

# **BUAN 6335.502 Organizing for Business Analytics Platforms**

## **Modern Data Platform for MeddGGenius: A Cloud-Based Data Lakehouse Solution**

**Professor – Mandar Samant, TA – Janani Dhanasekaran**

### **Executive Summary**

This report outlines a conceptual design for a modern data platform at MeddGGenius, a regional medical center facing increasing competition and struggling with outdated legacy systems. The current infrastructure, based on distributed RDBMS and a Teradata warehouse, is inadequate for handling the growing volume and variety of healthcare data, particularly unstructured data like medical images and sensor data.

To address these challenges, MeddGGenius has appointed a new Chief Data Officer, Ms. Angel Datsmart, to lead the digital transformation. The proposed solution is a cloud-based data lakehouse architecture that combines the flexibility of a data lake with the structure of a data warehouse. This unified platform will:

1. Consolidate all types of data into a single, accessible repository
2. Support both OLTP for critical operations and OLAP for advanced analytics
3. Enable seamless data integration and real-time processing
4. Ensure HIPAA compliance and meet NIST cybersecurity standards
5. Facilitate predictive analytics and machine learning capabilities

By implementing this modern data platform, MeddGGenius aims to:

- Enhance patient care through improved decision-making
- Streamline operations and reduce costs
- Eliminate data silos and establish a single source of truth
- Enable proactive identification of patient needs
- Maintain a competitive edge in an evolving healthcare landscape

This strategic initiative will empower MeddGGenius to leverage its data assets effectively, drive innovation in patient care, and solidify its position as a leading regional medical center.

### **1. Introduction and Problem Statement**

MeddGGenius is the leading regional medical center that nowadays suffers from a very outdated data infrastructure. Currently, this organization relies on distributed RDBMS systems and the

Teradata warehouse, both of which can barely cope with the dynamic character of contemporary healthcare data.

### **Limitations of the Legacy System:**

These systems are poorly equipped to deal with the ever-increasing volume and variety of healthcare data, especially unstructured data, which includes medical images, clinical notes, and real-time sensor data from IoT devices. This becomes a limitation for MeddGGenius in leveraging rich insights from diverse data sources.

### **Data Silos:**

This has promoted data silos at every level in a department and system, making the unification of a patient's information hard to retrieve; this further hinders efficient analysis and decision-making processes.

### **Challenges**

The infrastructure of MeddGGenius is so aged that operational efficiency, innovation, and productivity suffer a lot. Specifically, the organization has to address some major issues:

1. **Maintenance Costs:** The operation and maintenance of legacy systems are extremely expensive. Hardware parts, software updates, and skilled labor all contribute to high maintenance costs, diverting the budget from possible investment in modernization or the decommissioning of obsolete systems.

2. **Integration Challenges:** Most data management platforms, such as Electronic Health Records and imaging systems, work in silos. This is the biggest barrier to sharing data smoothly and thus hinders the development of a complete picture of patient data, which is highly essential for informed decisions.

3. **Performance Issues:** Aging systems are generally slow, experience frequent downtimes, and have limited scalability. These shortcomings affect staff productivity and patient care, causing delays in retrieving critical records or processing diagnostic data.

4. **Project Implementation Delays:** Inflexible infrastructure slows down the adoption of new technologies and processes. Issues of compatibility and inefficiencies in the system slow down the integration of advanced analytics tools or even updated EHR modules to a speed that keeps pace with industry advancements.

5. **Lack of Real-Time Processing:** Timely processing of data is crucial in healthcare. Existing systems fail to process real-time data streams from IoT devices or analytics from the emergency

room, which can put patients at risk and also hinder the adaptation of advanced solutions such as predictive analytics.

These challenges collectively form a significant barrier to operational excellence and patient-centric innovation at MeddGGenius. Addressing these issues requires implementing a modern, scalable data platform to enhance care delivery and maintain competitiveness in the evolving healthcare landscape.

### **Goals and Objectives:**

To address these challenges and maintain its competitive edge, MeddGGenius aims to:

1. Create a unified data platform capable of handling both OLTP and OLAP workloads
2. Develop a system that can efficiently manage structured and unstructured data
3. Improve patient care through enhanced data-driven decision-making.
4. Comply with the set regulations of HIPAA and NIST standards on cybersecurity

The organization has just hired a Chief Data Officer, Ms. Angel Datasmart, to drive this digital transformation. The CDO's vision is to develop a modern data platform that unifies all types of data into one place and is accessible by staff, patients, and data scientists alike.

This transformation is very significant to ensure that MeddGGenius overcomes present limitations, improves operational efficiencies, and innovates for patient care. The advanced analytics on the new platform, with support for predictive modeling and machine learning use cases, need to assure data security and compliance within the highly regulated healthcare sector. In order to meet these goals, MeddGGenius will have to consider a cloud-based solution that is scalable, flexible, and provides advanced data management capabilities. This approach will enable the organization to break down data silos, implement real-time processing, and leverage advanced analytics to improve patient outcomes and operational efficiency.

## **2. Proposed Solution**

AWS for a Cloud-Based Data Lakehouse Architecture to Power MeddGGenius will provide a modern, scalable, and compliant data platform. This caters to the unified needs of the organization in structured, unstructured data handling with support for advanced analytics and assurance of regulatory compliance.

### **Solution High-Level Architecture**

#### **1. Data Ingestion Layer:**

AWS Kinesis: This handles streaming in real time from IoT devices, medical equipment, and time-sensitive sources.

AWS S3: Serves as the central data lake, storing raw data in its native

## **2. Storage Layer:**

AWS S3: Provides the centralized data lake for storing raw data in its native format. S3 storage classes, such as Standard, Intelligent-Tiering, and Glacier, optimize costs based on data access patterns.

## **3. Data Processing and Transformation Layer:**

AWS Glue: Performs ETL jobs, data cataloging, and schema discovery.

## **4. Data Warehouse and Analytics Layer:**

AWS Athena: Enables SQL-based querying of data stored in S3.

## **5. Machine Learning and Advanced Analytics:**

AWS SageMaker: Supports the development, training, and deployment of machine learning models. Amazon Comprehend Medical: It processes unstructured medical text data.

## **6. Data Governance and Security:**

AWS Lake Formation: It provides a data governance approach to manage fine-grained access control. AWS IAM: It provides identity and access management across the platform for users.

AWS KMS: Ensures encryption of data at rest and in transit.

## **7. Visualization and Reporting:**

Amazon QuickSight: Build interactive dashboards and reports for different stakeholders.

## **8. Monitoring and Compliance:**

AWS CloudTrail: Logs API usage and resource changes for auditing.

Amazon CloudWatch: Monitor the entire infrastructure and provide observability.

## **Data Flow:**

1. Data is ingested from various sources using Kinesis (real-time) or batch processes.
2. Raw data is stored in S3, organized into different buckets based on data type and sensitivity.
3. AWS Glue catalogs the data and performs necessary ETL jobs.
4. Processed data is stored back in S3 or loaded into Sagemaker for structured analytics.
5. SageMaker accesses the processed data for machine learning tasks.

6. QuickSight will generate visualizations and reports based on the analyzed data.

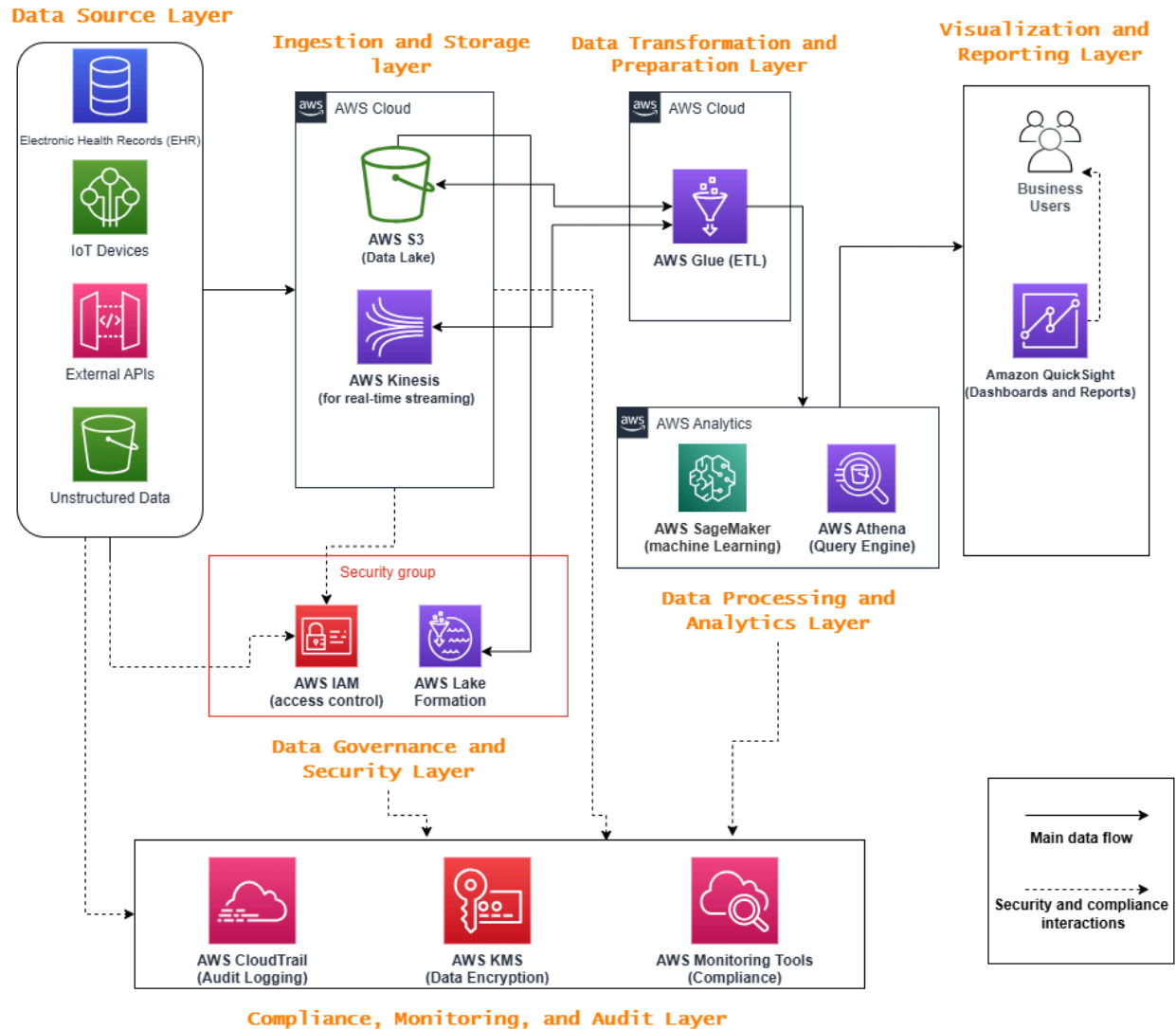
7. Lake Formation and IAM ensure proper access controls throughout the data lifecycle.

This architecture provides MeddGGenius with a scalable, secure, and compliant data platform for a variety of healthcare data types, thus enabling advanced analytics for better patient care and operational efficiency.

## Architecture Diagram

### MeddGGenius Cloud-Based Data Lakehouse Architecture

Proposed data platform leveraging AWS services for unified healthcare data management



### Detailed Breakdown of the Proposed Architecture for MeddGGenius

The architecture for the cloud-based data lakehouse solution on AWS can be broken into seven main layers. Here's a step-by-step breakdown, including the flow of data and the role of each AWS service:

## 1. Data Sources Layer

**Purpose:** Collect data from multiple sources in structured and unstructured formats.

**Components:**

- **Electronic Health Records (EHR):** Structured patient information like demographics, diagnoses, and treatments.
- **IoT Devices:** Continuous streaming data from patient monitoring devices (e.g., heart rate, oxygen levels).
- **External APIs:** Data from external systems, such as public health APIs or third-party healthcare apps.
- **Unstructured Data:** Medical images (X-rays, MRIs), clinical notes, and sensor data.

**Flow:** Data from these sources flows into the central data lake via **AWS Kinesis** (real-time data) or **AWS S3** (batch processing).

To expand,

### Diverse Data Sources in Healthcare

Modern healthcare generates a vast array of data, and we should aim to capture a wide range to maximize the platform's potential. Here are some key sources that we can think of:

1. **Electronic Health Records (EHRs):** These remain a cornerstone of patient data, containing structured information like demographics, medical history, diagnoses, medications, allergies, and lab results.
2. **IoT Devices:** The rise of wearable sensors and smart medical devices provides real-time physiological data (heart rate, blood pressure, glucose levels), activity data, and treatment adherence information. Examples include:
  - **Wearable fitness trackers:** Capture patient activity levels and sleep patterns.
  - **Smart insulin pens:** Record dosage and timing of insulin administration.
  - **Remote patient monitoring (RPM) devices:** Continuously track vital signs for early detection of health deteriorations.
3. **Medical Imaging:** Radiology images (X-rays, MRIs, CT scans) are essential for diagnosis and treatment planning. This also includes pathology slides and cardiology images.
4. **Clinical Notes:** Physician notes, discharge summaries, and other textual documentation provide valuable context not always captured in structured fields.

5. **Lab Information Systems (LIS):** These systems house lab results, including blood tests, microbiology cultures, and pathology reports.
6. **Pharmacy Systems:** Data on medication dispensing, patient prescriptions, and potential drug interactions.
7. **Administrative Systems:** Patient registration, scheduling, billing, and claims data offer insights into operational efficiency and revenue cycle management.
8. **External APIs:**
  - Public Health Databases:** Access data from sources like the CDC for disease surveillance, population health trends, and public health alerts.
  - Health Information Exchanges (HIEs):** Securely share patient data with other healthcare providers involved in a patient's care.
  - Third-Party Apps:** Integrate with patient-facing apps or wellness programs to gather data on patient-reported outcomes, lifestyle choices, and engagement in their care.
9. **Genomics Data:** As personalized medicine advances, integrating genomic data can inform diagnosis, treatment selection, and risk assessment.
10. **Social Determinants of Health (SDOH) Data:** Non-medical factors like socioeconomic status, access to transportation, and environmental conditions influence health outcomes. Incorporating SDOH data allows for a more holistic understanding of patient needs.

## Data Formats and Considerations

These data sources come in various formats:

- **Structured Data:** Found in EHRs, LIS, and administrative systems, this data is organized in a predefined manner (tables with rows and columns).
- **Unstructured Data:** Includes clinical notes, medical images, and sensor data. Requires more complex processing techniques (natural language processing, image analysis) to extract meaningful information.
- **Semi-structured Data:** Often found in external APIs or data feeds, this data may have some organization (like JSON or XML) but lacks the rigid structure of relational databases.

## Key Considerations for MeddGGenius:



- **Data Volume and Velocity:** Healthcare data is growing exponentially. The platform must handle high volumes of batch data and real-time streams from IoT devices.
- **Data Variety:** The ability to ingest and process different data formats is essential.
- **Data Veracity:** Ensuring data accuracy and completeness is crucial for reliable insights. Implement data quality checks and validation rules during ingestion.
- **Data Security and Privacy:** HIPAA compliance is essential. Secure data transmission, storage, and access control are non-negotiable.
- **Interoperability:** Seamless data exchange with existing systems and external partners is vital for care coordination and comprehensive analysis.

## Data Ingestion Mechanisms

- **AWS Kinesis:** Ideal for real-time data streams from IoT devices and patient monitoring systems, enabling immediate response to critical events.
    - **Kinesis Data Streams:** Handles high-volume streaming data.
    - **Kinesis Firehose:** Simplifies loading streaming data into S3 or other AWS services.
  - **AWS S3:** Central storage for all data, including batch uploads from EHRs, imaging systems, and external sources.
    - **S3 Transfer Acceleration:** Optimizes data transfer speeds for large datasets.
    - **S3 Storage Classes:** Choose cost-effective storage tiers based on data access patterns.
  - **AWS Database Migration Service (DMS):** Facilitates migrating data from existing databases to the data lakehouse.
  - **AWS DataSync:** Accelerates moving data between native storage and S3.
  - **API Gateways:** Securely connect to external APIs for data exchange.
- 

## 2. Ingestion and Storage Layer

**Purpose:** This layer serves as the foundation for ingesting, organizing, and securely managing raw data from a wide range of sources, including IoT devices, EHR systems, and external APIs. It is designed to handle both batch and real-time data efficiently, and integrate structured, semi-structured, and unstructured datasets in a seamless way. By utilizing AWS S3 as the central data lake and AWS Kinesis for real-time streaming, the layer supports use cases ranging from

historical data analysis to time-sensitive decision-making, such as patient monitoring or operational alerts.

This ingestion and storage layer also makes sure that data is organized and categorized logically for easier accessibility and processing. Real-time data from IoT devices streams into Kinesis and is quickly stored in S3. On the other hand, batch data such as EHRs and imaging files are directly uploaded into S3 buckets for secure storage and analysis. Together, these tools provide a scalable, durable, and highly reliable storage solution that meets the needs of healthcare data workflows.

By focusing on both efficiency and security, this layer prepares that all ingested data is readily available for downstream analytics. It plays an important role in ingesting raw data with the advanced analytics and reporting layers.

### Key Components:

- **AWS S3 (Data Lake):**
  - Acts as a centralized repository for storing raw, structured, and unstructured data.
  - Supports various file formats, including JSON, CSV, Parquet, and images, enabling flexibility in data management.
  - Enables cost-effective and scalable storage, ensuring all data is readily available for downstream processes.
- **AWS Kinesis:**
  - Facilitates real-time data streaming from sources like IoT devices, APIs, and external systems.
  - Ideal for time-sensitive use cases such as monitoring patient vitals or streaming healthcare sensor data.
  - Provides durable storage and near-instantaneous ingestion of streaming data into AWS S3.

### Flow:

#### 1. Real-time Streaming Data:

- IoT devices (e.g., wearable health monitors or smart medical equipment) and external APIs generate real-time data.
- Data streams directly into **AWS Kinesis**, where it is immediately processed and pushed into **AWS S3** for centralized storage.
- Example: heart rate, blood pressure, or device usage logs.

#### 2. Batch Data:

- Batch data sources, such as EHR, medical imaging files (e.g., X-rays or MRIs), and unstructured documents (e.g., clinical notes), are uploaded directly to **AWS S3** via secure APIs or batch upload tools.

- The data is categorized into appropriate S3 buckets for easier management and

### Relationship with Other Layers:

#### 1. Data Transformation & Preparation Layer

- The Ingestion and Storage Layer (AWS S3, Kinesis) provides raw data for the Data Transformation Layer (AWS Glue) to clean, structure, and prepare data for analytics.
- Then the AWS Glue pulls data directly from S3 or processes real-time streams from Kinesis for further enrichment.

#### 2. Ingestion & Storage Layer and Security Group

- The Security Group (AWS IAM, Lake Formation) manages access control and enforces data governance for S3 and Kinesis, ensuring only authorized users or systems can interact with ingested data.
- This layer protects sensitive data and prevents unauthorized access while maintaining compliance.

#### 3. Ingestion & Storage Layer and Compliance, Monitoring & Audit Layer

- Compliance Tools (CloudTrail, KMS) monitor and log all interactions with S3 and Kinesis, ensuring auditability and adherence to regulatory requirements.
- KMS encryption secures stored data, while monitoring tools flag anomalies to ensure compliant operations.

### Key Considerations for MeddGGenius:

- **Scalability:** Ensure the system can scale to handle increasing data volumes, from batch uploads to real-time IoT streams, without performance degradation.
- **Diverse Data Handling:** Support structured (e.g., EHRs), unstructured (e.g., medical images, clinical notes), and semi-structured (e.g., JSON, XML) data formats.
- **Data Quality:** Implement automated validation and enrichment processes to ensure reliable, clean data for downstream analytics and decision-making.
- **Data Security & Compliance:** Maintain HIPAA and regulatory compliance with encryption, strict access control, and continuous monitoring for unauthorized access.
- **Seamless Integration:** Enable smooth data exchange with existing systems, external APIs, and downstream tools, ensuring operational efficiency.
- **Real-Time Processing:** Ensure real-time ingestion and processing capabilities for critical use cases like patient monitoring and diagnostics.
- **Data Governance:** Establish clear policies for data ownership, retention, and access, ensuring proper management throughout the data lifecycle.

### 3. Data Governance and Security Layer

**Purpose:** Ensuring robust data governance and security is paramount for MeddGGenius to protect sensitive healthcare information and comply with regulatory standards. The Data Governance and Security Layer is designed to organize, secure, and manage data access effectively.

#### Key Components:

##### 1. AWS Lake Formation:

- **Data Cataloging and Permissions Management:** AWS Lake Formation simplifies the process of cataloging data, managing permissions, and enforcing governance policies across various AWS services. It allows for the centralization of permissions management for data resources, including databases and tables, within the AWS Glue Data Catalog.
- **Fine-Grained Access Control:** Lake Formation enables the implementation of fine-grained access controls, such as column-level and row-level security, ensuring that users access only the data necessary for their roles.

##### 2. AWS Identity and Access Management (IAM):

- **Role-Based Access Control (RBAC):** AWS IAM facilitates the creation and management of roles with specific permissions, ensuring that only authorized users can interact with sensitive data. This aligns with the principle of least privilege, a critical aspect of HIPAA compliance.
- **Multi-Factor Authentication (MFA):** Implementing MFA adds an extra layer of security, requiring users to provide multiple forms of verification before accessing sensitive information.

#### Data Flow:

Data stored in Amazon S3 is cataloged and secured using AWS Lake Formation, which enforces compliance with HIPAA and NIST standards. Lake Formation manages fine-grained access permissions, while AWS IAM controls role-based access, ensuring that only authorized personnel can access sensitive data.

#### Compliance Considerations:

##### 1. HIPAA Compliance:

AWS provides a secure environment that supports the processing, maintenance, and storage of protected health information (PHI), enabling covered entities and their business associates to comply with HIPAA regulations.

## 2. NIST Cybersecurity Framework:

By leveraging AWS services, organizations can align with the NIST Cybersecurity Framework, implementing risk management controls to protect data confidentiality, integrity, and availability.

By integrating AWS Lake Formation and AWS IAM, MeddGGenius can establish a robust data governance and security framework that ensures compliance with regulatory standards and protects sensitive healthcare data.

## 4. Data Transformation and Preparation Layer

### Key Purpose

This layer focuses on:

1. **Data Cleaning:** Removing inconsistencies, duplicates, and errors from the raw data to enhance data quality.
2. **Data Transformation:** Standardizing formats, converting unstructured data into structured forms, and applying domain-specific transformations.
3. **Data Preparation:** Organizing the data in a format suitable for analytics or machine learning, reducing the complexity for end-users and downstream services.

### Flow of Operations

1. **Data Extraction:**
  - AWS Glue connects to data sources such as S3 and Kinesis to retrieve raw data. For example, it could extract clinical notes, IoT-generated patient monitoring data, or external API data streams.
2. **Data Transformation:**
  - Raw data undergoes cleaning processes, such as handling missing values, resolving inconsistencies, and deduplication.
  - Unstructured data (e.g., clinical notes, images, or logs) is converted into structured formats like tabular datasets or JSON, suitable for analysis.
  - Specific transformations may include:
    - Mapping healthcare codes (e.g., ICD-10 to SNOMED).
    - Aggregating IoT sensor data into hourly or daily summaries.
    - Formatting timestamps into standardized formats.
3. **Data Standardization:**

- Standardization processes ensure interoperability. For instance, medication records from different systems might be harmonized using common terminologies like RxNorm.
- Sensitive data fields are encrypted or masked to comply with privacy standards like HIPAA.

#### 4. Data Loading and Preparation for Analytics:

- The transformed and standardized data is written back to S3 or passed to analytics services such as **AWS SageMaker** for machine learning or **AWS Athena** for SQL-based querying.

### Benefits

1. **Automation:** AWS Glue's automation minimizes manual intervention, reducing the effort and time needed for data preparation.
2. **Scalability:** The service scales to handle vast amounts of healthcare data from multiple sources without performance degradation.
3. **Flexibility:** The ability to handle structured, semi-structured, and unstructured data ensures compatibility with diverse healthcare datasets.
4. **Enhanced Analytics:** By transforming raw data into analyzable formats, this layer unlocks the potential of advanced analytics, such as predicting patient outcomes or identifying trends in healthcare operations.

### 5. Data Processing and Analytics Layer

**Purpose:** Perform advanced analytics and predictive modeling. The Data Processing and Analytics Layer is where MeddGGenius truly transforms raw data into actionable insights. Building upon the provided foundation, let's explore the capabilities of AWS SageMaker and Athena, along with other tools and techniques to maximize this layer's potential.

### AWS SageMaker: Powering Machine Learning

SageMaker is a comprehensive machine learning service that empowers developers and data scientists to build, train, and deploy models at scale. Here's how MeddGGenius can leverage it:

#### ● Predictive Modeling:

- **Patient Risk Stratification:** Develop models to predict patient risk for readmission, complications, or specific diseases based on historical data, demographics, and real-time physiological data from IoT devices.
- **Treatment Optimization:** Predict treatment response and personalize treatment plans based on individual patient characteristics and genomic information.

- **Early Disease Detection:** Train models on medical images (X-rays, MRIs) and clinical notes to detect early signs of diseases like cancer or cardiovascular conditions.
- **Image and Text Analysis:**
  - **Medical Image Analysis:** Use computer vision models to analyze X-rays, MRIs, and pathology slides for automated diagnosis, anomaly detection, and quantitative measurements.
  - **Clinical Natural Language Processing (NLP):** Extract key information from clinical notes, such as symptoms, diagnoses, and treatment plans, to enhance structured data and enable deeper analysis.
- **SageMaker Features for MeddGGenius:**
  - **SageMaker Studio:** Provides a unified interface for data exploration, model building, training, and deployment.
  - **SageMaker Autopilot:** Automates model development for users with limited machine learning expertise.
  - **SageMaker Ground Truth:** Streamlines data labeling for supervised learning tasks.
  - **SageMaker Model Monitor:** Tracks model performance over time and alerts on drift or accuracy degradation.

## AWS Athena: Querying the Data Lakehouse

Athena enables interactive querying of data residing in S3 using standard SQL. This is crucial for MeddGGenius to perform ad-hoc analysis and gain insights from its data lakehouse.

- **Use Cases:**
  - **Patient Cohort Analysis:** Identify patient groups based on specific criteria (e.g., diagnosis, demographics, treatment history) for research, clinical trials, or targeted interventions.
  - **Population Health Management:** Analyze trends in disease prevalence, resource utilization, and treatment outcomes across different populations.
  - **Operational Reporting:** Generate reports on key performance indicators (KPIs), such as hospital readmission rates, average length of stay, and cost of care.
- **Athena's Advantages:**
  - **Serverless:** No infrastructure to manage, scales automatically based on query volume.
  - **Cost-effective:** Pay only for the data scanned by queries.

- **Integration with other AWS services:** Seamlessly integrates with Glue Data Catalog for data discovery and with QuickSight for data visualization.

## 6. Visualization and Reporting Layer

### Purpose:

This layer serves as one of the final stages of the MeddGGenius cloud-based data lakehouse architecture. Its primary purpose is to transform processed data analytics into real-time, interpretable visualizations and reports for business users and stakeholders to make data-driven decisions. By leveraging business intelligence tools such as AWS QuickSight, this layer utilizes the complex healthcare data processed throughout the data lakehouse to generate interactive dashboards and reports. These dashboards are capable of integrating data from multiple sources, such as predictive outputs from AWS SageMaker, real-time streaming data from IoT devices, and query results from AWS Athena, providing a comprehensive view of healthcare operations. As a key component in stakeholder engagement, this layer plays an essential role in delivering live and accurate monitoring of critical data metrics relevant to the needs of business and healthcare leaders. Additionally, its ability to configure automated alerts and enable secure, role-based access ensures compliance with healthcare regulations and facilitates proactive decision-making in environments where time, accessibility and accuracy are crucial.

### Amazon QuickSight:

- **Interactive Dashboards:** Create real-time, visually rich dashboards that enable stakeholders to monitor patient health, operational metrics, and predictive analytics outcomes.
- **Integration with AWS Services:** Seamlessly integrates with AWS S3, SageMaker, and Athena, allowing direct visualization of processed data and machine learning outputs.
- **Customizable Visualizations:** Offers a wide range of chart types tailored to healthcare-specific use cases, such as ICU monitoring or readmission rate analysis.
- **Real-Time IoT Data Integration:** Displays real-time streaming data from patient IoT devices, enabling continuous monitoring and timely intervention.

### Flow:

- **Data Integration:**
  - Processed data is retrieved from upstream layers, including machine learning outputs from SageMaker and SQL-based queries from Athena.
  - QuickSight leverages machine learning outputs from SageMaker, allowing predictive insights to be directly visualized.
  - QuickSight queries structured data stored in S3 via Athena, enabling users to perform interactive SQL-like operations directly within the dashboards.



- Visualization Tool:
  - AWS QuickSight accesses these insights from SageMaker and Athena to create interactive dashboards and visual reports.
  - QuickSight's in-memory data engine optimizes performance for large-scale, real-time queries, ensuring quick updates for dashboards.
  - Interactive dashboards and static reports are generated in QuickSight based on stakeholder needs.

**Examples:**

- Dashboards tracking real-time ICU patient vitals, such as heart rate and oxygen levels.
- Visualizations of patient readmission risks, displaying trends and predictions based on SageMaker models.
- Resource management dashboards showing hospital bed occupancy and staffing efficiency metrics.

**7. Compliance, Monitoring, and Audit Layer****Purpose:**

The Compliance, Monitoring, and Audit Layer is a crucial part of the modern data platform for MeddGGenius. This is the layer that guarantees the adherence of the system to regulatory compliance, allowing for the reinforcement of robust data security. This layer utilizes key AWS services to provide a complete solution to track, encrypt, and monitor all data interactions.

**Flow**

The Compliance, Monitoring, and Audit Layer is a critical component of MeddGGenius's modern data platform, ensuring regulatory compliance and robust data security. This layer leverages key AWS services to create a comprehensive system for tracking, encrypting, and monitoring all data interactions.

AWS CloudTrail serves as the cornerstone for compliance and auditing, meticulously logging every API call and system event across the AWS infrastructure. This provides an immutable audit trail essential for HIPAA compliance, allowing MeddGGenius to track who accessed what data, when, and from where.

AWS Key Management Service (KMS) manages encryption keys for data at rest and in transit, implementing envelope encryption to secure sensitive healthcare data. This approach allows for fine-grained access control to encryption keys, ensuring only authorized personnel can access sensitive information.

The data flow within this layer is designed to be both secure and transparent:

- Every data interaction is automatically logged by CloudTrail in real-time.
- KMS ensures all data is encrypted using robust algorithms that meet or exceed HIPAA and NIST standards.
- Monitoring tools continuously analyze access patterns and system behaviors, alerting administrators to potential security breaches.

### 3. Data Migration and Management Strategy

The Data Migration and Management Strategy for MeddGGenius's transition to a modern data platform is designed to be comprehensive, minimizing disruption while ensuring data integrity and security. This strategy encompasses several key components:

#### **Phased Approach:**

The migration will be performed in pre-planned phases so that the transition is smooth and does not disturb the critical operations of MeddGGenius. This will give ample time for thorough validation and testing at each stage before the new system is put into full implementation.

#### **Data Assessment:**

The first step will be a comprehensive evaluation of the existing data. It means a very thorough inventory of all data sources, including EHRs, claims data, patient visit records, and unstructured data such as medical images and clinical notes. The assessment will pinpoint the quality, inconsistencies, and potential risks in the data. From this analysis, a detailed roadmap will be created that outlines the steps to be taken in cleansing, transforming, and migrating the data onto the new platform.

#### **Prioritization:**

The migration will be initiated with the most important datasets required for patient care and core operations. This may include all current patient records, active treatment plans, and recent medical imaging data. Less critical historical data or administrative records will be scheduled for

later phases. This prioritization ensures that the most important information is available in the new system as quickly as possible.

**ETL Processes:**

AWS Glue will be used to create strong ETL pipelines, which can handle structured data from existing databases and unstructured data from various sources. Standardizing data formats, resolving inconsistencies, and preparing the data for the new lakehouse architecture require custom-developed transformations. The ETL processes will also include data quality checks and enrichment steps to improve the overall value of the migrated data.

**Data Validation:**

At each step in the migration process, proper validation needs to be followed. Automation of checks should be done on completeness, accuracy, and consistency. Reconciliation processes are required to compare the migrated data with source systems to ensure no loss or corruption of information while transferring. This also involves user acceptance testing with key stakeholders to ensure that migrated data meets all operational and analytical requirements.

**Data Governance:**

A strong data governance framework will be implemented to control and secure the data. Implemented strict access using AWS Lake Formation and IAM to ensure that only the right people have access to sensitive patient information, thus abiding by the HIPAA regulations in place.

**4. Data Retention Policy**

The data retention policy of MeddGGenius is designed to guarantee full compliance with regulatory requirements, especially HIPAA, while optimizing storage costs and ensuring data accessibility. Key components of the policy include the following:

**Regulatory Compliance:**

The data retention strategy will be very well planned to ensure that all the relevant regulations are followed, but most importantly, it is HIPAA compliant. It involves retaining patient health information for at least six years from the actual date of creation or the last effective date, whichever is later. The policy also takes into consideration other laws and regulations that may apply to healthcare data retention, including state-specific requirements or other federal mandates.

**Retention Periods:**

The retention periods for various types of data are different and depend on legal requirements, operational needs, and research values. Patient medical records are usually kept for seven years

from the last visit or treatment date, which is in line with the general period stated in the project statement. However, some records may be retained for more extended periods, such as those of minors. Administrative records, financial data, and research-related information each have their own retention schedules, carefully determined to balance compliance, operational efficiency, and potential future value.

**Storage Tiers:**

AWS S3's diverse storage options are leveraged to optimize cost-effectiveness without compromising data accessibility. The policy implements a tiered storage approach:

S3 Standard for frequently accessed, current patient data

S3 Intelligent-Tiering for data with changing access patterns

S3 Glacier for long-term archival of older records

This tiered approach ensures that data are stored in the most cost-effective manner based on their access frequency and importance while still maintaining rapid retrieval capabilities when needed.

**Data Disposal:**

The policy will have in place robust procedures for secure data disposal at the end of the retention period. This shall involve:

Automated deletion through S3 Lifecycle policies

Ensuring cryptographic erasure for sensitive data using AWS Key Management Service

Keeping proper logs of all the data disposal actions for auditing purposes; performing periodic reviews to identify and securely dispose of data that has exceeded its retention period.

## 5. Unified Data Analytics Approach

- Use Case Examples:

Here's how the data lakehouse can help MeddGGenius:

- Preventive Care: Find patients who might be at risk for chronic conditions and help them out early.
- Medical Imaging Analysis: Use fancy computer vision to analyze X-rays and MRIs for faster and more accurate diagnoses.
- Patient Vital Data Tracking: Keep an eye on patient vitals in real-time using IoT devices and catch any problems quickly.
- Inventory Management: Make sure MeddGGenius always has the right medical supplies and meds on hand without wasting anything.
- Staff Scheduling: Schedule staff more efficiently by looking at past data and predicting how many patients they'll need to care for.
- Outcome Prediction: Predict how patients will do with certain treatments and help doctors make the best decisions.

Through data processing and analytics, the data lakehouse uses AWS Sagemaker and AWS Athena for machine learning and analytics. The data analytics and processing layer is designed to provide predictive analytics modeling using healthcare data to monitor patients and operational activities. MeddGGenius can leverage these services towards use cases such as tracking patients with chronic conditions, analyzing X-rays and MRIs, patient vital tracking, inventory management, scheduling, and decision-making.

**Preventive Care:** AWS Sagemaker will train predictive models with patient data, including medical records. MeddGGenius can use this model to predict risk factors within patients and attentively monitor their patients while providing treatments to these patients. Machine learning can use algorithms to determine the likelihood of a patient developing a condition based on existing factors. AWS SageMaker would be able to identify patterns and unlock insights while providing personalized care to patients. AWS Athena is supplemented to query patient data using the data lakehouse and analyze patients with high-risk conditions.

**Patient Vital Data:** Through IoT devices such as smart medical devices or blood pressure tracking systems, vitals would be monitored to track their health status. These devices would track blood pressure, heart rate, or oxygen levels, which are processed in real-time. AWS Athena would store IoT data while collecting data in real-time to store and process for querying and reporting, providing streamlined analytics and insights. AWS SageMaker would use their model to detect early conditions such as blood pressure or cardiac arrest and when they pose a risk. Healthcare providers can provide personalized care to each patient and improve preventive care, leading to improved patient outcomes.

**Medical Imaging and Analysis:** Computer vision allows for the analysis of X-rays and MRIs. AWS SageMaker would train deep learning models and use neural networks to analyze imaging data to detect early signs of conditions based on X-ray vision. Deep learning models would help radiologists identify areas where the bone structure is experiencing pain or conditions. Using imaging data, radiologists can diagnose patients and provide treatments to patients, assisting in preventive care. AWS Athena can provide aggregated reports based on image data to track the accuracy of the diagnosis or time series analysis where conditions are critical. Using these services would lead to higher accuracy in diagnosis and improved personalized care.

**Staff Scheduling and Inventory Management:** Staff would be scheduled efficiently using AWS SageMaker to develop predictive models in areas where there may be more assistance. By analyzing historical data such as patient appointments, seasonal patterns where illnesses such as flu or colds become common, or patient conditions, MeddGGenius can efficiently utilize staff in areas where staff is needed, reducing unnecessary overstaffing and improving schedule systems. The inventory management process would be improved by analyzing historical data on medication consumption and the number of patients to provide the right amount of stock for

prescriptions. SageMaker can predict the supply-demand and ensure enough stock is in hand. Meanwhile, Athena uses inventory data to monitor supply levels and analyze items that are less utilized or items that are quickly being depleted. This data querying will allow MeddGGenius to quickly stock up on items and schedule orders at appropriate times.

## **6. Compliance Monitoring and Reporting**

### **1. HIPAA and NIST Compliance:**

The data lakehouse will be constructed to adhere to all HIPAA and NIST security standards. This will include implementing appropriate technical, physical, and administrative safeguards to protect the privacy and security of protected health information (PHI). We will also regularly review and update our security measures to ensure they are aligned with the latest regulatory requirements.

### **2. Data Encryption:**

We will employ AWS encryption services to encrypt data at rest and in transit. This will ensure that data remains confidential, even if it is intercepted or accessed by unauthorized individuals. We will use encryption keys that are managed by AWS Key Management Service (KMS) to ensure the highest level of security.

### **3. Access Control:**

We will implement strict access controls to ensure that users can only access data relevant to their roles. This will be achieved through the use of role-based access control (RBAC) and other security measures. We will also regularly review and update our access control policies to ensure they remain effective.

### **4. Audit Logging:**

We will maintain detailed logs of data access and modifications. This will help us to track and investigate any suspicious activity. We will also use these logs to generate reports for regulatory compliance purposes.

### **5. Compliance Monitoring:**

We will leverage AWS CloudTrail and other tools to continuously monitor compliance with regulatory requirements. This will help us to identify any potential areas of non-compliance. We will also use these tools to generate reports for regulatory compliance purposes.

### **6. Regulatory Reporting:**

The platform will facilitate the creation of reports for various regulatory bodies. We will be able to aggregate data from multiple sources. We will utilize reporting tools to generate customized reports. These reports will be tailored to the specific requirements of

each regulatory body. We will also provide tools to allow users to generate their own reports. We will also provide a self-service portal for users to access reports and other compliance-related information.

## 7. Conclusion

MeddGGenius can overcome its current challenges and prepare for future needs with the cutting-edge cloud-based data lakehouse solution on AWS. This platform will be modern, scalable, and compliant, empowering MeddGGenius to achieve the following key benefits:

**Improved Patient Care:** The data lakehouse will enable MeddGGenius to utilize advanced analytics and machine learning algorithms to uncover valuable insights from patient data. This will lead to more precise diagnoses, personalized treatment plans, and improved overall patient outcomes.

**Effective Data Management:** The solution will provide MeddGGenius with a centralized repository for storing and analyzing diverse healthcare data, including electronic health records, medical images, genomic data, and more. This will streamline data management processes, improve data accessibility, and facilitate collaboration among healthcare professionals.

**Enhanced Operational Efficiency:** By leveraging the cloud-based data lakehouse, MeddGGenius can automate manual tasks, reduce data processing time, and improve the efficiency of healthcare operations. This will enable the organization to allocate more resources to patient care and strategic initiatives.

**HIPAA Compliance and NIST Cybersecurity Standards:** The data lakehouse solution will be designed to ensure compliance with HIPAA regulations and meet NIST cybersecurity standards, protecting patient data from unauthorized access, use, or disclosure.

**Innovation in Healthcare Delivery:** The data lakehouse will serve as a foundation for innovation in healthcare delivery. MeddGGenius will be able to develop new applications, services, and tools that leverage patient data to improve care coordination, enhance patient engagement, and drive better health outcomes.

Overall, the proposed cloud-based data lakehouse solution will empower MeddGGenius to unlock the full potential of its data, transforming patient care, achieving strategic objectives, and positioning the organization as a leader in healthcare innovation.

## 8. References

1. Amazon Web Services. (n.d.). Amazon QuickSight. Retrieved November 27, 2024. <https://aws.amazon.com/quicksight/?amazon-quicksight-whats-new.sort-by=item.additionalFields.postDateTime&amazon-quicksight-whats-new.sort-order=desc>
2. Amazon Web Services. (n.d.). Population health applications with Amazon HealthLake, Part 1: Analytics and monitoring using Amazon QuickSight. Retrieved December 2, 2024. <https://aws.amazon.com/blogs/machine-learning/population-health-applications-with-amazon-healthlake-part-1-analytics-and-monitoring-using-amazon-quicksight/>
3. AWS Architecture Blog: AWS Architecture Blog (amazon.com)
4. Azure Architecture Blog. (n.d.). Azure Architecture Blog – Microsoft Community Hub. Retrieved November 21, 2024. <https://techcommunity.microsoft.com/category/azure/blog/azurearchitectureblog>
5. High Scalability. (n.d.). Real-world architectures: Examples like WhatsApp and Netflix. Retrieved November 25, 2024. <http://highscalability.com/>
6. Microsoft. (n.d.). *Analytics end-to-end with Azure Synapse*. Retrieved December 2, 2024. <https://learn.microsoft.com/en-us/azure/architecture/example-scenario/dataplate2e/data-platform-end-to-end>
7. Turck, M. (2021). Trends and technology landscape for data, AI, and ML architectures. Retrieved December 2, 2024. <https://mattturck.com/data2021/>