# Big Data Pipeline Architecture for Healthcare Analytics

Dipika Bogati

Radhika Srivastava

Sailesh Kafle

Sushant Nepal

# Introduction

Design a scalable, secure **data pipeline** to process and analyze 50TB of healthcare data

Handle structured (CSV), semi-structured (IoT JSON), and unstructured (medical imaging) data

Manage **50TB** of historical records and real-time data from **10,000 IoT** devices generating **1GB/hour**
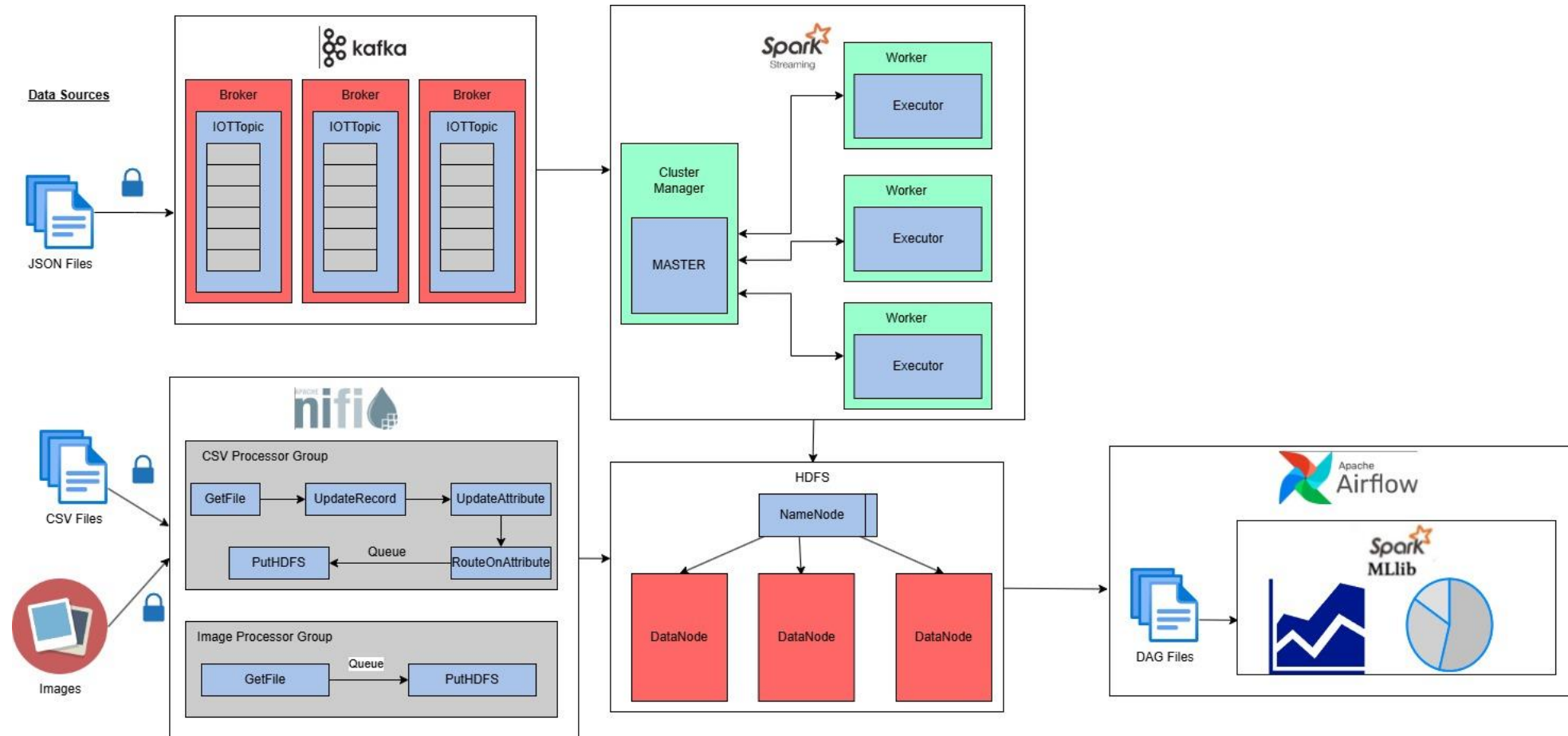
Ensure **HIPAA** compliance while managing **inconsistencies** (missing fields, duplicates)

Build a Hadoop-based ecosystem to ingest, store, and analyze this **data securely and efficiently**

# Architecture Design

# Data Ingestion and Storage: Apache NIFI

Handling batch ingestion of structured patient records (CSV)

Unstructured medical imaging files (binary format)

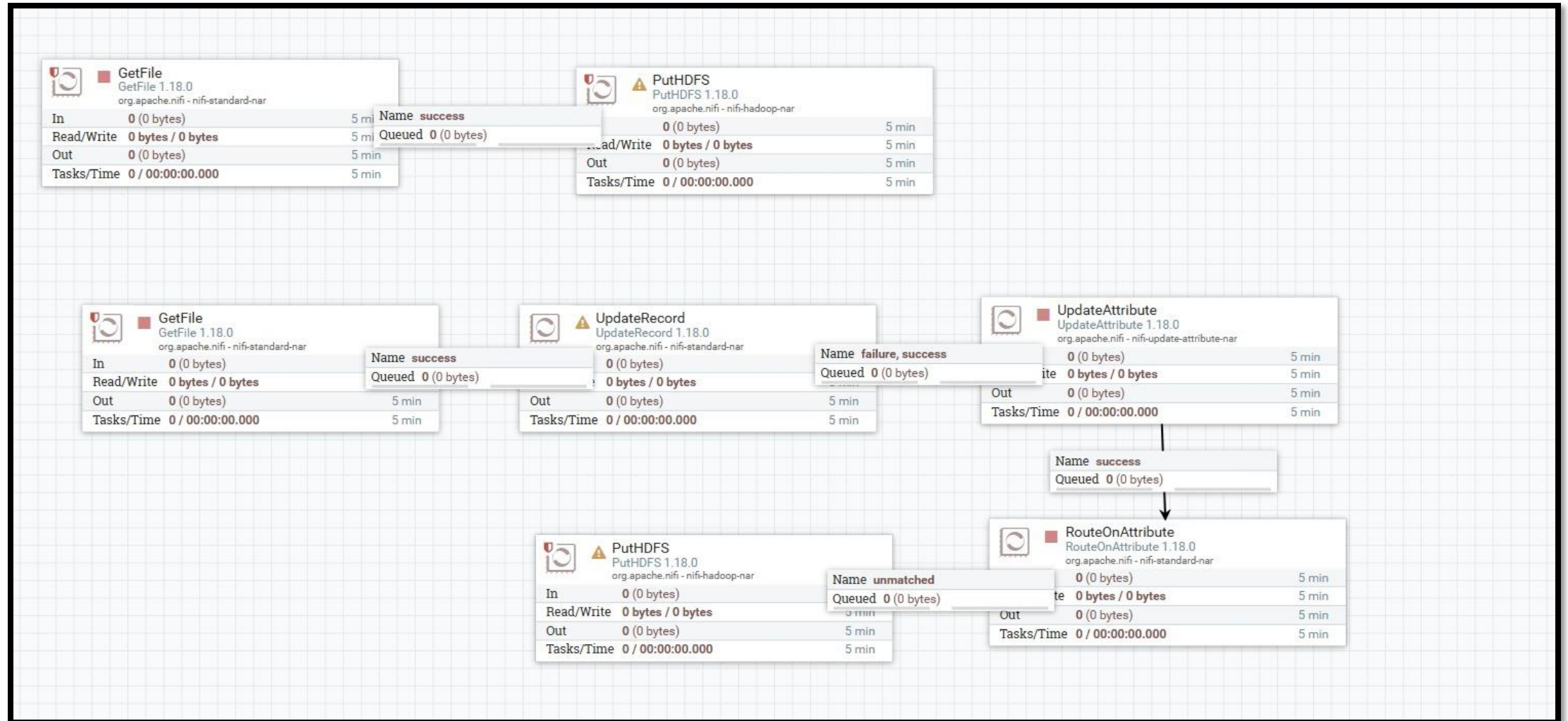Efficient handling of large volume CSV files (50 TB)

Supports binary file ingestion without loss of data quality
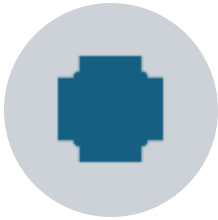
Scalable and fault-tolerant data flow management

Easily integrates with HDFS for data storage

Provides backpressure and retry mechanisms for stable data ingestion

# Apache NIFI Design Flow

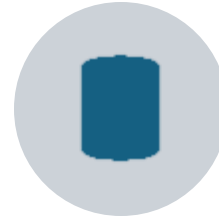# Data Ingestion and Storage: Apache Kafka

Ingesting high-speed real-time JSON data from 10,000 IoT devices.

Ensures durability and fault tolerance by replicating data across multiple brokers.

Decouples data producers (IoT devices) from consumers (Spark Streaming), allowing for asynchronous processing.

Data is stored reliably in Kafka topics for processing by downstream system

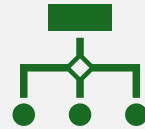Kafka partitions topics for horizontal scalability and parallel processing.

Guarantees no data loss with retention capabilities, even during processing delays.

# Apache Kafka

 Producer : IOT devices

 Kafka Cluster:

3 brokers as 3 replication

One topic (IOTTopic)
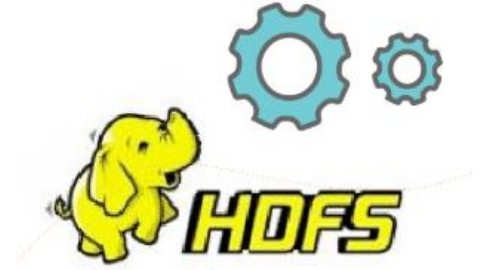
6 partitions

 Consumer : Spark Streaming

# Spark Streaming

- Consumer of the Apache Kafka.

- Processing real-time data.

- Provides data validation, cleaning, and real-time
  - Anomaly detection (e.g., heart rate spikes, low oxygen levels)
  - Null Value
  - Duplicate Value

- Processes and writes cleaned data to HDFS in Parquet format for further analysis.
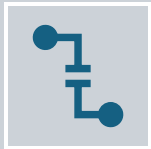
# Data Ingestion and Storage:
# Hadoop Distributed File System

Storing structured and unstructured data for large datasets (50 TB patient records, IoT data, and medical images) for scalability and reliability.

Ensures high-availability by removing secondary-point of failure and adding secondary name node.

Added warm tier and cold tier to archieve the old data that is more cost – friendly and does not delete the old data as well.

# Apache AirFlow

Ensures efficient workflow management at scale with built-in fault tolerance

Automates the schedule of Apache SparkMLLib jobs on the period of time.

Monitoring of the visualization pipeline is efficient and easy in Apache Airflow.

# SparkML

SparkML jobs are scheduled to run hourly to process and analyze the latest data

The program retrieves the processed data from HDFS for analysis

SparkML is used to build and train machine learning models on healthcare data to generate insights

The models provide intelligence that helps in solving healthcare problems, such as identifying patient trends and predicting health risks

Visualizations based on model outputs are generated to monitor and track health trends over time
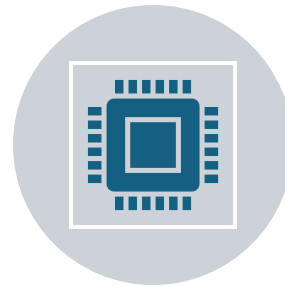
# Security Consideration

Role-based access control for ingestion, HDFS permissions for data storage, and detailed audit logging

Masking of fields (e.g., patient names, SSNs) during ingestion and processing, ensuring privacy protection.

Transparent Data Encryption (TDE) for data at rest in HDFS, ensuring protection against storage device breaches.

Real-time monitoring of the system, with centralized secret management and alerts for suspicious activities
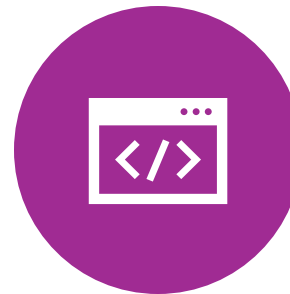
# Conclusion

Warm tier and Cold tier for old data to reduce storage and cost.

Minimum privileges, masking PII data, SSL/TLS data ingestion maintain the HIPAA policy.

Decided the resources in such a way that there will be no resource over-head and efficient for the data flow.

Pipeline is scalable as resources can be added in the system without designing if needed in the future.

Any Questions?