**ITCS-6190**
**Cloud Computing for Data Analytics**
**LLM - Detect AI Generated Text**

**Team Members**

Apoorve Bhargava
Soumik Paul
Juhi Jadhav
Raj Kumar
Sai Kumar Reddy, Bommareddy

**Dataset: [LLM - Detect AI Generated Text](#)**

**AWS Services Used**
AWS S3
AWS Athena
AWS Glue
AWS Quicksight
AWS Sagemaker

**About the data**

The LLM (Large Language Model) Detect AI Generated Text dataset is likely designed for training models to identify and differentiate between human-generated and AI-generated text. It encompasses a diverse collection of text samples, featuring instances of both human and machine-generated content. This dataset aims to enhance the capabilities of language models in discerning the unique patterns and characteristics associated with AI-generated text, aiding in the development of effective detection algorithms. Researchers and developers use such datasets to advance the field of natural language processing by creating models that can reliably identify content produced by sophisticated language models like GPT (Generative Pre-trained Transformer).

There are 4 CSV files in the dataset with respective columns

1. **Train_Essays.csv:  Id, prompt_id, text, generated**
2. **Train_Prompts.csv: prompt_id, prompt_name, instructions, source_text**
3. **Test_essays.csv : id, prompt_id, text**
4. **Sample_submission.csv : id, generated**

# Creation of Bucket and storing the dataset into the bucket

# Data Preparation using AWS Glue

## Step 1: AWS Glue is used to fetch the data from S3 to AWS Glue

ⓘ You can now create Apache Iceberg tables in the AWS Glue Data Catalog. To learn more, visit the documentation ↗  Create table  ✕

AWS Glue > Tables

## Tables

A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition.

### Tables (4)

View and manage all available tables.

Last updated (UTC) November 19, 2023 at 20:41:00  ↻  Delete  Add tables using crawler  Add table

Q Filter tables

| ☐ | Name | ▲ | Database | ▼ | Location | ▼ | Classification | ▼ | Deprecated | ▼ | View data | Data quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | sample_submission_csv | | group7db | | s3://mybucketgroup7/sample | | CSV | | - | | Table data | View data quality |
| ☐ | test_essays_csv | | group7db | | s3://mybucketgroup7/test_es | | CSV | | - | | Table data | View data quality |
| ☐ | train_essays_csv | | group7db | | s3://mybucketgroup7/train_e | | CSV | | - | | Table data | View data quality |
| ☐ | train_prompts_csv | | group7db | | s3://mybucketgroup7/train_p | | CSV | | - | | Table data | View data quality |

Notebooks
Job run monitoring
**Data Catalog tables**
Data connections
Workflows (orchestration)
▼ Data Catalog
Databases
**Tables**
Stream schema registries
Schemas
Connections
Crawlers
Classifiers
Catalog settings
▼ Data Integration and ETL
ETL jobs
Visual ETL
Notebooks
Job run monitoring
Interactive Sessions
Data classification tools
Sensitive data detection
Record Matching
Triggers
Workflows (orchestration)
Blueprints
Security configurations
▶ Legacy pages

What's New ↗
Documentation ↗
AWS Marketplace

🔵 Enable compact mode
🔵 Enable new navigation

# Data Transformation using Visual ETL

☰ Untitled job ✎  ⚠ Job has not been saved  🔵 Try new UI  Actions ▾  Save  Run

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

**Data source properties - S3**

Name
Amazon S3

S3 source type  Info
◉ S3 location
Choose a file or folder in an S3 bucket.
○ Data Catalog table

S3 URL
Q s3://mybucketgroup7/train_essays.csv  ✕  View ↗  Browse S3
☑ Recursive
Read files in all subdirectories.

Data format
CSV

Delimiter
Comma (,)

Escape character - optional
Enter a character to use for escaping
The character which immediately follows is used as-is, except for a small set of well-known escapes (\n, \r, \t, and \0)

Quote character
Double quote (")

☑ First line of source file contains column headers
☐ Records in source files can span multiple lines
ⓘ Infer schema

▶ Additional options

Data source - S3 bucket: Amazon S3
Transform - SQL Query: SQL Query
Transform - SQL Query: SQL Query
Transform - SQL Query: SQL Query
Transform - SQL Query: SQL Query

**Data preview** (200)  Info  READY ⓘ  ↻  End session  Previewing 4 of 4 fields

Q Filter sample dataset

| id | ▼ | prompt_id | ▼ | text | ▼ | generated |
|---|---|---|---|---|---|---|
| | | | | Cars. Cars have around since they became famous in the 1900s, when Henry Ford created and built the first ModelT. Cars have played a major role in our every day lives since then. But now, people are starting to question if limiting c | | |

## Step 2: Preview the data on AWS Athena
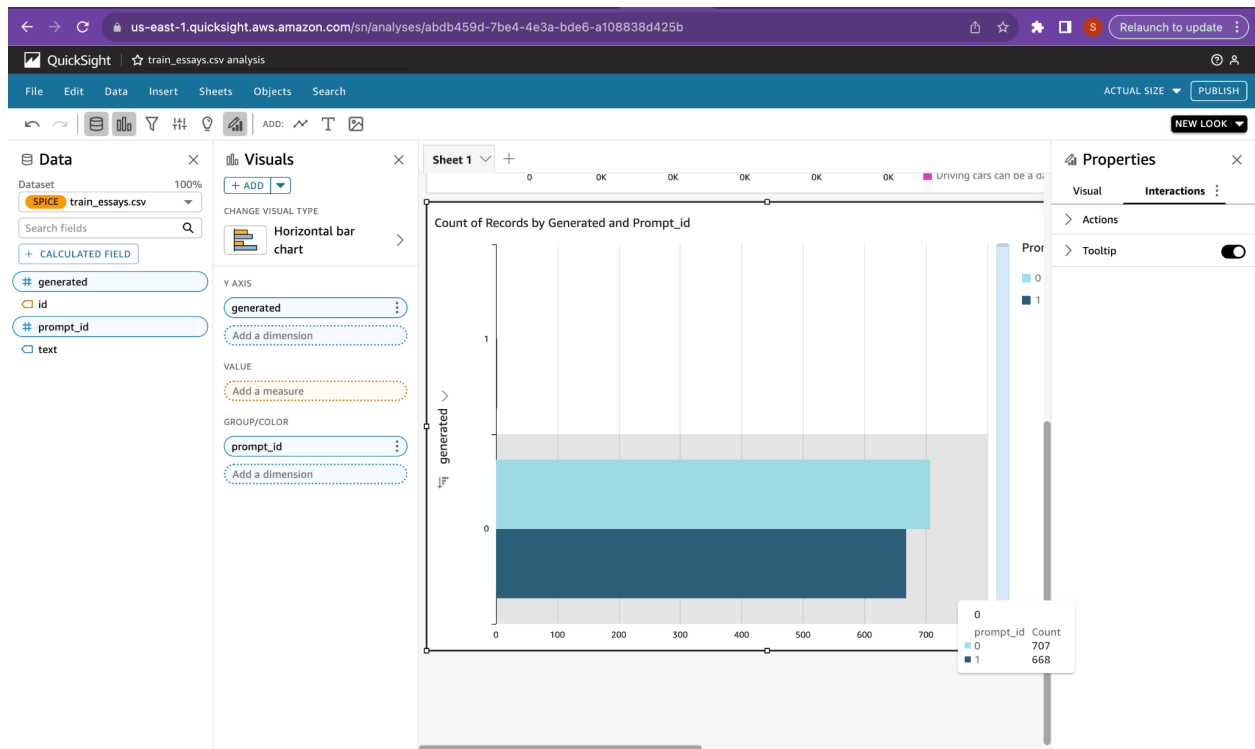Run a few SQL queries on AWS Athena to understand the datasets.

# Data Visualization on Quick Sight

The dataset has been visualized on AWS Quick Sight to get a better understanding of the data
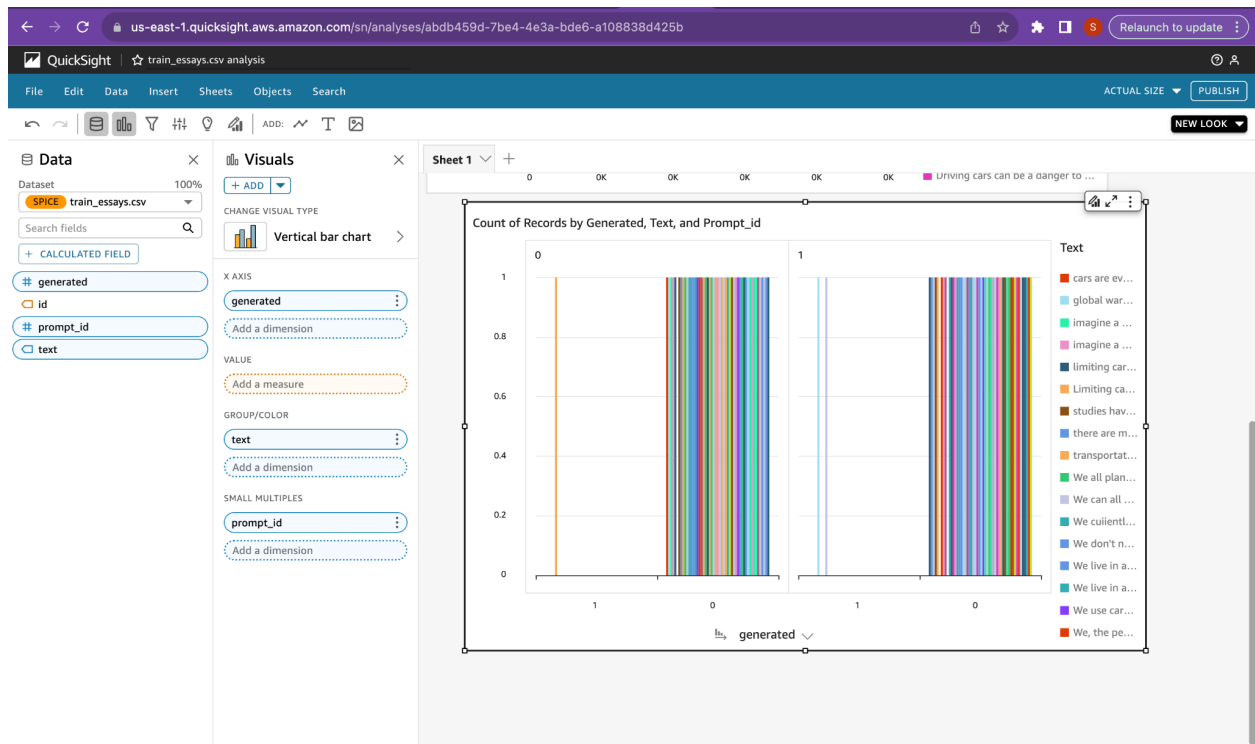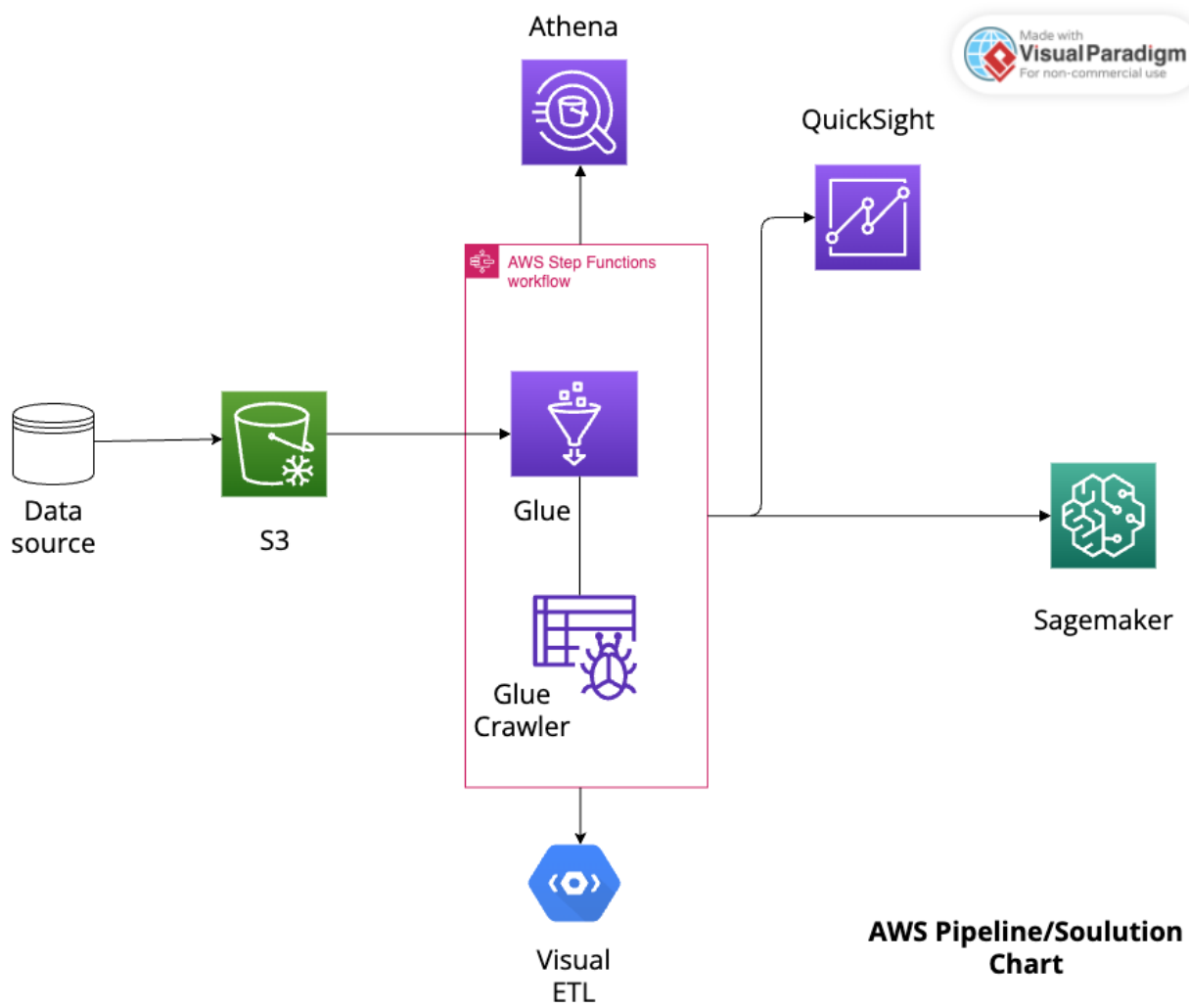
## Count of Records by IDs and Text

# Count of Records by Generated and Prompt_Ids



# Count of Records by Generated Text, and Prompt_Ids

**AWS Pipeline Solution Chart**

Athena

QuickSight

AWS Step Functions
workflow

Glue

Data
source

S3

Glue
Crawler

Sagemaker

Visual
ETL

**AWS Pipeline/Soulution
Chart**