MAESTRÍA EN INTELIGENCIA ARTIFICIAL Y CIENCIA DE DATOS

CASO DE EXTRACCION TRANSFORMACION Y CARGA DE DATOS

Yoniliman Galvis Aguirre 22500214



CASO

Definición: Dataset de base de datos industrial en formato csv ubicado en github

Trata de la producción de 8 trenes, 51 variantes 516 días, 14876 lotes y 236530 registros

Se usa archivos yaml para cargar certificados.

Se toma el csv original a un pandas df y se preprocesa para buscar los datos hurfanos, como es inferior a 2% se eliminan

Se crea la base de datos para ETL en postgresql

Se crea tabla para el dataset

Se hace un vaceado a la tabla creada

Usando METABASE se conecta a posgresql Se leen las tablas

Se crean las preguntas (consultas)

Se crean dashboards

Obtener Datos (ingesta):

Transformar datos

Dataset:

Extraer datos:

Carga

Se carga configuracion de archivo yaml

Se eliminan archivos dde datos mas viejos que 7 dias Se descarga el achivo .raw desde github, se sube a un pandas df.

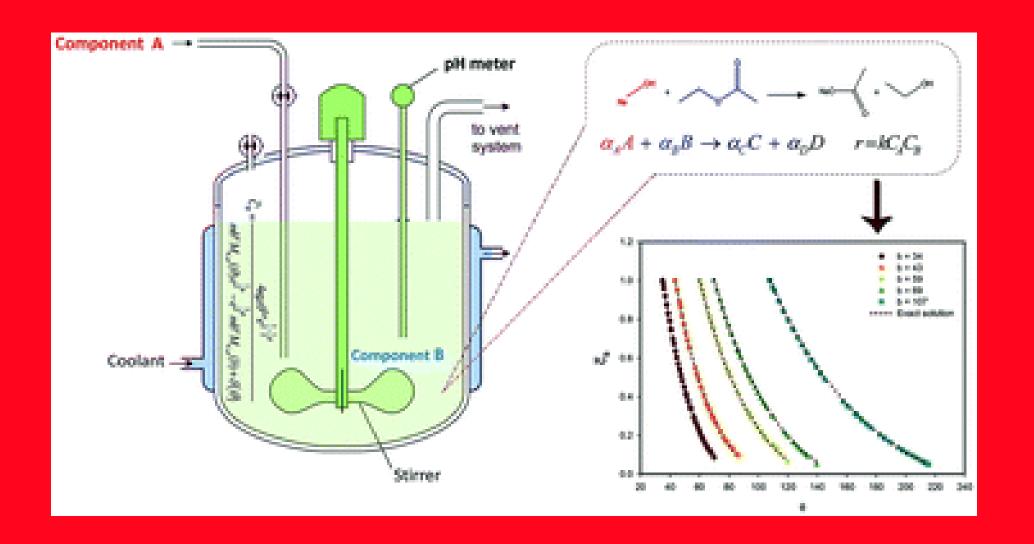
Se guarda en csv en la carpeta de data

Se crea un pandas profiling para analizar los datos.

Se usa archivos yaml para cargar certificados.

crea tablas de staging usando un archivo .sql en la base de datos

Se preprocesan, organizan y preprocesan los datos, se populan las tablas de stagin



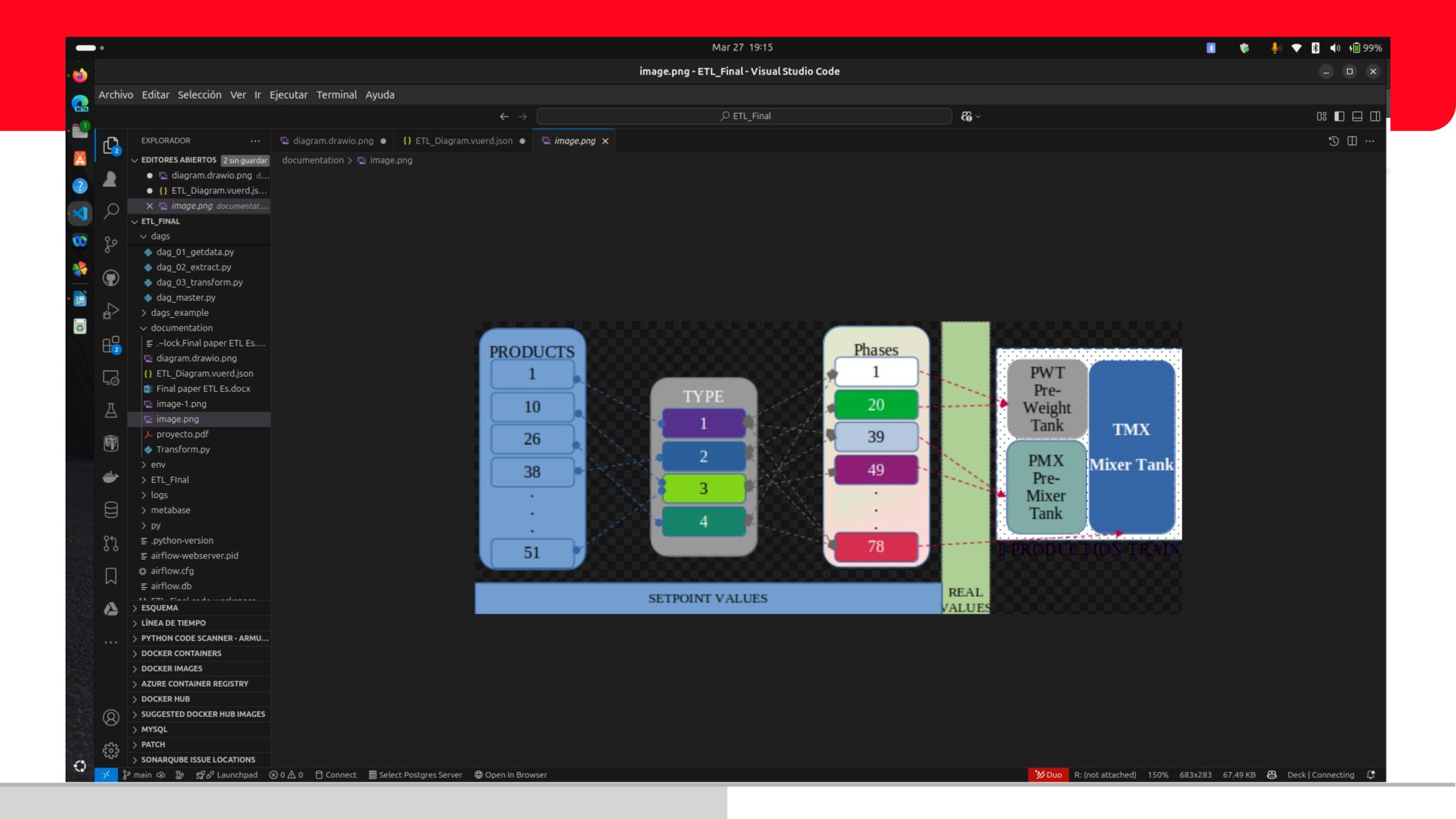


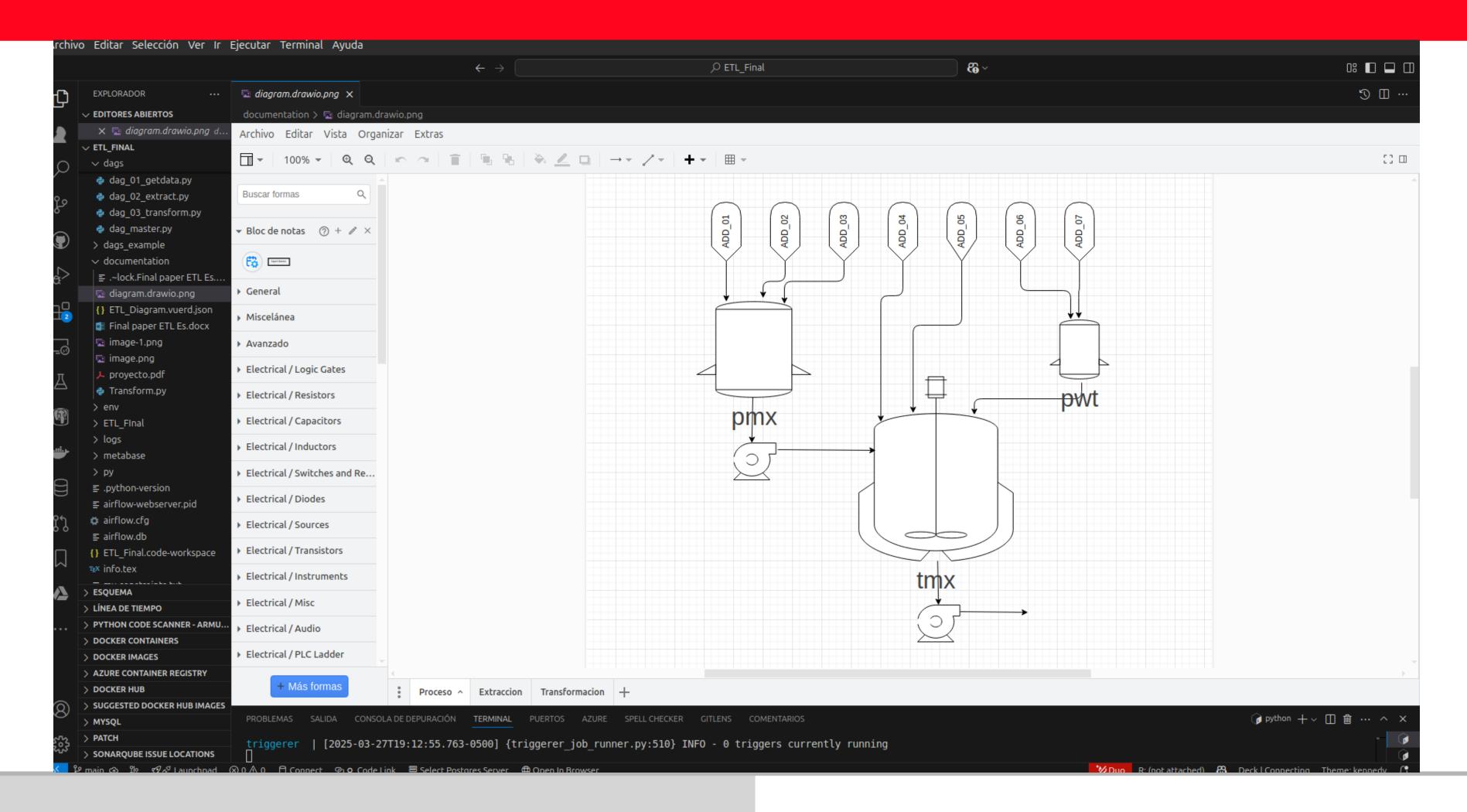
Los procesos de fabricación por lotes son sistemas con una variabilidad compleja y, especialmente en procesos de múltiples productos, garantizar la calidad del producto requiere pruebas de calidad constantes que demandan tiempo de producción, costos laborales y el uso de reactivos químicos que producen desechos que a menudo requieren un tratamiento costoso antes de que sus residuos puedan ser liberados al medio ambiente.

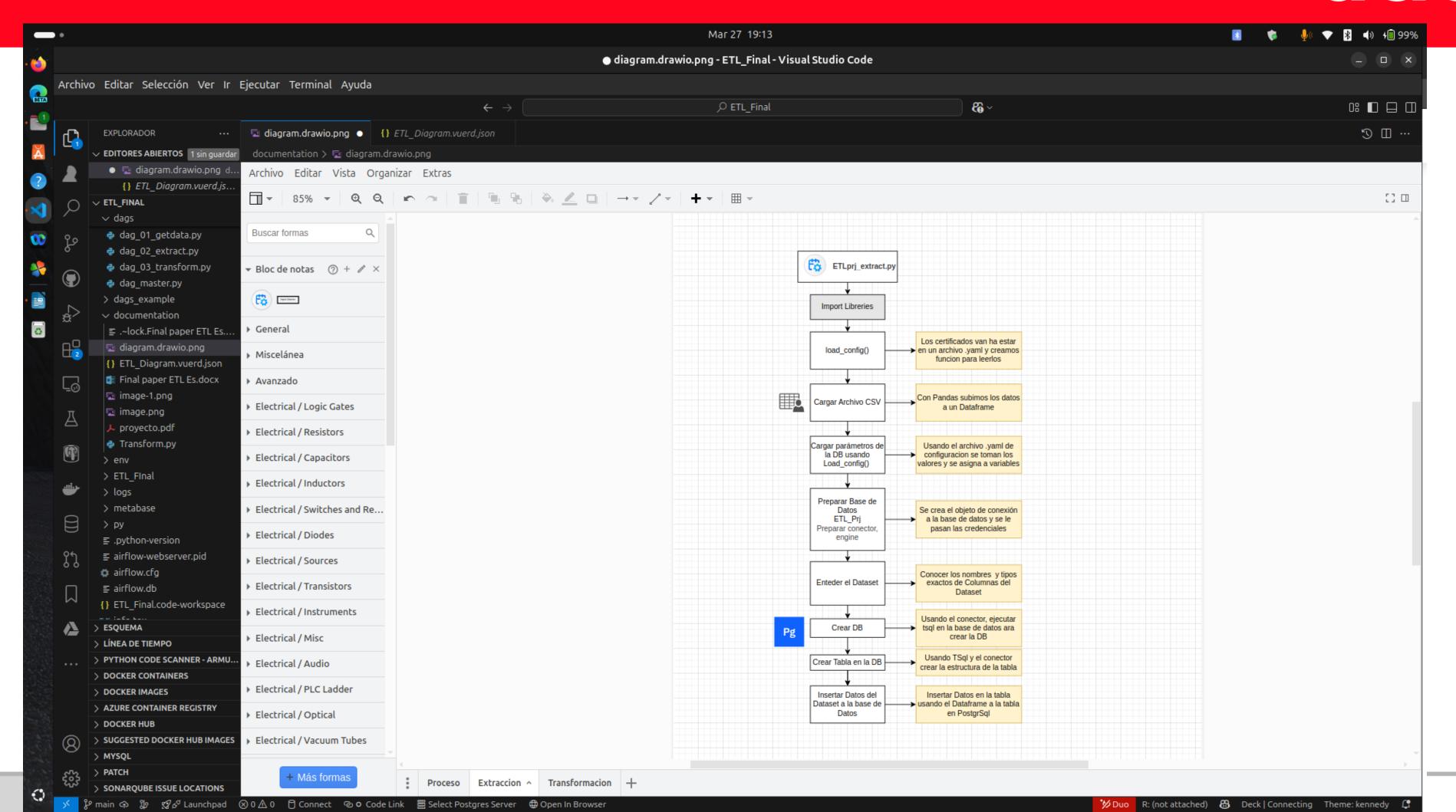
Está más que justificado utilizar herramientas de ETL, analítica de datos, inteligencia artificial y aprendizaje automático para desarrollar nuevas herramientas que mejoren los procesos de fabricación actuales con menores costos, menos horas de trabajo y un menor impacto en el medio ambiente, apuntando a un proceso sin residuos.

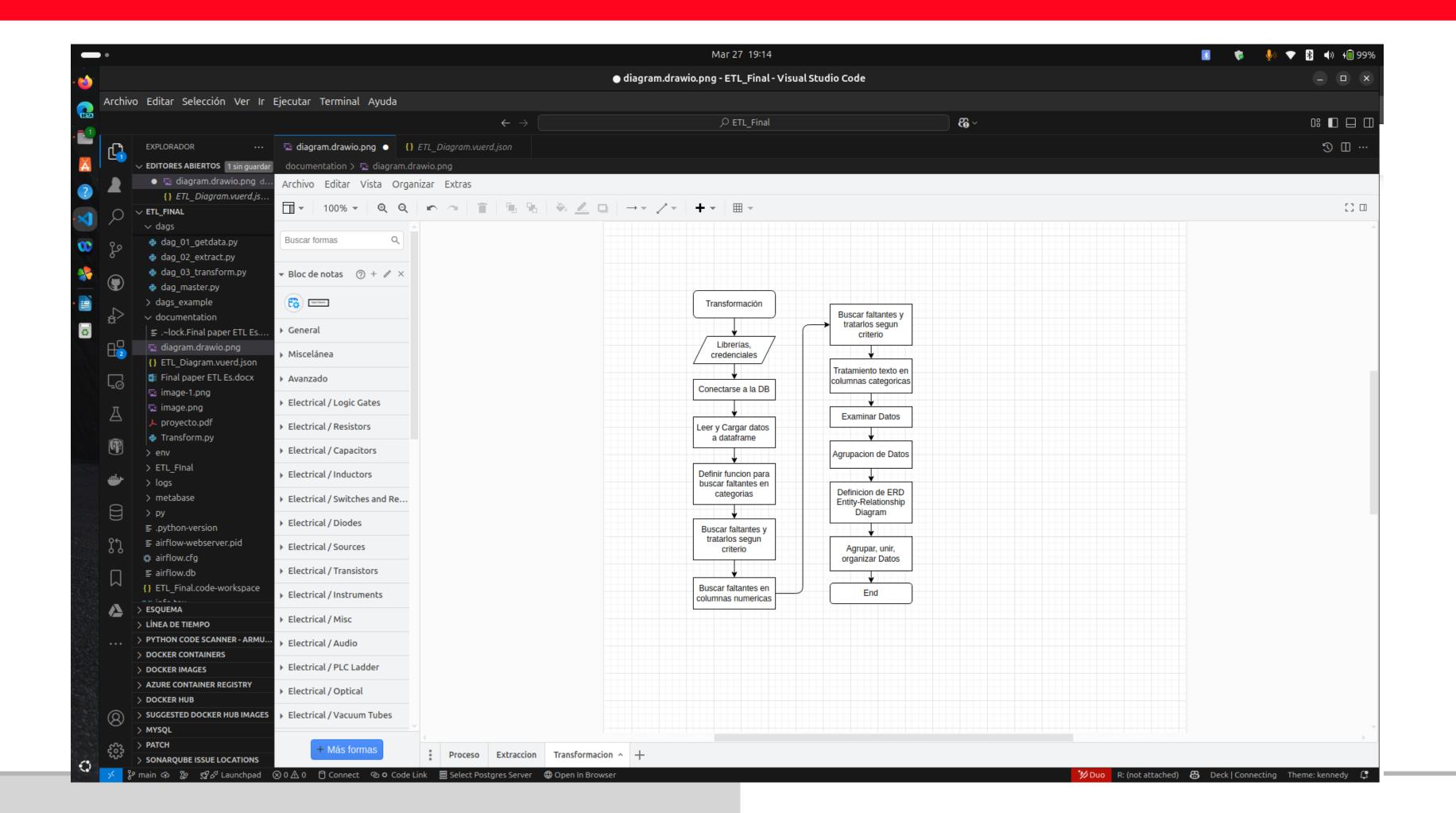
Los modelos ETL tambien permite desarrollar Dashboards que permitan analizar el comportamiento de los procesos para poder mejorarlos y reducir los fuera de especificaciones y tambien desarrollar mejores planes de produccion que disminuyan la cantidad de recursos necesarios para la fabricación.

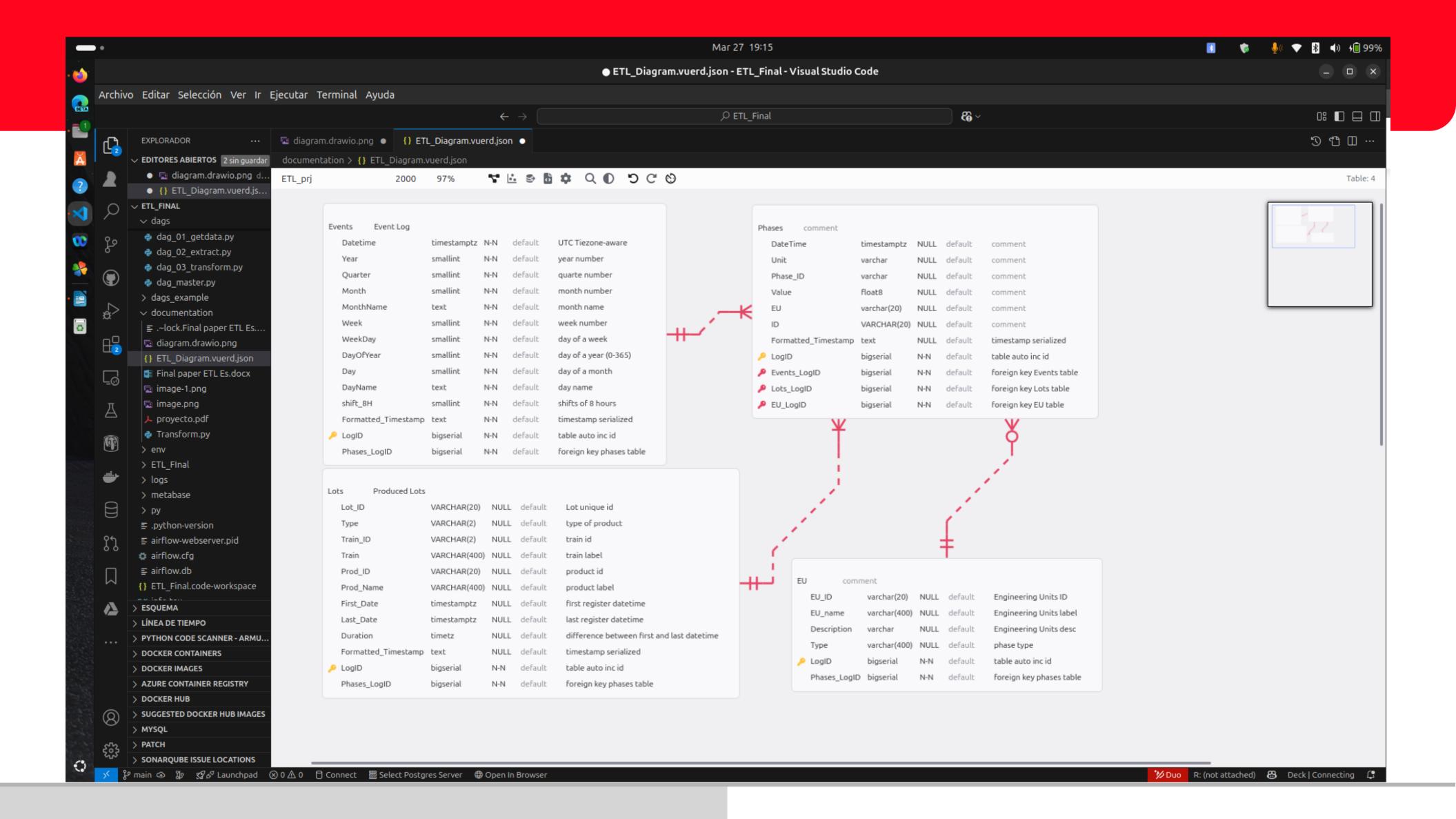
DATASET udo

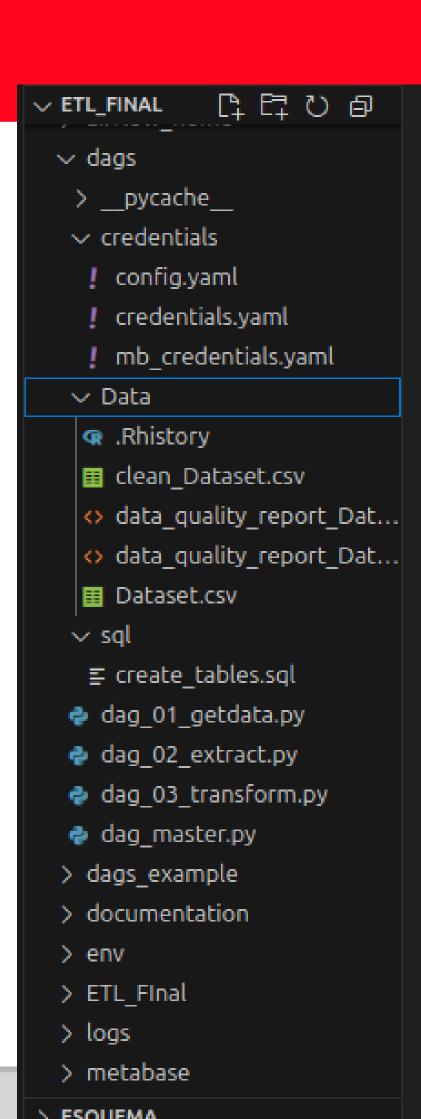












```
create_tables.sql-
Archivo Editar Selección Ver Ir Ejecutar Terminal Ayuda

✓ EDITORES ABIERTOS 2 sin guardar

                                dags > sql > = create tables.sql
                                       CREATE TABLE IF NOT EXISTS EU

    tiagram.drawio.png d...

        • {} ETL_Diagram.vuerd.js...
                                         EU ID
                                                       varchar(20) UNIQUE,
        EU name
                                                       varchar(400),
     ▽ETL_FINAL [ユロット]
                                         Description varchar(400),
                                         Type
                                                       varchar(400),
      dags
                                         LogID
                                                       bigserial
                                                                    NOT NULL,
        > __pycache__
                                         Phases LogID bigserial
                                                                    NOT NULL,
        credentials
                                         PRIMARY KEY (LogID)
         ! config.yaml
        ! credentials.yaml
         ! mb_credentials.yaml
                                       COMMENT ON COLUMN EU.EU ID IS 'Engineering Units ID';
                                       COMMENT ON COLUMN EU.EU name IS 'Engineering Units label';

∨ Data

                                       COMMENT ON COLUMN EU.Description IS 'Engineering Units desc';
        • Rhistory
                                       COMMENT ON COLUMN EU. Type IS 'phase type';
        clean_Dataset.csv
                                       COMMENT ON COLUMN EU.LogID IS 'table auto inc id';
        data_quality_report_Dat..
                                       COMMENT ON COLUMN EU. Phases LogID IS 'foreign key phases table';
        data_quality_report_Dat...
        ■ Dataset.csv
                                       CREATE TABLE IF NOT EXISTS Events

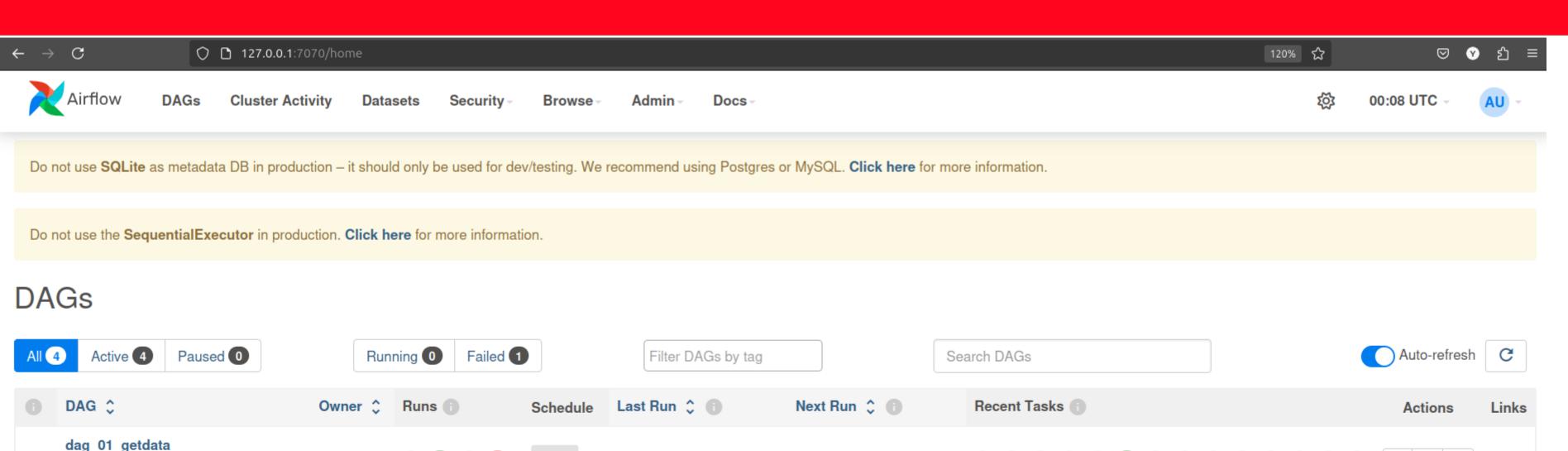
√ sql

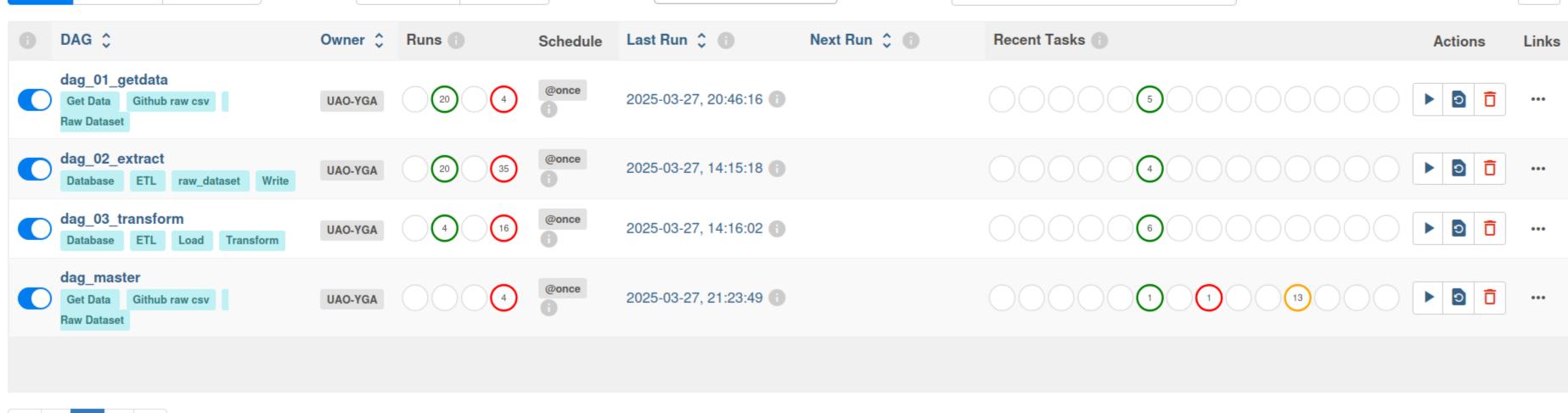
                                         Datetime
                                                              timestamptz NOT NULL UNIQUE,

≡ create_tables.sql

                                 22
                                         Year
                                                              smallint NOT NULL,
       dag_01_getdata.py
                                         Quarter
                                                              smallint
                                                                          NOT NULL,
       dag_02_extract.py
                                                              smallint
                                         Month
                                                                          NOT NULL,
       dag_03_transform.py
                                         MonthName
                                                                           NOT NULL,
       dag_master.py
                                                              smallint
                                         Week
                                                                          NOT NULL,
       > dags_example
                                         WeekDay
                                                              smallint
                                                                          NOT NULL,
                                         DayOfYear
                                                              smallint
                                                                          NOT NULL,
       > documentation
                                                              smallint
                                         Day
                                                                          NOT NULL,
      > env
                                         DayName
                                                                           NOT NULL,
       > ETL_Final
                                                              smallint
                                         Hour
                                                                          NOT NULL,
       > logs
                                         Minute
                                                              smallint
                                                                          NOT NULL,
       > metabase
                                         Second
                                                              smallint
                                                                          NOT NULL,
                                                              smallint
     > ESQUEMA
                                         shift 8H
                                                                          NOT NULL,
                                         Formatted Timestamp text
                                                                           NOT NULL,
      > LÍNEA DE TIEMPO
                                         LogID
                                                                          NOT NULL,
                                                              bigserial
     > PYTHON CODE SCANNER - ARMU...
                                         Phases LogID
                                                              bigserial NOT NULL,
      > DOCKER CONTAINERS
                                         PRIMARY KEY (LogID)
     > DOCKER IMAGES
     > AZURE CONTAINER REGISTRY
     > DOCKER HUB
                                       COMMENT ON TABLE Events IS 'Event Log';
                                       COMMENT ON COLUMN Events.Datetime IS 'UTC Tiezone-aware';
     > SUGGESTED DOCKER HUB IMAGES
                                       COMMENT ON COLUMN Events. Year IS 'year number';
     > MYSQL
     > PATCH
                                       COMMENT ON COLUMN Events. Month IS 'month number';
      > SONARQUBE ISSUE LOCATIONS
                                       COMMENT ON COLUMN Events.MonthName IS 'month name':
```







Showing 1-4 of 4 DAGs

