# San Francisco 311 Requests

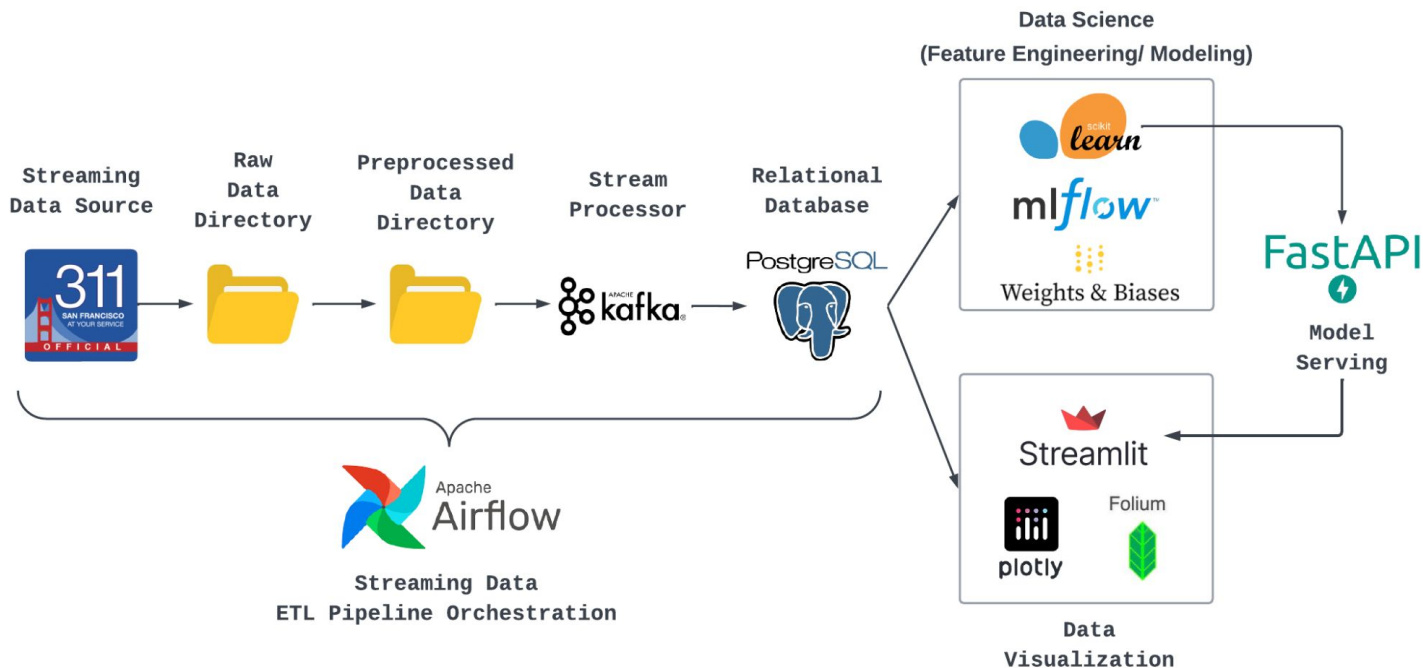# Data Streaming ETL Pipeline

# Project Overview

**Summary**  **San Francisco** faces the challenge of efficiently managing and addressing the diverse array of incidents reported through its **311 system**. With a continuous influx of cases, there is a pressing need to **streamline the allocation of resources** and **enhance response times to improve service quality**.

**Objectives**

1. **Streaming Data ETL Pipeline:** Leveraging skills from Data Streaming course to build a real-time streaming data pipeline.

2. **Data Visualization:** Mapping 311 cases across the city to identify hotspots and patterns.

3. **Closure Time Prediction:** Using machine learning to predict case closure times, considering factors like case type and location.

# Architecture Diagram



Streaming Data Source → Raw Data Directory → Preprocessed Data Directory → Stream Processor → Relational Database

Streaming Data
ETL Pipeline Orchestration

Data Science
(Feature Engineering/ Modeling)

Model Serving

Data Visualization

# Dataset

- Open source dataset managed by San Francisco Government

- **Dataset Creation:** October, 2011

- **Update frequency:** Daily (multiple times per hour)

- **Dataset Size:** 2.2 GB

- **Number of Records:** 6.6 M

- **Features:** 48 (interpretable: 15)



Source: https://data.sfgov.org/City-Infrastructure/311-Cases/vw6y-z8j6/explore

# Tools & Packages

- **Data Orchestration:** apache-airflow, sodapy

- **Data Streaming:** avro, confluent_kafka, fastavro

- **Database:** pgAdmin, psycopg2, sqlalchemy

- **Modeling & Serving:** fastapi[all], mlflow, scikit-learn

- **Data Visualization:** geopandas, folium, numpy, pandas, plotly, streamlit

- **Version Control:** git, Github

# Airflow - ETL Pipeline Orchestration

```
∨ airflow
  ∨ dags
    > __pycache__
    ∨ scripts
      > __pycache__
      🐍 __init__.py
      ⚙ .env
      🐍 consumer.py
      🐍 entities.py
      🐍 helper_functions.py
      🐍 producer.py
      🐍 schemas.py
    🐍 data_ingestion_dag.py
    🐍 data_loading_dag.py
    🐍 operator_config.py
  > logs
  ≡ airflow-webserver.pid
  ⚙ airflow.cfg
  🐍 webserver_config.py
```

**DAG 1**

```python
with DAG(
    dag_id="sf_311_data_ingestion",
    schedule='15 0 * * *',
    description='Fetch, preprocess, and produce data to Kafka cluster',
    default_args=default_args) as dag:

create_raw_data_dirs_ops >> fetch_data_ops >> save_raw_data_ops
save_raw_data_ops >> extract_cols_ops >> create_processed_data_dirs_ops >> save_processed_data_ops
save_processed_data_ops >> produce_data_ops
```

**DAG 2**

```python
with DAG(
    'sf_311_data_loading',
    default_args=default_args,
    description='Consume data from Kafka and upsert to PostgreSQL Database',
    schedule_interval='30 0 * * *',
):
consume_data_ops
```

# Airflow - ETL Pipeline Orchestration



**DAG 1 - Data Ingestion, Preprocessing, Produce**

**Schedule: every night at 00:15 AM**

**DAG 2 - Consume to Database**

**Schedule: every night at 00:30 AM**

# Kafka - Data Streaming

## Kafka Topic



## Kafka Schema Registry - Avro Serializer

# PostgreSQL Database

## SQL Table Creation

```
CREATE TABLE IF NOT EXISTS historical_311_request
(
    request_id character varying PRIMARY KEY NOT NULL,
    requested_datetime character varying NOT NULL,
    updated_datetime character varying NOT NULL,
    status_description character varying,
    agency_responsible character varying,
    service_type character varying,
    service_subtype character varying,
    address character varying,
    street character varying,
    supervisor_district character varying,
    neighborhood character varying,
    police_district character varying,
    latitude double precision,
    longitude double precision,
    source character varying
);
```

## pgAdmin - PostgreSQL Tool

# Scikit-Learn/ FastAPI - Modeling & Serving

## Scikit learn

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import accuracy_score, confusion_matrix
import joblib

# Splitting data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Model Training and Hyperparameter Tuning with Random Forest
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 5, 10, 15]
}

rf = RandomForestClassifier(random_state=42)
grid_search = GridSearchCV(estimator=rf, param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=-1)
grid_search.fit(X_train, y_train)

best_params = grid_search.best_params_

# Training the model with best parameters
best_rf = RandomForestClassifier(**best_params, random_state=42)
best_rf.fit(X_train, y_train)

# Predictions on test set
y_pred = best_rf.predict(X_test)

# Metrics Calculation
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

print(f"Accuracy: {accuracy}")
```

```
Accuracy: 0.586056549770947
```

## FastAPI

# Streamlit/ Plotly/ Folium - Visualization



**Page 1 -**

**Real Time**

**Data Visualization**

**Dashboard**



**Page 2 -**

**Request Time**

**Prediction**