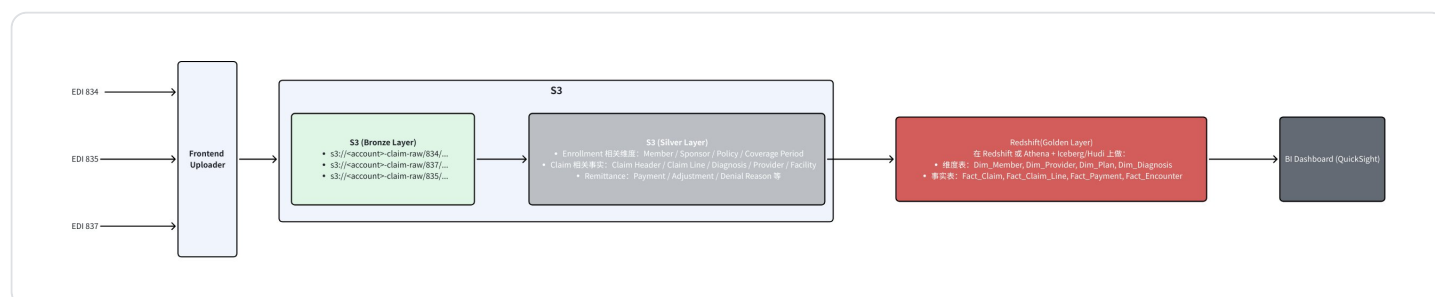


Claim Management System — Logic Optimization & Operability Enhancement



1. 目标重申（用一句话对齐）

在 AWS 上搭一个 **湖仓一体（Lakehouse）** 的 *Claim Data* 平台，

- 上游拿到的是 EDI 834 / 835 / 837（或已转换/可转换为 CSV）
- **S3 作为唯一原始存储层（Raw Layer）**，原始只能是 CSV
- 满足合规（HIPAA：加密、访问控制、审计、数据血缘）
- 下游既能做 **OLAP 分析 / BI**，也能支撑后续 **实时/批处理风控 & 报表**

下面按“分层 + 组件”来规划。

2. 总体架构视图（从左到右）

1.1 数据流主线（业务视角）

1. 外部/上游系统

- Policy Admin / Claim System / Clearing House 提供：
 - EDI 834 (Enrollment)
 - EDI 837 (Claims)
 - EDI 835 (Remittance)
- 形式可以是一个我自建的前端，可选文件上传格式（支持x12/csv），

2. Ingestion & Landing（落地 + 安全）

- 统一通过 **Ingestion Bucket (S3 Raw)** 落地：

- `s3://<account>-claim-raw/834/...`

- `s3://<account>-claim-raw/837/...`

- `s3://<account>-claim-raw/835/...`

- 要求:

- 对于CSV: 直接存储到raw

- 对于X12:

- 上游上传 **X12** 到 `claim-edi-archive` (WORM + 合规保留策略)。

- S3 Event / Step Functions 触发 **解析 Lambda / Glue Job**:

- 读取 X12 → 按你定义好的 mapping 转成一到多个 CSV (member / claim_header / claim_line / payment ...)

- 解析结果 CSV 写入 `claim-lake/raw/...`, 这里才是 **数据湖的 raw 层**。

3. 标准化 & 清洗 (Bronze → Silver)

- 使用 Glue / EMR / Lambda 将原始 CSV 转成 **结构化表**:

- Enrollment 相关维度: Member / Sponsor / Policy / Coverage Period

- Claim 相关事实: Claim Header / Claim Line / Diagnosis / Provider / Facility

- Remittance: Payment / Adjustment / Denial Reason 等

- 输出到 **S3 Silver 层** (仍然是 CSV 或 Parquet, 建议 Parquet + 分区)。

4. 数仓建模 (Silver → Gold)

- 在 **Redshift 或 Athena + Iceberg/Hudi** 上做:

- 维度表: Dim_Member, Dim_Provider, Dim_Plan, Dim_Diagnosis

- 事实表: Fact_Claim, Fact_Claim_Line, Fact_Payment, Fact_Encounter

- 形成一个 **可分析的星型模型**, 支撑:

- Claim Cost 分析

- Provider Performance

- Fraud & Overutilization 检测

- Enrollment & Retention 分析

5. 消费层 (BI & API)

- BI: QuickSight (SPICE) 看 dashboard (按 payer/provider/member 分析 claim)。

- Ad-hoc 分析: Athena / Redshift Query Editor。

- 下游服务: 通过 API Gateway + Lambda 或 Redshift Data API 暴露查询结果 (比如给 Risk Engine / Care Management 系统)。

3. 分层设计（Lakehouse 三层 + Catalog）

你刚好要求“数湖一体”，我们可以按**经典 Bronze / Silver / Gold + Catalog**来设计。

2.1 Bronze（Raw / Landing Layer —— 必须是 S3 + CSV）

目的：合规 + 不可篡改 + 数湖入口

存储：

- S3: `claim-raw-bucket`
- Key 命名示例：
 - `raw/834/year=2025/month=11/day=21/source=clearinghouse/file_20251121_001.csv`
 - `raw/837/year=2025/...`
 - `raw/835/year=2025/...`

要求：

1. 安全性：

- 开启 **S3 Versioning + MFA Delete**（防止误删/篡改）。
- 开启默认 **S3 加密（SSE-KMS）**，KMS CMK 单独 Key（例如 `kms-key-claim-raw`）。
- **Bucket Policy** 严控：只允许特定 VPC Endpoint / 角色访问。

2. 审计 & 血缘：

- 使用 **CloudTrail + S3 Access Logs** 跟踪访问。
- 每个文件落地时写一条 **元数据记录**（Glue Table 或 DynamoDB）：
 - `file_name, file_type(834/837/835), source_system, ingest_time, checksum, record_count, uploader`

这一层你可以理解为“合规黑盒”：只进不改，只读解，强加密，可审计。

2.2 Silver（Standardized / Clean Layer）

目的：从“行级 CSV”变成“业务实体级的结构化表”

存储：

- 专门一个 S3 bucket 或同 bucket 不同 prefix：
 - `s3://claim-lake/silver/834_member/`
 - `s3://claim-lake/silver/837_claim_header/`
 - `s3://claim-lake/silver/837_claim_line/`

- `s3://claim-lake/silver/835_payment/`

特点:

- Schema 已经和**业务实体**对齐:
 - 834 → Member, Coverage, Enrollment Event
 - 837 → Claim, Claim Line, Diagnosis, Service Provider
 - 835 → Payment, Adjustment, Remittance Advice
- 数据做了基础清洗:
 - 类型转换 (日期、金额、枚举字段)
 - 去重、基础校验 (必填、PK/FK 合法性)
 - 记录解析错误到 **Error Table / Error Path** (比如 `silver/_error/...`)

计算引擎:

- 批处理:
 - **AWS Glue Job (Spark)**: 批量每天/每小时跑, 解析 CSV → 结构化表。
- 实时/准实时:
 - 如果将来需要 near real-time, 可以考虑 Kinesis + Lambda 先入 raw, 然后 Glue Streaming/EMR Streaming 转入 Silver。

2.3 Gold (Warehouse / Serving Layer)

1. Redshift (推荐做 Enterprise Data Warehouse)

- 把 Silver (S3 上的 Parquet) 通过 Glue Catalog 映射到 Redshift 外部表或者直接 COPY 进内部表。
- 建模为典型星型:
 - Dimensions:
 - `dim_member`
 - `dim_provider`
 - `dim_plan`
 - `dim_diagnosis`
 - `dim_facility`
 - Facts:
 - `fact_claim` (一行一 claim)
 - `fact claim line` (一行一行项目)

- `fact_payment` (从 835)
 - `fact_eligibility` (可从 834 enrollment 状态推)
-

4. 数据目录 & 元数据：Glue Data Catalog

要做“数湖一体”，**Catalog 是中枢神经**。

3.1 Catalog 作用

1. 记录所有 **数据库/表/字段/分区** 信息 (Raw/Silver/Gold) 。
2. 让 Athena / Glue / Redshift Spectrum / EMR 用同一套 schema。
3. 后续可以接 **Data Governance 工具** (标签、PII 标注、访问控制) 。

3.2 拆分建议

在 Glue Catalog 中建立几个 Database:

- `claim_raw_db`
 - `tbl_834_raw_csv`
 - `tbl_837_raw_csv`
 - `tbl_835_raw_csv`
 - `claim_silver_db`
 - `tbl_member`
 - `tbl_claim_header`
 - `tbl_claim_line`
 - `tbl_payment`
 - `claim_gold_db`
 - `dim_member`
 - `dim_provider`
 - `dim_plan`
 - `fact_claim`
 - `fact_payment`
-

5. 合规与安全设计 (HIPAA 角度)

4.1 存储安全（静态）

- 所有 S3 bucket 开启 **SSE-KMS**，使用独立的 CMK（raw/silver/gold 可共用或分开）。
- Redshift 使用 **KMS 加密**，且在 VPC 私网中。
- 对含 PHI 的字段打标签（比如 `phi_pii=true`），将来可配合 Column Masking / Row-Level Security。

4.2 访问控制（访问中）

- 全面使用 **IAM Role + ABAC/Tag-Based Policy**:
 - 比如 `role-claim-analyst` 只能访问 `gold` 层的某些表。
 - `role-claim-etl` 才能访问 `raw/silver` + KMS key decrypt。
- 使用 **Lake Formation**（如果走 S3 Lake）做细粒度权限（表、列、行级）。

4.3 审计 & 监控

- **CloudTrail**：记录所有数据平面/控制平面操作。
- **CloudWatch + CloudWatch Metric/Log**:
 - Glue Job 失败/成功次数
 - ETL 延迟、文件处理量
 - S3 访问错误/拒绝统计
- 可以搭一个简单的 “**Data Pipeline Health Dashboard**”（QuickSight 或 CloudWatch Dashboard）。

6. Ingestion 层设计（文件如何进入 S3 Raw）

虽然你现在问题主要在数据中心/湖仓本身，但 **入口设计** 对合规也很关键：

1. SFTP / FTPS：

- 使用 **AWS Transfer Family(SFTP)** → 直接写到 S3 Raw Bucket。

2. API Upload：

- API Gateway + Lambda → 验证签名/格式 → 写入 S3 Raw。

3. 内网对接：

- 如果是同一 VPC 内其他系统，可以通过 VPC Endpoint 直接 putObject 到 S3。

所有入口统一一个原则：**只写入 CSV 到 Raw**，禁止绕过。

7. 处理 & 调度 (Workflow Orchestration)

为了系统可管理，需要一个统一调度：

- **Step Functions / Managed Airflow (MWAA)** 二选一：
 - 典型流程（以日批为例）：
 - i. 监听新文件（S3 Event → SQS → StepFunctions/Airflow Trigger）。
 - ii. 任务 1：调用 Glue Job 解析 834/837/835 → 写入 Silver。
 - iii. 任务 2：数据质量检查（记录数对比、必填字段检查，异常报警）。
 - iv. 任务 3：将 Silver 数据加载/物化到 Gold（Redshift/ Iceberg）。
 - v. 任务 4：刷新 QuickSight SPICE dataset 或触发通知。

先直接按你这套架构，把“会用到什么技术栈”拆成三块讲：

1. **AWS Services 全家桶**（按分层对应一下）
 2. **Version Control / IaC / DevOps 技术栈**
 3. **文件上传网站的技术选型**（前后端 + 在 AWS 上怎么落地）
-

8. 这里会用到哪些 AWS Services?

1) 存储 & 数据湖 (Raw / Silver / Gold)

- **Amazon S3**
 - Raw: `claim-raw-bucket`
 - Silver / Gold: `claim-lake-bucket`
 - 配置：
 - Versioning + (可选) Object Lock / WORM
 - SSE-KMS 加密
 - S3 Access Logs (记录访问行为)
- **AWS KMS**
 - KMS CMK：
 - `kms-claim-raw`
 - `kms-claim-silver-gold`
 - 用于 S3 + Redshift + Glue Job 的加密解密

如果你要做长期归档 / 法规保留：

- **S3 Glacier / Glacier Deep Archive**
 - 存历史 X12 原文、老的 Raw CSV 归档
-

2) Ingestion & API 入口

- **AWS Transfer Family (SFTP)**
 - 让 Clearing House / 上游系统通过 SFTP 直接丢文件到 S3 Raw。
 - **Amazon API Gateway**
 - 做“文件上传网站”的后端入口（REST/HTTP API），前端调用 API 拿 pre-signed URL 或直接 POST。
 - **AWS Lambda**
 - 文件到达 S3 → 触发 Lambda：
 - 校验文件名/格式/元数据
 - 调用解析逻辑（或者丢到 SQS/Kinesis）
 - **Amazon SQS**
 - S3 Event → SQS → Step Functions / Lambda 做异步处理、重试、DLQ。
 - （可选）**Amazon Kinesis Data Streams / Firehose**
 - 如果后面有 near real-time claim event 流式场景，可以把事件打入 Kinesis。
-

3) 计算 & ETL & Lakehouse

- **AWS Glue**
 - Glue Crawler：建 Raw / Silver 表的 schema（注册到 Glue Data Catalog）。
 - Glue Job（Spark）：
 - 解析 CSV → 标准化成 Member / Claim / Payment 等实体（Silver）。
 - 从 Silver 聚合 / 转换 → Gold 表（给 Redshift/Catalog）。
 - Glue Workflow：可当简单的 pipeline orchestration。
- （可选）**Amazon EMR / EMR Serverless**
 - 如果后面有大规模批处理、机器学习、复杂 Spark Job，用 EMR 更灵活。
- **Amazon Athena**
 - 直接查 S3 上的 Raw / Silver / Gold 表（通过 Glue Catalog）。
 - 给数据科学 / ad-hoc 查询用。

4) 数据仓库 / Gold 层 Serving

- **Amazon Redshift (推荐用 Redshift Serverless)**

- 建 star schema:

- `dim_member / dim_provider / dim_plan / dim_diagnosis / dim_facility`
- `fact_claim / fact_claim_line / fact_payment / fact_eligibility`

- 集成方式:

- Spectrum 外部表直接扫 S3 (Silver / Gold Parquet)
- 或定期 COPY / MERGE 到 Redshift 内部表 (性能更好)

5) Catalog & 权限治理

- **AWS Glue Data Catalog**

- Database:

- `claim_raw_db / claim_silver_db / claim_gold_db`

- **AWS Lake Formation**

- 表/列级权限控制 (谁能查 PHI、谁只能看汇总)。
- 再往上可以做 Tag-based access (例如: `phi=true` 的列只有特定角色可见)。

6) 安全、合规 & 审计

- **AWS IAM**

- Role 设计:

- `role-claim-ingestion` 只负责一个方向的写入 Raw
- `role-claim-etl` 负责 Glue / Lambda / KMS 解密等
- `role-claim-analyst` 只看 Gold / 部分 Silver

- **AWS CloudTrail**

- 记录所有控制平面 & 数据平面操作 (谁创建了 bucket、谁查了什么表)。

- **Amazon CloudWatch**

- 监控 Glue Job / Lambda / Step Functions 的运行情况。
- 告警: 延迟、失败率、异常访问。

- (可选) **Amazon Macie**
 - 自动扫描 S3，识别 PII/PHI，帮你验证合规风险。
-

7) BI & 报表 & Dashboard

- **Amazon QuickSight**
 - 直接连：
 - Athena (查 S3 Gold / Redshift)
 - Redshift
 - 开 SPICE 做加速，给业务用户看 dashboard。
-

8) 网络 & 入口安全

- **Amazon VPC + VPC Endpoint**
 - 确保 Glue / Lambda / Redshift 访问 S3 走内网。
 - **Amazon CloudFront**
 - 前端静态网站 (React 构建结果) + API Gateway 的边缘加速。
 - **AWS WAF**
 - 给 API Gateway / CloudFront 加 web 防火墙 (防 SQL 注入、常见攻击)。
 - **Amazon Route 53**
 - 域名 + DNS + ACM 证书 (HTTPS)。
-

9. Version Control / IaC / DevOps 技术栈

你这套平台可以完全按“代码化的平台”来做：

1) Version Control

- **Git + GitLab / GitHub**
 - Repo 分组建议：
 - `infra-claim-platform`：只放 Terraform / CDK / CloudFormation
 - `etl-claim-glue-jobs`：Glue Job / Spark / PySpark 代码
 - `api-claim-upload-portal`：前端 + 后端代码
 - Branch 策略：

- `main` / `dev` / feature branches
 - 合并到 `main` 触发 prod 环境部署
-

2) Infra as Code (IaC)

- **Terraform** (你刚提到的, 强烈推荐)
 - 模块划分:
 - `network` : VPC, Subnet, Security Group, VPC Endpoint
 - `security` : IAM Role, KMS, CloudTrail, Config
 - `data-lake` : S3 Raw/Silver/Gold, Glue Catalog, Lake Formation
 - `warehouse` : Redshift, Workgroup, Namespace
 - `app` : API Gateway, Lambda, Cognito, CloudFront, Route 53
 - Backend:
 - S3 + DynamoDB 做 Terraform Remote State & State Lock。
 - (可选) **AWS CDK**
 - 用 TypeScript/Python 写复杂资源 (比如 CloudFront + S3 + Lambda@Edge 的组合)。
 - 也可以只用 Terraform, 就不用 CDK 了, 看你团队习惯。
-

3) CI/CD

- **GitHub Actions / GitLab CI/CD**
 - Pipeline 示例:
 - Infra Pipeline** (针对 `infra-claim-platform` repo)
 - `terraform fmt` / `terraform validate`
 - `terraform plan` (在 MR/PR 上做 review)
 - `terraform apply` (合并到 main 时自动 apply 到 prod)
 - ETL Pipeline** (针对 Glue Job 代码)
 - 单元测试 (PyTest)
 - 打包 Zip/Whl 上传到 S3 / ECR
 - 调用 `aws glue update-job` / `create-job` API
 - App Pipeline** (文件上传网站)

- 前端: `npm test` / `npm build` → 把 build 结果 sync 到 S3 + 触发 CloudFront Invalidatation
 - 后端: 打包 Lambda / Docker 镜像 → 部署到 Lambda / ECS
-

4) 配置 & Secrets 管理

- **AWS Systems Manager Parameter Store**
 - 存: 环境配置 (例如 S3 bucket name, 业务开关)。
 - **AWS Secrets Manager**
 - 存: 数据库密码、第三方 API Key 等敏感信息。
 - 使用 IAM Role + KMS 加密访问这些参数。
-

5) 数据质量 & 测试 (可选但强烈推荐)

- **Great Expectations / Soda Core**
 - 对 Silver → Gold 过程加 data quality check, 例如:
 - record count 对比
 - PK 不为空
 - 金额范围合理
 - 通过 Glue Job / Lambda 调用。
-

10. 文件上传网站应该用什么技术?

你这个“文件上传网站”主要需求是:

- 用户 (上游系统或 internal ops) 登录
- 上传 EDI X12 或 CSV
- 写入 S3 Raw (带元数据)
- 触发后续 ETL

推荐一个 **偏工程化 + Serverless** 的组合:

架构概览

前端: React SPA → S3 + CloudFront

后端: API Gateway + Lambda (Python/FastAPI 或 Node.js/NestJS 风格)

鉴权: Cognito

1) 前端技术栈

- **React + TypeScript**
 - UI 框架：Ant Design / MUI / Chakra UI 随便挑一个。
 - 功能：
 - 登录（Cognito Hosted UI 或 SDK）
 - 文件选择 / 拖拽上传
 - 调用后端 API 拿 pre-signed URL
 - 展示上传进度、历史上传记录
 - 部署：
 - 使用 `npm build` → 将 build 静态资源部署到 S3
 - 使用 CloudFront 作为 CDN + HTTPS 入口
-

2) 后端（API 层）

- **Amazon API Gateway**
 - 提供 REST API：
 - `POST /upload-request`：生成 pre-signed URL + 记录 metadata
 - `GET /uploads`：查询历史文件
- **AWS Lambda**
 - 运行语言建议：
 - Python（FastAPI + Mangum / 纯 handler，和你 ETL 生态统一）
 - 或 Node.js + NestJS / Express
 - 主要逻辑：
 - 校验调用者身份（通过 Cognito JWT）
 - 生成 S3 pre-signed URL（指定 bucket、key、文件大小限制、Content-Type）
 - 把一条上传记录写入 **DynamoDB**：
 - `file_id, file_name, file_type, uploader, upload_url_expire_at, status`

- 上传完成后可以用 S3 Event 通知 Lambda，更新 DynamoDB 中 status = `INGESTED`，再丢消息到 SQS/Step Functions 继续 ETL。
-

3) 鉴权 & 用户管理

- **Amazon Cognito**
 - User Pool：管理用户登录 / 注册 / MFA。
 - Identity Pool（可选）：映射到不同 IAM Role，控制他们可以上传哪些 payer 的数据。
 - 前端用 Cognito SDK 获取 ID Token / Access Token，带入 API Gateway。
-

4) 部署 & 运维

- 前端：
 - CI: GitHub Actions / GitLab CI
 - 自动：
 - build → sync to S3 → CloudFront Invalidation
 - 后端：
 - CI：
 - 单测 (PyTest / Jest)
 - 打包 Lambda zip
 - `aws cloudformation deploy` 或 Terraform 自动更新 Lambda + API Gateway
 - 监控：
 - CloudWatch Logs：记录 API 请求、Lambda 错误
 - CloudWatch Alarm：错误率、延迟、5xx 异常等
-

如果你要一句“总结式 stack 介绍”给简历/文档用

你可以这样写一版简短 summary（英文）：

Built a HIPAA-compliant claim Lakehouse on AWS using S3 (raw/silver/gold), Glue, Athena, and Redshift Serverless, governed by Glue Data Catalog + Lake Formation, with end-to-end ingestion orchestrated by S3 Events, Lambda, Step Functions, and SQS.

Infrastructure is fully managed as code via Terraform and GitLab CI/CD, while a React + API Gateway + Lambda + Cognito upload portal enables external partners to securely submit X12/CSV claim feeds directly into the raw layer.

如果你愿意，我可以下一步帮你：

- 把上面这些技术栈梳成一张“**组件清单表**”（列：功能模块 / AWS Service / 非 AWS 技术 / 备注），
- 或者直接帮你写一段“项目架构说明书”的「技术栈章节」。