and Lakehouse Architectures

The growing volume and variety of Medicaid data have led to the widespread adoption of data lakes and lake house architectures. Traditional data warehouses, while effective for structured reporting, often struggle to handle unstructured and semi-structured data at scale.

- Data Lakes: Store raw data in its native format, providing flexibility for diverse data types and reprocessing needs

- Lakehouse Architectures: Combine the scalability of data lakes with the performance and structure of data warehouses

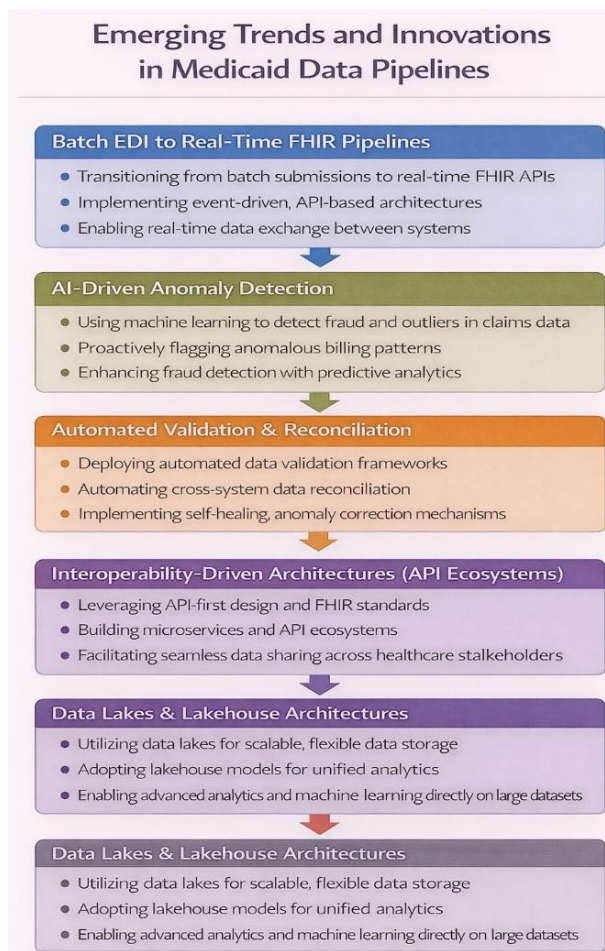
Lakehouse platforms enable:

- Unified storage for structured and unstructured data
- Advanced analytics and machine learning directly on large datasets
- Reduced data duplication and improved cost efficiency

These architectures support modern use cases such as predictive analytics, population health management, and longitudinal data analysis, making them a cornerstone of next-generation Medicaid data platforms.

### 7.6. Concluding Perspective

Emerging trends in Medicaid data pipelines reflect a broader shift toward real-time, intelligent, and interoperable healthcare systems. The integration of FHIR-based APIs, AI-driven analytics, and modern data architectures is transforming how data are collected, processed, and utilized. Organizations that embrace these innovations will be better positioned to meet regulatory requirements, improve operational efficiency, and deliver higher-quality healthcare services. Ultimately, these advancements are paving the way for a more connected, data-driven healthcare ecosystem.



**Figure 5. Emerging Technologies and Future Trends in Medicaid Data Pipelines**

**Table 5: Emerging Trends and Innovations in Medicaid Data Pipelines**

Trend / Innovation	Key Concepts	Technical Implementation	Impact on Medicaid Data Pipelines
Transition to Real-Time FHIR Pipelines	Shift from batch EDI (837) to real-time data exchange	FHIR APIs, event-driven architectures, streaming platforms (e.g., Kafka)	Enables near real-time reporting, faster decision-making, improved care coordination
AI-Driven Anomaly Detection	Fraud detection, outlier identification, predictive analytics	Machine learning models (supervised/unsupervised), anomaly detection algorithms	Enhances fraud detection, improves data quality, reduces financial losses
Automated Data Validation & Reconciliation	Rule-based validation, cross-system reconciliation, self-healing pipelines	Workflow automation, rule engines, ML-assisted validation systems	Reduces manual effort, improves accuracy, accelerates reporting cycles
Interoperability-Driven Architectures	API ecosystems, microservices, FHIR-based interoperability	API gateways, microservices architecture, standardized data exchange protocols	Facilitates seamless data sharing across providers, payers, and systems
Data Lakes Architecture	Raw data storage, schema-on-read flexibility	Cloud storage (e.g., S3, ADLS), distributed file systems	Supports large-scale data ingestion and flexible analytics
Lakehouse Architecture	Unified analytics combining data lakes and warehouses	Delta Lake, Apache Iceberg, unified query engines	Enables high-performance analytics, reduces data duplication, supports ML workloads
Hybrid Processing Models	Combination of batch (EDI) and real-time (FHIR) pipelines	Integration layers, data translation (EDI ↔ FHIR), orchestration tools	Ensures backward compatibility while enabling modernization
Advanced Analytics & AI Integration	Predictive modeling, population health insights, risk scoring	AI/ML platforms, big data analytics tools, cloud AI services	Drives data-driven decision-making and proactive healthcare management

### 8. Case Study / Conceptual Framework

A case study or conceptual framework is valuable in research on Medicaid claims and encounter reporting because it translates regulatory expectations and technical principles into a practical, understandable system design. While previous sections discuss standards, interoperability rules, governance, and pipeline layers in detail, this section brings those elements together into a single operational view. It demonstrates how a Medicaid reporting pipeline can be organized end to end, beginning with source data capture and concluding with compliant submission and decision-support reporting.

In the context of Medicaid, a conceptual framework is especially important because the reporting environment is not shaped by technology alone. It is influenced by CMS reporting requirements, managed care obligations, interoperability mandates, data quality expectations, and financial accountability. As a result, a reference architecture must do more than move data from one system to another. It must support validation, traceability, compliance monitoring, resubmission workflows, and analytics, all within a scalable and auditable structure.

This section presents an example Medicaid pipeline architecture using the sequence **Ingestion** → **Validation** → **Transformation** → **Storage** → **Submission** → **Reporting**. It also maps major CMS requirements to technical components, showing how policy and architecture must align to produce a compliant and high-performing reporting ecosystem.

### 8.1. Example Architecture for a Medicaid Pipeline

A Medicaid claims and encounter pipeline can be conceptualized as a sequence of tightly connected layers, each responsible for a distinct set of functions. Although actual implementation may vary by state Medicaid agency, managed care organization, or technical platform, the overall design logic remains consistent.

At a high level, the architecture begins with upstream source systems such as provider EHRs, billing platforms, pharmacy systems, and managed care feeds. These sources generate raw transactional and clinical data in different formats, including EDI transactions, flat files, proprietary extracts, and increasingly, FHIR-based API payloads. Because these systems often differ in structure, timing, and data quality, the pipeline must be capable of absorbing heterogeneous inputs while applying consistent downstream controls.

The proposed architecture follows six major stages:

#### 8.1.1. Ingestion

The ingestion stage is responsible for collecting data from all participating source systems. This includes batch ingestion of 837 institutional and professional claims, pharmacy transactions, and managed care encounter files, as well as streaming or API-based ingestion from modern interoperability platforms. The ingestion layer should support secure file transfer, API gateways, message queues, and scheduling tools. This stage performs the first level of technical control. Files and payloads are checked for structural validity, schema compliance, completeness of required fields, and source authentication. In practice, the ingestion layer acts as the pipeline's intake gate, ensuring that corrupted, incomplete, or unauthorized data do not flow unchecked into downstream processing.

#### 8.1.2. Validation

After ingestion, data move into a validation layer where business rules, compliance checks, and data quality rules are applied. This stage determines whether the incoming data meet required standards for completeness, timeliness, accuracy, and logical consistency. Examples include verifying valid diagnosis and procedure code combinations, checking beneficiary identifiers, ensuring provider information is valid, and confirming that service dates fall within acceptable reporting timelines. Validation is one of the most critical stages in Medicaid reporting because errors at this point have direct implications for federal compliance, financial reconciliation, and audit outcomes. The validation layer should also generate exception logs, quality scores, and rule-failure reports so that rejected or questionable records can be reviewed, corrected, and resubmitted.

#### 8.1.3. Transformation

Once data pass validation, they move into transformation and standardization. In this stage, heterogeneous source formats are converted into a unified canonical model suitable for storage, analytics, and submission. Code mapping occurs here, such as aligning ICD diagnosis codes, CPT or HCPCS procedure codes, NDC drug codes, and other source-specific elements to internal enterprise data models or analytical vocabularies such as RxNorm and ATC for medication analysis. Transformation also includes normalization of field names, date formats, provider identifiers, member identifiers, and service categories. If the pipeline supports both legacy EDI workflows and modern FHIR-based interoperability, this stage may also include translation logic between EDI structures and FHIR resources. This makes the transformation layer the architectural bridge between fragmented source systems and a unified Medicaid reporting model.

#### 8.1.4. Storage

After transformation, the data are stored in layered repositories designed for different use cases. Raw source files and payloads are retained in a data lake for auditability and reprocessing. Standardized and curated records are loaded into structured repositories such as enterprise warehouses or lake house environments for reporting, reconciliation, and regulatory analysis. This stage may also include Master Data Management functions, which create consistent identities for members, providers, organizations, and plans across systems. Without this layer, the same entity may appear multiple times under slightly different identifiers or formats, reducing data reliability and making cross-system reporting difficult. The storage layer must therefore support both historical retention and efficient retrieval, especially in environments where longitudinal access to multiple years of data is required.

### 8.1.5. Submission

The submission layer converts validated and curated records into CMS-compliant reporting outputs. This may include preparing encounter data for submission to state systems, encounter data systems, or T-MSIS-aligned reporting environments. It may also include packaging output into required EDI structures or state-specific templates. A mature submission layer includes automated pre-submission checks, acknowledgment tracking, rejection management, and resubmission workflows. If a record is rejected by a downstream system, the pipeline should preserve the error context, notify the responsible team, and support a closed-loop correction process. This turns submission from a one-time outbound transfer into a governed operational workflow.

### 8.1.6. Reporting

The final stage of the framework is reporting and analytics. Here, the processed data are used for regulatory dashboards, operational monitoring, financial reconciliation, utilization analysis, managed care oversight, and quality performance review. This layer serves both compliance and business intelligence functions. Regulators and compliance teams may need reports on completeness, timeliness, submission acceptance rates, and audit readiness. Program administrators may need insights into claims volume, utilization patterns, provider behavior, and encounter lag. Finance teams may require reconciliation and risk adjustment support. Thus, the reporting layer converts a compliance pipeline into a strategic decision-support system.

## 8.2. Conceptual Flow of the Proposed Framework

The proposed framework can be summarized as a controlled flow of data and decisions:

**Source Systems → Ingestion → Validation → Transformation → Storage → Submission → Reporting**

This sequence is not merely linear. In practice, it contains multiple feedback loops. For example, validation failures may trigger correction workflows back to source owners. Submission rejections may result in targeted reprocessing from the storage layer. Reporting anomalies may reveal upstream data quality issues that require updates to validation rules. This feedback-oriented design is essential in Medicaid environments because reporting is rarely a one-pass process. Instead, it is iterative, compliance-sensitive, and operationally dependent on continuous quality improvement. A strong conceptual framework therefore includes both forward data movement and backward exception management.

## 8.3. Mapping CMS Requirements to Technical Components

One of the strongest ways to present this case study is to explicitly show how CMS expectations are translated into technical architecture. This helps demonstrate that the proposed design is not just technically sound, but also policy-aligned. Below is a structured mapping of major CMS-oriented requirements to corresponding technical components.

**Table 6: Mapping CMS Compliance Requirements to Technical Components in Medicaid Data Pipelines**

CMS / Compliance Requirement	Technical Component	Purpose in the Pipeline
Accuracy of encounter and claims data	Validation engine with business rules and referential checks	Ensures records conform to coding, logic, and reporting requirements
Completeness of required data elements	Ingestion checks and completeness scoring modules	Detects missing fields and prevents incomplete records from moving forward
Timely reporting and submission cycles	Workflow scheduler, SLA monitoring, alerting tools	Helps ensure records are processed and submitted within required deadlines
Auditability and traceability	Data lineage framework, audit logs, metadata repository	Supports CMS reviews, internal audits, and root-cause analysis
Standardized electronic submission	EDI generation services, schema validators, interface engines	Formats output into compliant submission structures

Managed care encounter oversight	MCO feed integration and reconciliation layer	Aligns managed care submissions with state and CMS expectations
Longitudinal data access and retention	Data lake, warehouse, and archival storage strategy	Supports historical reporting and multi-year analysis
Interoperability and API-driven exchange	FHIR API gateway, transformation layer, event streaming tools	Enables integration of modern interoperability workflows into the pipeline
Error correction and resubmission	Exception management queue and resubmission workflow engine	Supports rejected-record remediation and controlled reprocessing
Reporting quality and operational oversight	Dashboards, analytics mart, compliance scorecards	Provides visibility into submission performance and data quality outcomes

This mapping is useful because it shows that each regulatory expectation should correspond to a deliberate technical control. In other words, compliance should not be treated as a final checkpoint; it should be built directly into the architecture.

#### 8.4. Illustrative Conceptual Case Narrative

To make this framework even more concrete, it may help to describe it as a simplified case example. Consider a state Medicaid environment receiving data from multiple provider groups, pharmacy benefit systems, and managed care organizations. Each source sends data in different formats and on different schedules. Some providers submit institutional claims in 837I format, outpatient claims in 837P, pharmacy systems generate NCPDP transactions, and newer platforms expose portions of the data through FHIR APIs.

In the proposed architecture, all these data feeds enter a central ingestion layer. The system first verifies source identity, file integrity, schema compliance, and baseline field completeness. It then routes valid records into a rules engine where encounter completeness, diagnosis validity, provider-member relationships, and timeliness constraints are checked. Records that fail are diverted into an exception queue and assigned for remediation.

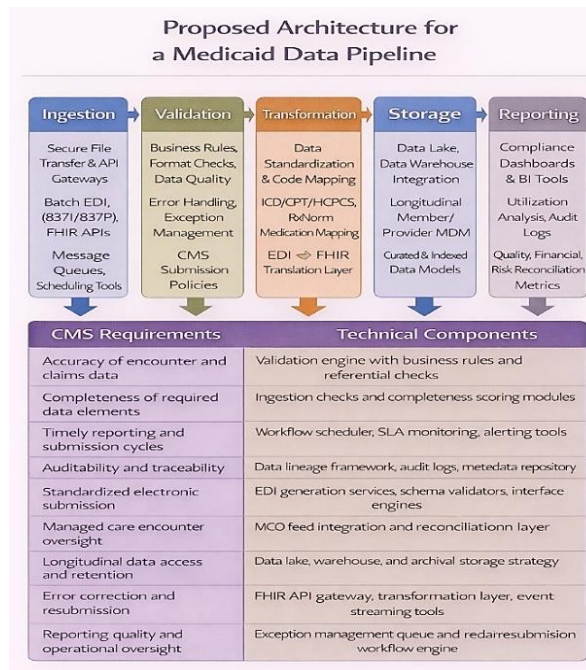
Records that pass are normalized into a canonical Medicaid data model. Drug-related data are classified using standardized medication vocabularies, clinical codes are harmonized, and identifiers are aligned with enterprise master data. These curated records are loaded into the raw and structured storage layers, where they are preserved for both compliance and analytical reuse. The submission engine then creates state- or CMS-compatible reporting outputs and transmits them through controlled workflows. Downstream acknowledgments are captured automatically. If any records are rejected, the system routes them into correction and resubmission workflows without losing lineage or audit context. Finally, dashboards provide visibility into acceptance rates, lag days, completeness scores, financial reconciliation status, and operational performance. This case narrative illustrates how the conceptual framework operates not just as a technical design, but as a compliance-centered reporting ecosystem.

#### 8.5. Why This Framework Matters

The value of this conceptual framework lies in its ability to connect fragmented policy, technical, and operational requirements into one integrated system design. Medicaid reporting environments are often challenged by disconnected source systems, manual quality checks, delayed submissions, and limited visibility into errors. A structured reference architecture addresses these weaknesses by introducing repeatable controls, transparent workflows, and scalable data management practices. This framework also reflects the ongoing modernization of healthcare data systems. It supports traditional EDI-based reporting requirements while leaving room for real-time interoperability through FHIR APIs, event-driven ingestion, and modern analytics environments. That makes it especially relevant in a period where Medicaid agencies and managed care organizations must operate across both legacy and emerging technologies.

#### 8.6. Concluding Perspective

The case study and conceptual framework presented here demonstrate how a Medicaid claims and encounter pipeline can be designed as an integrated, policy-aligned architecture. By organizing data movement into the stages of Ingestion, Validation, Transformation, Storage, Submission, and Reporting, the framework creates a clear operational model for handling complexity, improving data quality, and supporting compliance. Just as importantly, mapping CMS requirements to technical components shows that effective Medicaid pipeline design is not only about engineering efficiency. It is also about building a system that is auditable, traceable, interoperable, and resilient enough to support reimbursement integrity, program oversight, and regulatory trust.



**Figure 5. Proposed End-to-End Medicaid Data Pipeline Architecture for CMS Compliance and Reporting**

### 9. Discussion

The design of Medicaid claims and encounter data pipelines involves navigating a complex set of trade-offs between performance, compliance, interoperability, and data quality. As healthcare systems evolve toward real-time data exchange and advanced analytics, organizations must carefully evaluate architectural decisions that balance legacy constraints with modern innovation. This discussion synthesizes key insights from the proposed pipeline architecture and highlights critical considerations related to batch vs. real-time processing, interoperability regulations, and the role of data quality in financial and clinical outcomes.

#### 9.1. Trade-offs Between Batch vs. Real-Time Pipelines

One of the most significant architectural decisions in Medicaid data systems is the choice between batch processing and real-time (streaming) pipelines. Traditional Medicaid workflows have been built around batch processing using EDI transactions such as 837 claims. Batch pipelines are highly reliable, well-understood, and aligned with existing regulatory submission cycles. They are particularly effective for large-scale data aggregation, financial reconciliation, and compliance reporting. However, batch processing introduces inherent latency, as data are processed at scheduled intervals rather than continuously. This delay limits the ability to perform real-time analytics, detect anomalies early, or support time-sensitive decision-making.

In contrast, real-time pipelines enabled by FHIR APIs and event-driven architectures allow continuous data ingestion and processing. This supports use cases such as immediate eligibility verification, real-time prior authorization, and proactive fraud detection. Real-time systems enhance responsiveness and enable more dynamic healthcare delivery models. However, real-time pipelines introduce additional complexity. They require advanced infrastructure, including streaming platforms, API gateways, and low-latency processing frameworks. They also demand more sophisticated error handling, monitoring, and synchronization mechanisms. Furthermore, regulatory systems such as CMS reporting frameworks still rely heavily on batch-oriented submissions, meaning that real-time capabilities must often coexist with legacy processes.

As a result, many organizations adopt a hybrid architecture, combining batch pipelines for compliance and reporting with real-time pipelines for analytics and interoperability. This approach balances reliability with innovation but requires careful coordination between different processing models.

## 9.2. Impact of Interoperability Regulations on System Design

Interoperability regulations particularly those driven by the Centers for Medicare & Medicaid Services have fundamentally reshaped Medicaid data pipeline design. Requirements such as FHIR-based APIs, patient access mandates, and payer-to-payer data exchange have shifted the focus from isolated data systems to connected, API-driven ecosystems.

These regulations introduce several key implications for system design:

- **API-First Architecture:** Systems must be designed around standardized interfaces, enabling seamless data exchange across stakeholders. This requires the integration of API gateways, authentication mechanisms, and service orchestration layers.
- **Data Standardization:** Interoperability depends on consistent data representation. Pipelines must support standardized vocabularies and schemas, such as FHIR resources, to ensure compatibility across systems.
- **Real-Time Data Access:** Regulations increasingly emphasize timely data availability, pushing pipelines toward streaming and event-driven architectures.
- **Security and Privacy Controls:** As data sharing increases, so does the need for robust security measures, including encryption, access control, and auditability.

While these regulations promote transparency and improve care coordination, they also increase the complexity of system design. Organizations must invest in new technologies and skills to meet interoperability requirements while maintaining compatibility with legacy systems.

## 9.3. Importance of Data Quality in Financial and Clinical Outcomes

Data quality is a critical determinant of both financial performance and clinical effectiveness in Medicaid programs. Inaccurate, incomplete, or delayed data can have far-reaching consequences, affecting reimbursement, compliance, and patient care.

From a financial perspective, data quality directly influences:

- Capitation rate calculations and risk adjustment models
- Claims adjudication and reimbursement accuracy
- Detection of fraud, waste, and abuse

Errors in data can lead to overpayments or underpayments, financial penalties, and increased audit scrutiny. Poor data quality also undermines trust in reporting systems and can result in reduced federal funding or corrective action requirements.

From a clinical perspective, high-quality data enable:

- Accurate patient records and care coordination
- Population health analytics and outcome measurement
- Evidence-based decision-making

Incomplete or inconsistent data can lead to gaps in care, misinformed clinical decisions, and reduced quality of service delivery. Given these impacts, data quality must be treated as a core design principle rather than a downstream concern. This requires embedding validation, governance, and monitoring mechanisms throughout the pipeline, ensuring that data remain accurate, complete, and timely at every stage.

## 9.4. Concluding Perspective

The discussion highlights that Medicaid data pipeline design is not simply a technical challenge but a strategic balancing act. Organizations must navigate trade-offs between batch reliability and real-time responsiveness, align with evolving interoperability regulations, and prioritize data quality as a driver of both financial integrity and clinical outcomes. Ultimately, the most effective pipeline architectures are those that integrate these considerations into a cohesive system leveraging hybrid processing models, API-driven interoperability, and robust data governance to deliver compliant, scalable, and high-quality data solutions.



**Figure 6. Key Design Considerations and Trade-Offs in Medicaid Data Pipelines**

**Table 7. Key Discussion Insights for Medicaid Data Pipeline Design**

Discussion Area	Key Considerations	Advantages	Challenges / Trade-offs	Implications for System Design
Batch Pipelines (EDI-Based)	Scheduled processing of claims (837), compliance-driven workflows	High reliability, mature standards, easier regulatory alignment	High latency, delayed insights, limited real-time responsiveness	Suitable for CMS reporting, requires optimization for large-scale batch processing
Real-Time Pipelines (FHIR-Based)	Event-driven, API-based data exchange	Immediate data availability, supports real-time analytics and care coordination	Increased complexity, higher infrastructure cost, integration challenges	Requires streaming architecture, API gateways, and low-latency systems
Hybrid Architecture (Batch + Real-Time)	Combination of batch reporting and real-time ingestion	Balances compliance and innovation, supports legacy and modern systems	Increased architectural complexity, synchronization challenges	Recommended approach for Medicaid systems transitioning to interoperability
Impact of Interoperability Regulations	FHIR mandates, API access, payer-to-payer exchange	Improved data sharing, enhanced care coordination, patient empowerment	Requires system redesign, higher security	Drives API-first, microservices-based, and interoperable architectures

			requirements, skill gaps	
Data Standardization Requirements	Use of standardized vocabularies (FHIR, ICD, CPT)	Ensures consistent data exchange and interoperability	Mapping complexity across systems, versioning issues	Requires strong transformation and schema alignment layers
Security & Privacy Requirements	Compliance with HIPAA and CMS policies	Protects patient data, ensures regulatory compliance	Increased implementation complexity and operational overhead	Requires encryption, RBAC, audit logs, and secure APIs
Data Quality in Financial Outcomes	Accurate claims, risk adjustment, reimbursement	Prevents over/underpayments, supports financial integrity	Poor data leads to penalties, audit risks, funding loss	Requires continuous validation, reconciliation, and monitoring
Data Quality in Clinical Outcomes	Accurate patient records, care coordination, analytics	Improves patient care, supports population health insights	Incomplete data leads to poor decision-making and care gaps	Requires high-quality, standardized, and timely data pipelines
Operational Complexity vs. Innovation	Balancing legacy systems with modern technologies	Enables gradual modernization without disruption	Integration complexity, higher maintenance effort	Requires modular, scalable, and flexible pipeline architecture

## 10. Conclusion

### 10.1. Summary of Key Findings

This study has presented a comprehensive analysis of end-to-end data pipeline design for Medicaid claims and encounter reporting, emphasizing the integration of technical architecture, regulatory compliance, and interoperability standards. The findings highlight that Medicaid data ecosystems are inherently complex, requiring the coordination of multiple stakeholders, diverse data sources, and stringent reporting requirements.

A layered pipeline architecture spanning ingestion, validation, transformation, storage, submission, and reporting emerges as a critical framework for managing this complexity. The study demonstrates that embedding validation, governance, and audit mechanisms throughout the pipeline significantly improves data quality, reduces reporting errors, and enhances compliance with federal requirements. Additionally, the analysis underscores the growing importance of hybrid processing models that combine traditional batch-based EDI workflows with real-time, API-driven data exchange. This hybrid approach enables organizations to maintain compatibility with legacy systems while adopting modern interoperability capabilities.

### 10.2. Importance of Robust Pipeline Design for Medicaid Programs

Robust pipeline design is not merely a technical necessity but a strategic requirement for Medicaid programs. A well-designed pipeline ensures that data are accurate, complete, and timely qualities that directly impact reimbursement accuracy, program integrity, and policy effectiveness. From a financial perspective, reliable pipelines support accurate claims adjudication, risk adjustment, and capitation rate calculations, reducing the risk of overpayments, underpayments, and regulatory penalties. From a clinical perspective, high-quality data enable better care coordination, improved population health management, and more informed decision-making.

Furthermore, robust pipelines enhance operational efficiency by automating data processing, reducing manual intervention, and enabling scalable data management. They also provide the transparency and traceability required for audits and regulatory oversight, strengthening trust between stakeholders and governing bodies such as the Centers for Medicare & Medicaid Services.

### 10.3. Future Outlook: Fully Interoperable, API-Driven Healthcare Ecosystems

Looking ahead, Medicaid data systems are expected to evolve toward fully interoperable, API-driven healthcare ecosystems. The increasing adoption of FHIR-based APIs, real-time data exchange, and interoperability mandates will continue to reshape how healthcare data are shared and utilized.

Future pipelines will likely be characterized by:

- Real-time, event-driven architectures that enable immediate data availability and proactive decision-making
- API-first ecosystems that facilitate seamless integration across providers, payers, and third-party applications
- Advanced analytics and AI integration, supporting predictive modeling, anomaly detection, and personalized healthcare insights
- Unified data platforms, such as lake house architectures, that combine scalability with high-performance analytics

As these trends mature, the distinction between data pipelines and healthcare platforms will increasingly blur, resulting in integrated systems that support both operational workflows and strategic analytics.

### 10.4. Concluding Perspective

In conclusion, the design of end-to-end data pipelines for Medicaid claims and encounter reporting is a foundational element of modern healthcare systems. By aligning technical architecture with regulatory requirements and emerging interoperability standards, organizations can build pipelines that are not only compliant but also scalable, efficient, and future-ready. Such pipelines will play a pivotal role in enabling data-driven healthcare delivery, improving program oversight, and advancing the broader goal of a connected and interoperable healthcare ecosystem.

## References

- [1] Centers for Medicare & Medicaid Services. (2023). *Medicaid managed care encounter data toolkit*. Retrieved from <https://www.medicaid.gov>
- [2] Centers for Medicare & Medicaid Services. (2023). *T-MSIS data quality and reporting guidance*. CMS.
- [3] Vorisek, C. N., et al. (2022). Fast Healthcare Interoperability Resources (FHIR) for healthcare research: Systematic review. *JMIR Medical Informatics*, 10(7), e35724. <https://doi.org/10.2196/35724>
- [4] Mandl, K. D., et al. (2020). SMART/HL7 FHIR bulk data access for population health. *npj Digital Medicine*, 3(1), 1–6.
- [5] Hong, N., et al. (2019). Developing a scalable FHIR-based clinical data normalization pipeline. *JAMIA Open*, 2(4), 570–579.
- [6] Li, Y., Wang, H., Yerebakan, H., Shinagawa, Y., & Luo, Y. (2023). Enhancing health data interoperability with large language models: A FHIR study. *arXiv preprint arXiv:2310.12989*.
- [7] U.S. Department of Health and Human Services. (2023). *Building data capacity for patient-centered outcomes research*. National Academies Press.
- [8] Adams, W. O., & Krukowski, M. T. (2019). Patient access to health data: The impact of FHIR on patient engagement. *Health Affairs*, 38(5), 773–779.
- [9] Cho, E., & Lee, Y. K. (2019). A comparative study of FHIR and HL7 v2 standards. *International Journal of Information Management*, 49, 412–420.
- [10] Chan, K. W. (2020). The impact of FHIR on personalized healthcare and health management. *BMC Health Services Research*, 20(1), 101–110.
- [11] Saini, V., Reddy, S. G., Kumar, D., & Ahmad, T. (2021). Evaluating FHIR's impact on healthcare data interoperability. *International Journal of Health Informatics*.
- [12] Chauhan, P., et al. (2023). Interoperable synthetic health data using FHIR for clinical decision support systems. *arXiv preprint arXiv:2308.02613*.
- [13] Health Level Seven International. (2023). *FHIR Release 4 (R4) specification*. HL7.
- [14] Office of the National Coordinator for Health Information Technology. (2023). *US Core Data for Interoperability (USCDI) Version 4*. ONC.
- [15] Kahn, M. G., et al. (2016). A harmonized data quality assessment framework for electronic health data. *eGEMS*, 4(1), 18–25.
- [16] Wang, Y., Kung, L., & Byrd, T. A. (2018). Big data analytics: Understanding its capabilities and potential benefits. *Technological Forecasting and Social Change*, 126, 3–13.
- [17] Health Insurance Portability and Accountability Act. (1996). *Public Law 104–191*. U.S. Government Publishing Office.