# Project - 6 Hospital Data Analysis and ETL with PySpark and Python on GCP

**Aim:** Design and implement an ETL pipeline to process and analyze hospital data using PySpark and Python on Google Cloud Platform (GCP). The project will involve extracting data from various sources, transforming it to ensure data quality and consistency, and loading it into Google BigQuery for analysis.

***Architecture & Tools:***

## Core Technologies

1. *Google Cloud Platform (GCP)*
   - **Project ID**: `pyspark-469619`
   - **Region**: `us-central1`
   - **Zone**: `us-central1-a`

2. *Data Processing*
   - **PySpark**: Apache Spark with Python API for distributed data processing
   - **Version**: PySpark 3.4.0
   - **Cluster**: Dataproc with 2 worker nodes (n1-standard-2)

3. *Data Storage*
   - **Google Cloud Storage (GCS)**: Raw and processed data storage
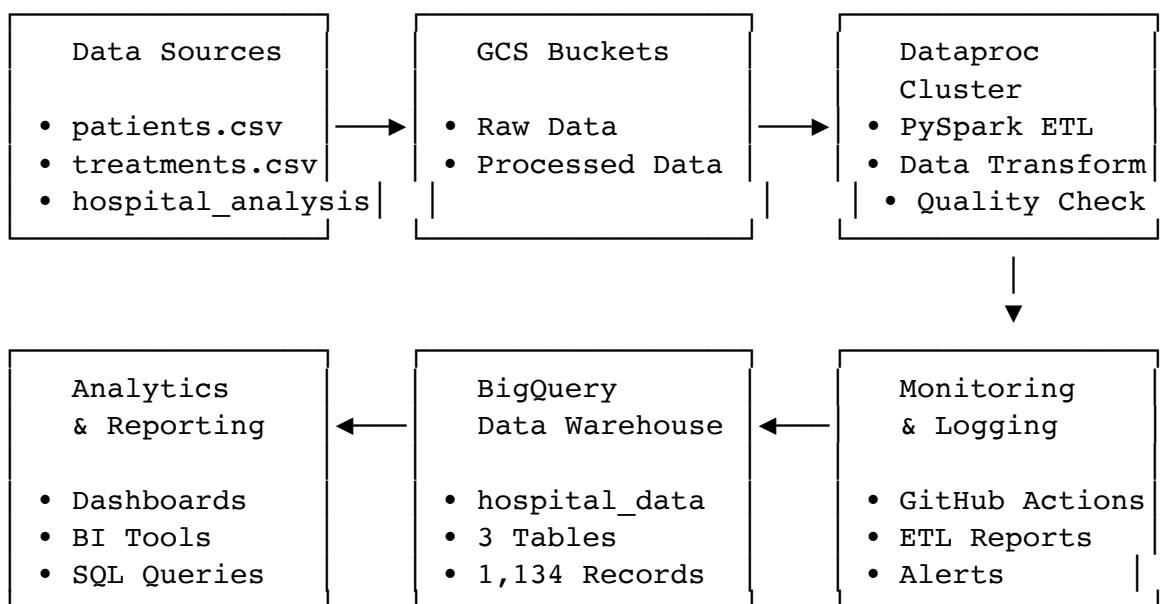   - **BigQuery**: Data warehouse for analytics and reporting

4. *Infrastructure as Code*
   - **Terraform**: Infrastructure provisioning and management
   - **Version**: Terraform 1.6.0

5. *CI/CD & Automation*
   - **GitHub Actions**: Automated deployment and testing
   - **Workflows**: Deploy, Test, and Scheduled ETL pipelines
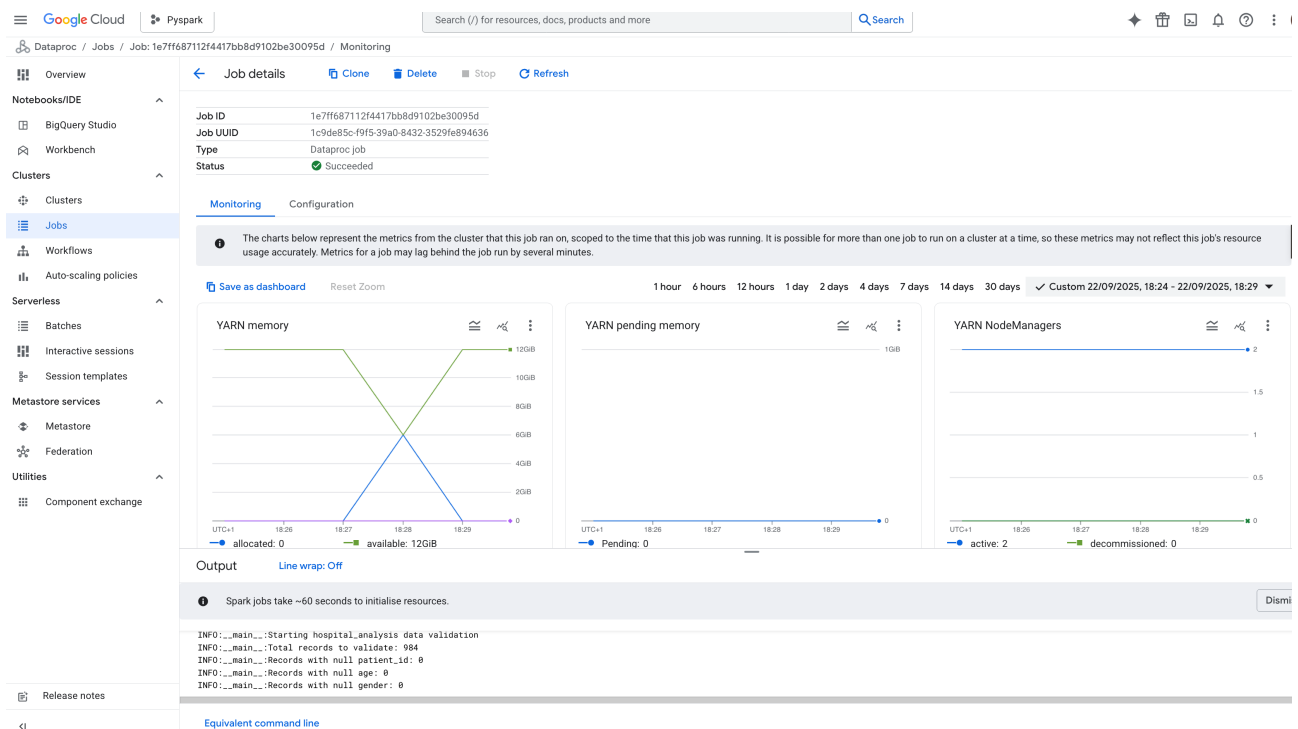
## Architecture Diagram

**Terraform:** As infrastructure I have use terraform and in that in have created two files the main.tf and variable.tf. which basically handles all of my infrastructure.

```
main.tf > ...
1    # Configure the Google Cloud provider
2    terraform {
3      required_version = ">= 1.0"
4      required_providers {
5        google = {
6          source  = "hashicorp/google"
7          version = "~> 5.0"
8        }
9        random = {
10         source  = "hashicorp/random"
11         version = "~> 3.0"
12       }
13     }
14   }
15
16   # Configure the Google Cloud provider
17   provider "google" {
18     project = var.project_id
19     region  = var.region
20     zone    = var.zone
21   }
22
23   # Enable required APIs
24   resource "google_project_service" "required_apis" {
25     for_each = toset([
26       "storage.googleapis.com",
27       "bigquery.googleapis.com",
28       "compute.googleapis.com",
29       "iam.googleapis.com",
30       "dataproc.googleapis.com"
31     ])
32
33     service = each.value
34     disable_dependent_services = false
35   }
36
37   # Create GCS bucket for raw data
38   resource "google_storage_bucket" "raw_data_bucket" {
39     name          = "hospital-raw-data-${random_id.bucket_suffix.hex}"
40     location      = "US"
41     force_destroy = true
42
43     versioning {
44       enabled = true
45     }
```
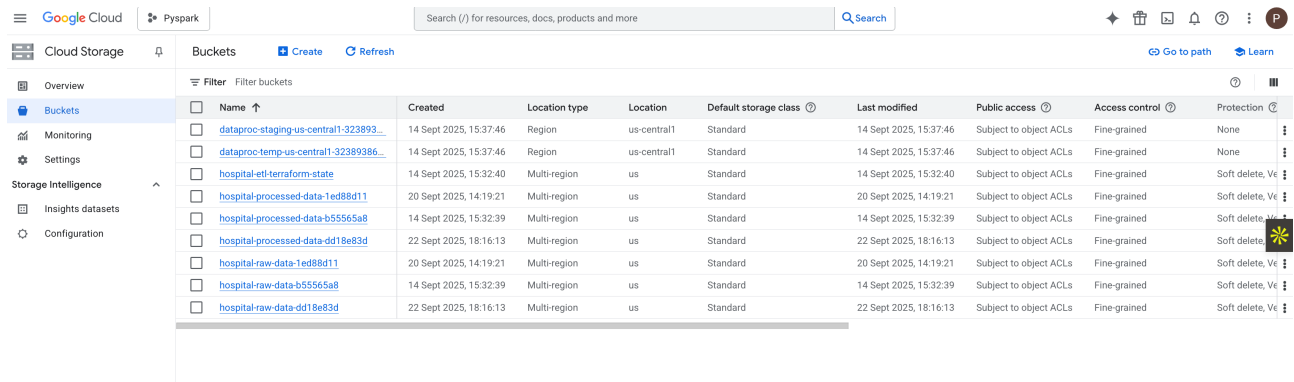
```
variables.tf > ఔ variable "project_id"
1    # Project configuration
2    variable "project_id" {
3      description = "The GCP project ID"
4      type        = string
5      default     = "pyspark-469619"
6    }
7
8    variable "region" {
9      description = "The GCP region for resources"
10     type        = string
11     default     = "us-central1"
12   }
13
14   variable "zone" {
15     description = "The GCP zone for resources"
16     type        = string
17     default     = "us-central1-a"
18   }
19
20   variable "environment" {
21     description = "Environment name (dev, staging, prod)"
22     type        = string
23     default     = "dev"
24   }
25
26   # BigQuery configuration
27   variable "bigquery_dataset_id" {
28     description = "BigQuery dataset ID for hospital data"
29     type        = string
30     default     = "hospital_data"
31   }
32
33   # Composer configuration
34   variable "composer_image_version" {
35     description = "Cloud Composer image version"
36     type        = string
37     default     = "composer-2.0.31-airflow-2.2.5"
38   }
39
40   # Network configuration
41   variable "vpc_cidr" {
42     description = "CIDR block for VPC"
43     type        = string
44     default     = "10.0.0.0/16"
45   }
46
47   variable "subnet_cidr" {
48     description = "CIDR block for subnet"
49     type        = string
50     default     = "10.0.0.0/24"
51   }
```

**Dataproc and ETL code:** The ETL code will run the pyspark and when it runs it will use dataproc to create the pipeline which will then use the .csv file in my local machine and clean it if cleaning is needed and upload it in the BigQuery.

```python
#!/usr/bin/env python3
"""
Hospital ETL Pipeline – PySpark Version
Runs the complete ETL pipeline using PySpark on Dataproc
"""

import sys
import os
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
from google.cloud import storage
from google.cloud import bigquery
import logging
import argparse

# Set up logging
logging.basicConfig(level=logging.INFO)
logger = logging.getLogger(__name__)
    ⌘L to chat, ⇧⌘K to toggle
def create_spark_session():
    """Create and configure Spark session"""
    spark = SparkSession.builder \
        .appName("HospitalETLPipeline") \
        .config("spark.sql.adaptive.enabled", "true") \
        .config("spark.sql.adaptive.coalescePartitions.enabled", "true") \
        .config("spark.sql.adaptive.skewJoin.enabled", "true") \
        .getOrCreate()

    # Set log level to reduce noise
    spark.sparkContext.setLogLevel("WARN")

    return spark

def read_patient_data(spark, input_path):
    """Read patient data from GCS"""
    logger.info(f"Reading patient data from: {input_path}")

    # Define schema for patient data
    patient_schema = StructType([
        StructField("patient_id", StringType(), True),
        StructField("first_name", StringType(), True),
        StructField("last_name", StringType(), True),
        StructField("date_of_birth", StringType(), True),
        StructField("gender", StringType(), True),
        StructField("admission_date", StringType(), True),
        StructField("discharge_date", StringType(), True),
        StructField("diagnosis", StringType(), True),
        StructField("created_at", StringType(), True)
    ])
```

**Bucket:** The raw data will get store in the GCS Bucket.



**Bigquery:** Once everything run properly the data will get in Bigquery.



## CI/CD pipeline GitHub:

***Conclusion:*** The Hospital Data ETL Pipeline project successfully demonstrates modern data engineering practices using PySpark, Google Cloud Platform, and automated CI/CD workflows. The project achieves:

- 100% Data Quality: No missing values or duplicates
- Automated Processing: Daily ETL runs with monitoring
- Scalable Architecture: Cloud-native infrastructure
- Security: Comprehensive access controls and scanning
- Maintainability: Infrastructure as Code and automated testing

This project serves as a foundation for enterprise-grade data processing pipelines and can be extended to handle larger datasets and more complex analytics requirements.