# Clinical Trials Intelligence Platform

## End-to-End Serverless Data Engineering Pipeline on Google Cloud Platform

# Index

# 1. Introduction

## 1.1 Project Background

Clinical trial data is distributed across multiple international registries and regulatory platforms. These sources differ in format, structure, and completeness, making centralized analytics difficult.

## 1.2 Business Problem

Manual ingestion and cleaning of registry data was:

- Slow
- Error-prone
- Non-scalable
- Not idempotent
- Lacking monitoring and automation

## 1.3 Project Objective

Design and implement a fully automated, idempotent, serverless, monitored, scalable data pipeline on Google Cloud Platform that integrates multi-source clinical trial data into a unified analytics warehouse.

---

# 2. Data Sources

## 2.1 ClinicalTrials.gov (API Subset)

- Accessed via REST API
- JSON format
- Condition-filtered extraction

## 2.2 WHO ICTRP

- Aggregated global registry dataset
- Structured downloadable format

## 2.3 EU Clinical Trials Register (EUCTR)

- Structured CSV/TXT datasets
- Drug trials in EU/EEA

## 2.4 ISRCTN Registry

- API / CSV export
- Global intervention studies

## 2.5 EMA Clinical Data

- Regulatory study reports
- Structured data extracts

---

# 3. Data Characteristics & Handling Strategy

## 3.1 Multi-Format Data

Sources provided:

- JSON (API responses)
- CSV
- Excel
- TXT

All formats were programmatically parsed and converted into Pandas DataFrames before normalization.

---

## 3.2 Inconsistent Schemas

Observed differences:

- Different column names for same field

- Nested vs flat structures
- Status value inconsistencies

Solution:

- Canonical schema mapping
- Defensive extraction using `.get()`
- Standardized field names before warehouse load

---

# 3.3 Missing Identifiers

Some records lacked:

- Secondary IDs
- Consistent registry identifiers

Solution:

- Used primary registry ID
- Validation warnings logged
- MERGE-based idempotent load

---

# 3.4 Free-Text Conditions

Medical conditions were:

- Non-standardized
- Multi-value strings
- No ontology mapping

Decision:
Stored raw values.
Deferred ontology harmonization to future enrichment phase.

---

# 3.5 Encoding Issues

Observed:

- Non-ASCII characters
- Special characters in sponsor names

Solution:
Explicit UTF-8 handling during ingestion.

---

## 3.6 Sponsor Inconsistencies

Example:

- "Pfizer"
- "Pfizer Inc."
- "Pfizer Ltd"

Decision:
Preserved raw sponsor text.
Entity resolution deferred.

---

# 4. System Architecture (GCP)

## 4.1 Architecture Overview

Serverless stack:

- IAM
- BigQuery
- Docker
- Artifact Registry
- Cloud Run (Jobs)
- Cloud Scheduler
- Cloud Logging
- Log-based Alerts

---

## 4.2 IAM & Security Configuration

Configured service accounts with:

- BigQuery Data Editor
- BigQuery Job User
- Cloud Run Invoker
- Logging Writer

Resolved 403 permission errors during setup.

---

# 4.3 BigQuery Data Warehouse

Dataset: `clinical_trials`

Tables:

- `stg_trials`
- `Analytics_trials`



Features:

- Partitioned by ingestion_date
- Idempotent MERGE
- Columnar storage optimization

---

# 4.4 Docker Containerization

Pipeline containerized using:

```
FROM python:3.10
WORKDIR /app
COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt
COPY . .
CMD ["python", "main.py"]
```

---

# 4.5 Artifact Registry

Used to store Docker images for Cloud Run execution.

---

# 4.6 Cloud Run (Compute Layer)

- Executes ingestion pipeline
- Serverless container execution
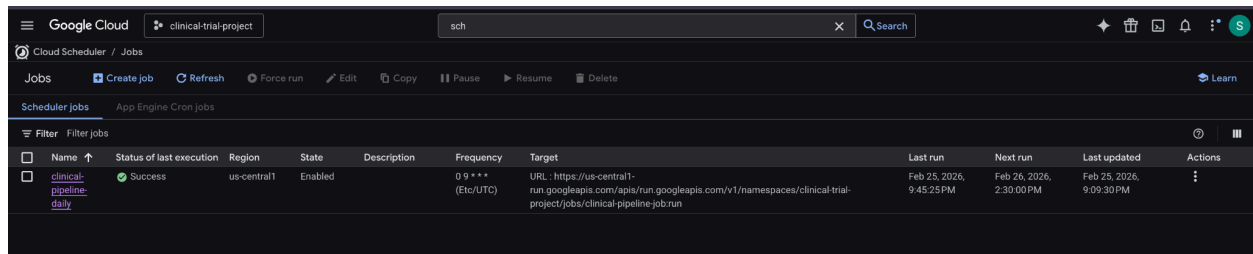- Auto-scaling
- Timeout configured
- Service account attached

# 4.7 Cloud Scheduler (Automation Layer)

- Daily trigger
- Cron: `0 9 * * *`
- HTTP trigger to Cloud Run Job

Timezone to be in UTC, 9am UTC is 1430 IST



# 4.8 Cloud Logging

Logs captured:

- Pipeline start/end
- Record counts
- Validation warnings
- Errors

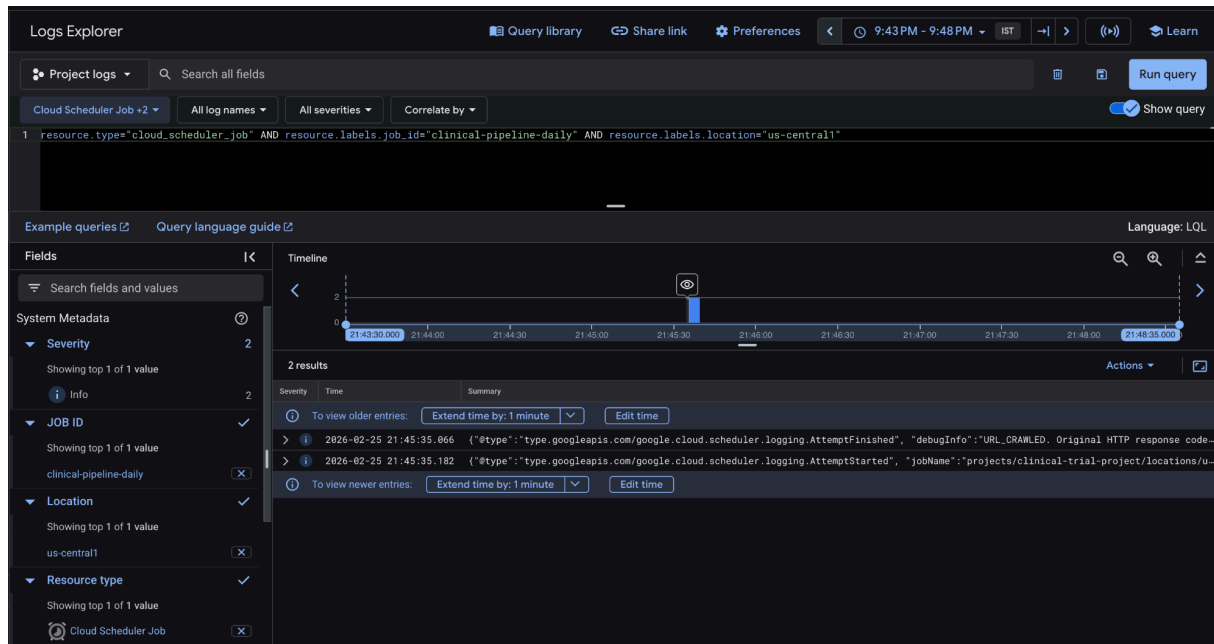# 4.9 Logging-Based Alerts

Monitoring policy created:

Query:

resource.type="cloud_run_revision"
severity>=ERROR

Triggers notification if job fails.

# 5. Data Pipeline Design

## 5.1 Ingestion Layer

- API calls
- File parsing
- Pagination handling

## 5.2 Cleaning & Normalization

- Canonical schema mapping
- Type casting
- Null handling

## 5.3 Validation Layer

Implemented in `validation.py`:

- Empty dataset detection
- Null ID warnings
- Record count checks

---

## 5.4 Idempotent Loading Strategy

BigQuery MERGE:

```
MERGE target T
USING staging S
ON T.registry_id = S.registry_id
WHEN MATCHED THEN UPDATE SET *
WHEN NOT MATCHED THEN INSERT ROW
```

Ensures safe reruns.

---

## 5.5 Partitioning Strategy

Partitioned by ingestion_date for:

- Cost reduction
- Faster scans
- Efficient incremental loads

---

# 6. Codebase Structure
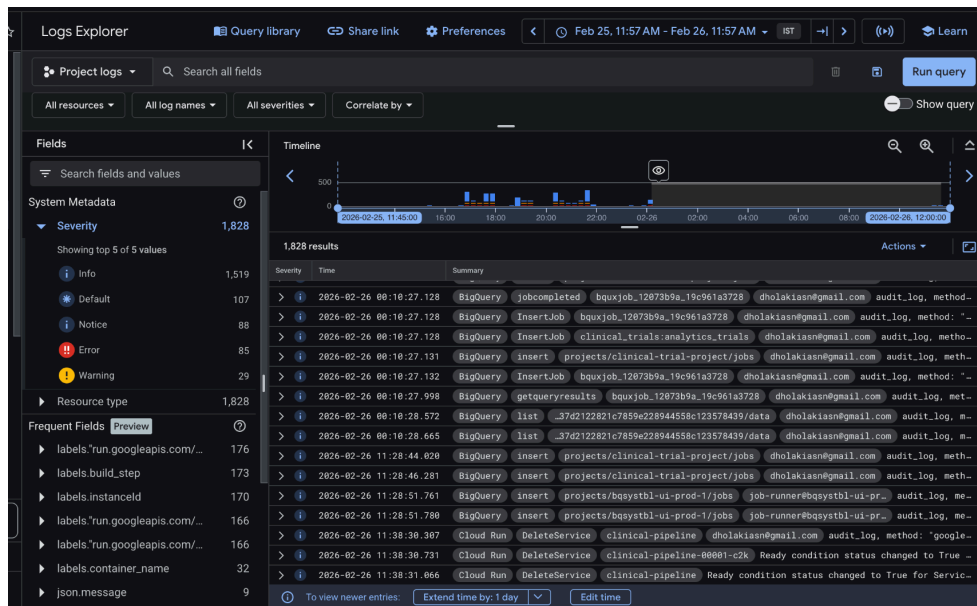
```
clinical-trial-pipeline/
│
├── main.py
├── validation.py
├── requirements.txt
├── Dockerfile
└── README.md
```
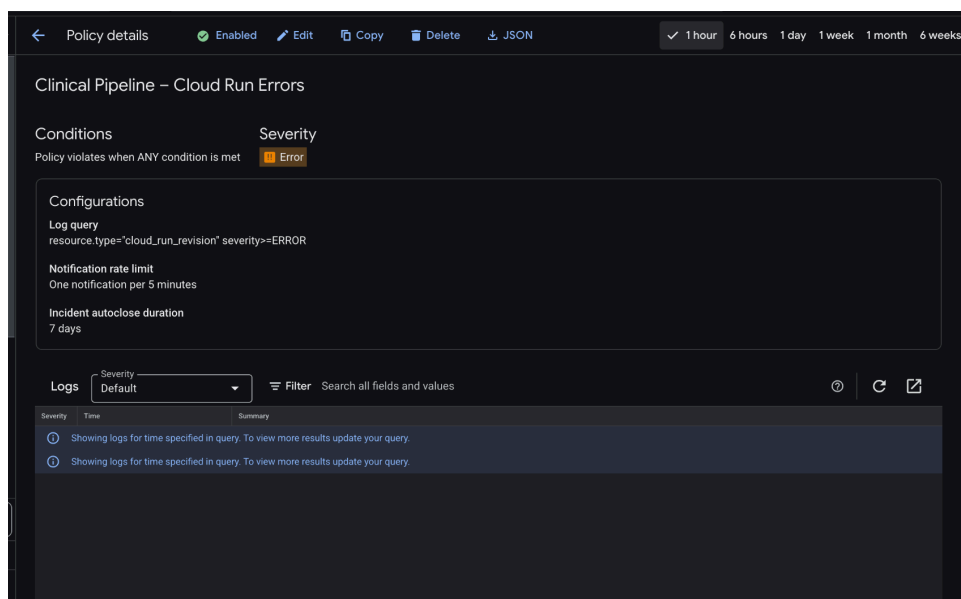
---

# 7. Monitoring & Observability

## 7.1 Log Capture

Cloud Run logs streamed to Cloud Logging.



## 7.2 Error Detection

Severity-based logging implemented.

# 7.3 Alert Policy

Log-based alert triggers on ERROR.





# 8. Analytics & Insight Demonstration

# 8.1 Sample BigQuery Queries

- Trials by year
- Trials by status
- Top sponsors
- Condition distribution

# 8.2 Looker Studio Dashboards

Visualizations:

- Top countries-
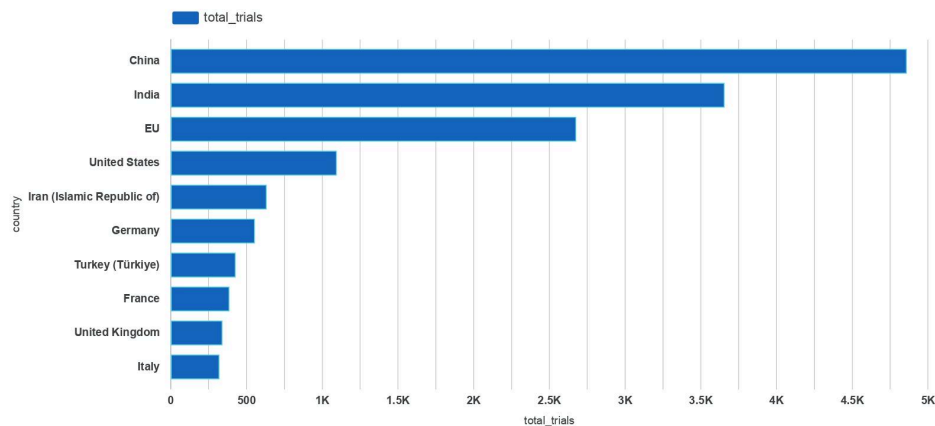  https://lookerstudio.google.com/reporting/e3ebfdf2-9e34-48f9-8b06-088484ae1464

## Top Countries by Trial Count

**Data Source: WHO ICTRP + EUCTR + ISRCTN + EMA**
**Updated: Dynamic (Cloud Run Pipeline)**
**Total Trials: 19,551**



- Top conditions-
  https://lookerstudio.google.com/reporting/cd51f9f3-8f28-4db9-a3ae-05846badaee5

# Top 10 Clinical Trial Conditions



All : 100%

Myocardial Infarction : 29.09%

O- Medical and Surgical :...

Diabetes M...

Cardioge...

Acute Myocardial Infarction : 21.17%

HIV Infections : 5.76%

Acute Coron...

Hypertensio...

ST Elevation Myocar...

Coronary Artery Disease : 4.6...

- Source distribution -
  https://lookerstudio.google.com/reporting/0c5b3228-b4f5-492d-b2d3-ea6c02a18471

# Clinical Trial Distribution by Registry Source



"It shows the distribution of integrated clinical trials by registry source. We can see the majority of the dataset comes from CTGov, CTRI, and ChiCTR, indicating strong representation from US, Indian, and Chinese registries.

# 9. Challenges Faced & Resolutions

## 9.1 Idempotency

Resolved via MERGE.

## 9.2 Scheduler Timezone Issues

Resolved via explicit timezone configuration.

## 9.3 IAM Permission Errors

Resolved via role assignment.

## 9.4 Docker Dependency Failures

Resolved via updated requirements.txt.

## 9.5 Cloud Run Startup Failures

Resolved via correct container CMD & configuration.

## 9.6 API Pagination

Resolved via nextPageToken looping.

---

# 10. Cost Optimization Strategy

- Fully serverless
- No idle compute
- Pay-per-execution model
- Partitioned warehouse
- No persistent VMs

# 11. Conclusion

This project successfully delivers:

- Multi-source ingestion
- Schema normalization
- Idempotent warehouse loading
- Automated scheduling
- Monitoring & alerting
- Scalable, serverless architecture
- Analytical dashboards

It transforms heterogeneous clinical trial data into a centralized, analytics-ready platform using modern cloud-native engineering principles.