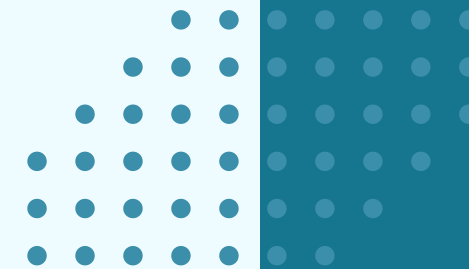# Diabetes Risk Analysis ETL Pipeline

Using Azure Data Factory and Azure Databricks
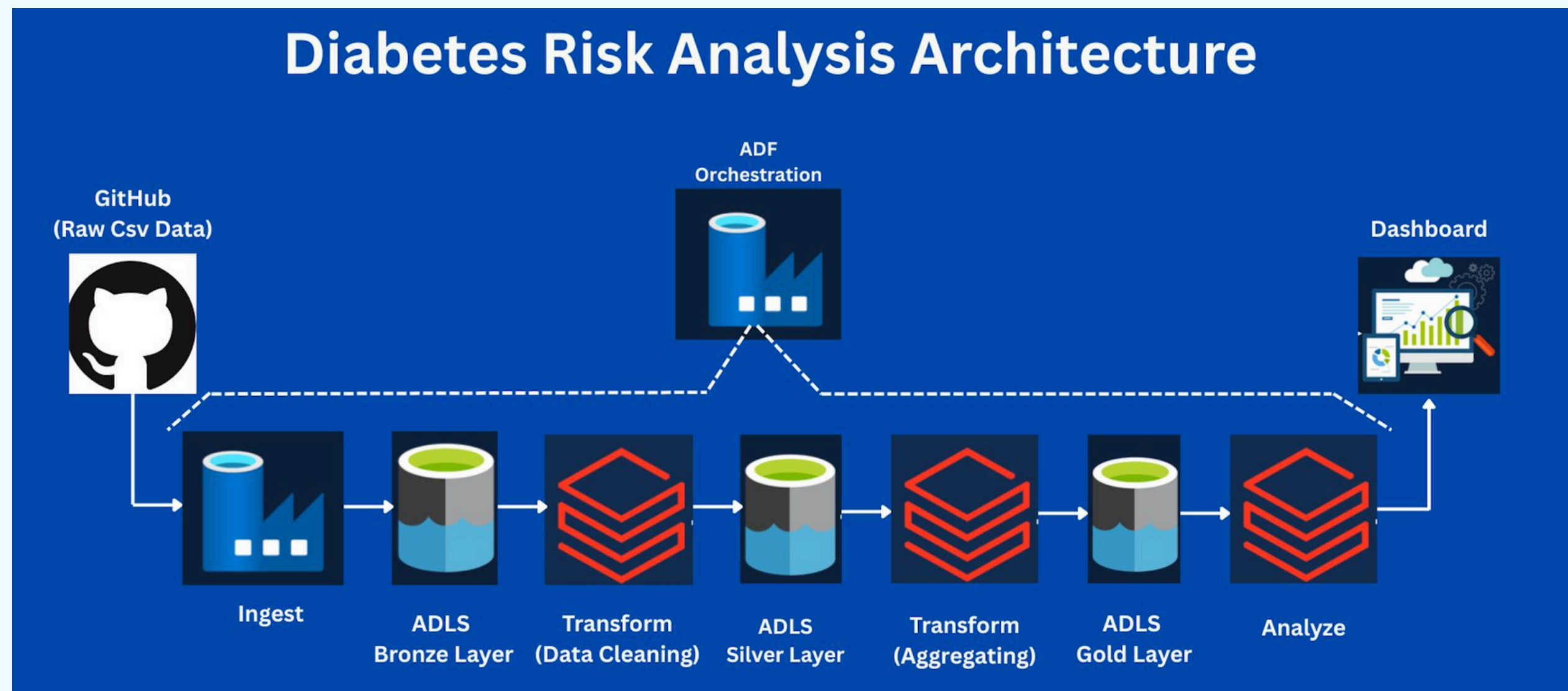
**PRESENTED BY**

Puneeth Kumar Amudala
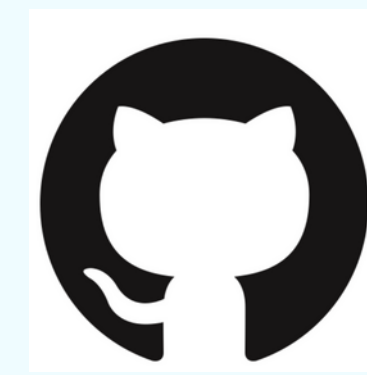
# Project Introduction

- DATASET FROM GITHUB (CSV)
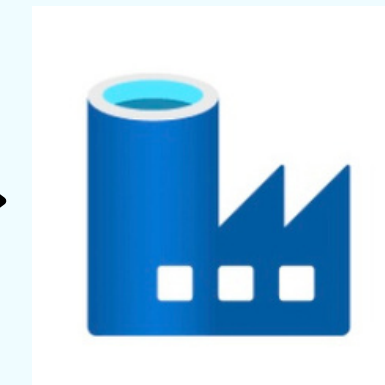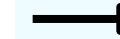- HEALTH INDICATORS RELATED TO DIABETES
- BATCH PROCESSING



## Diabetes Risk Analysis Architecture

GitHub (Raw Csv Data) → Ingest → ADLS Bronze Layer → Transform (Data Cleaning) → ADLS Silver Layer → Transform (Aggregating) → ADLS Gold Layer → Analyze → Dashboard

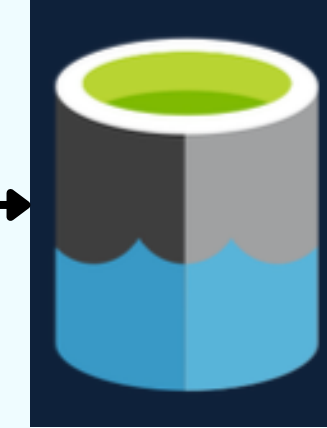ADF Orchestration

# Data Ingestion
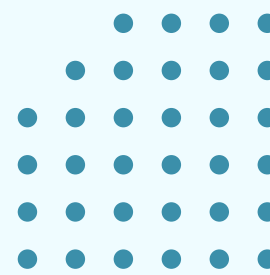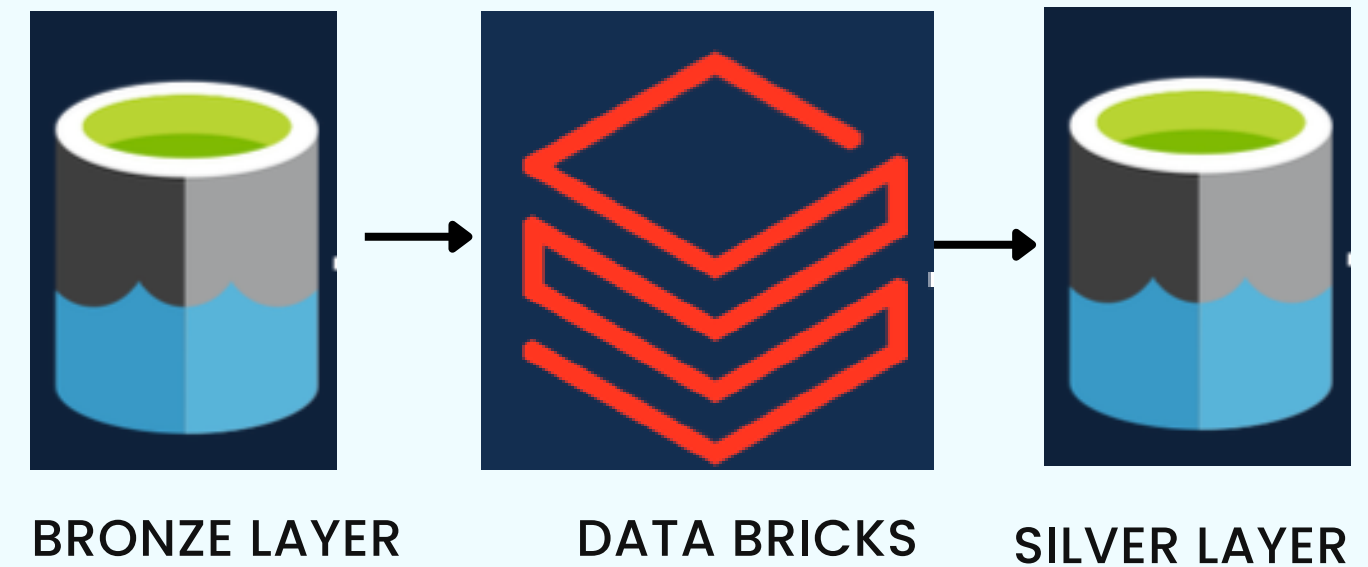**Tool: Azure Data Factory**



GITHUB → ADF → BRONZE LAYER

- PULLED FROM GITHUB USING HTTP CONNECTOR
- STORED RAW DATA IN ADLS BRONZE LAYER

# Data Transformation
## Tool: Azure Databricks

- BRONZE TO SILVER
- APPLIED SCHEMA
- REMOVED NULLS, DUPLICATES
- SAVED AS PARQUET IN SILVER LAYER



BRONZE LAYER     DATA BRICKS     SILVER LAYER

## 5. Load Cleaned Data to Silver Layer

21

```python
# Saving the cleaned data to the Silver layer in Parquet format, overwriting if the folder already exists
df_filtered.write.mode("overwrite").format("parquet").save("dbfs:/mnt/diabetes/silver/diabetes_cleaned")
```
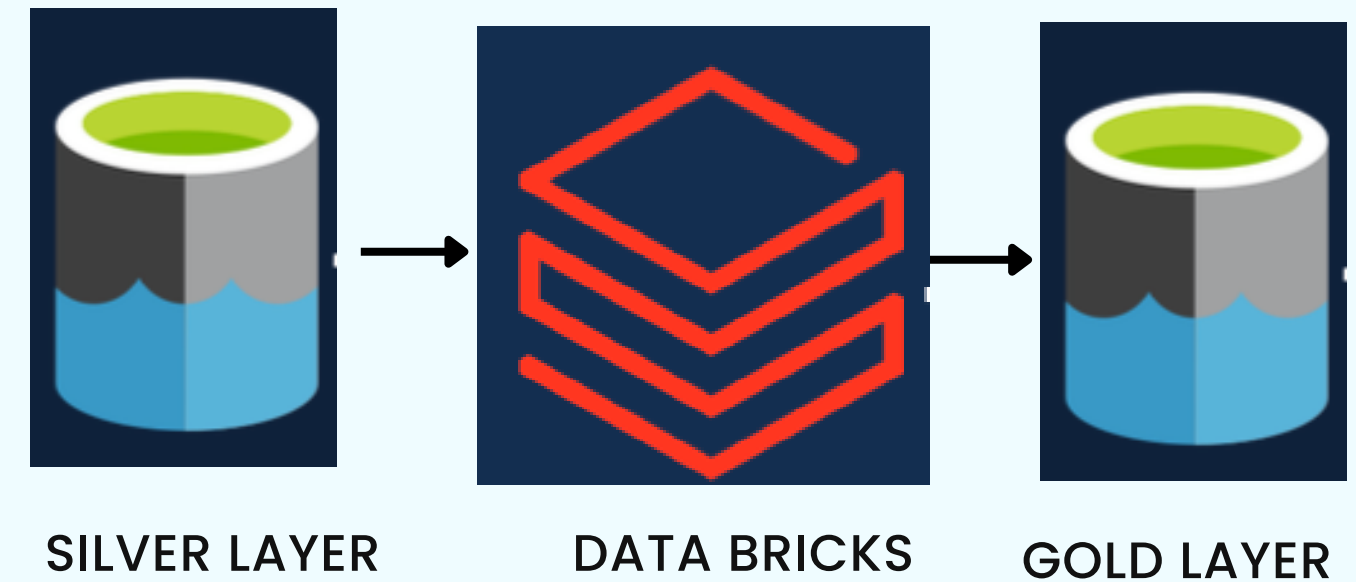
### Dropping Duplicated Values

```python
# Removing duplicate rows from the raw DataFrame
df= df_raw.dropDuplicates()
```

```python
# Filtering out records where BMI is less than 80
df_filtered= df.filter(col('BMI')< 80)
```

# Data Loading
## Tool: Azure Databricks

- BRONZE TO SILVER
- FEATURE ENGINEERING
- AGGREGATIONS
- SAVED AS DELTA IN GOLD LAYER



SILVER LAYER        DATA BRICKS        GOLD LAYER

### 3. Load the data into Gold Layer

```
25

# Group by AgeGroup, Sex, RiskLevel and calculate averages for health risk factors
df_gold = df_silver.groupBy("AgeGroup", "Sex","RiskLevel") \
    .agg(
        count("*").alias("TotalPatients"),
        round(avg("Diabetes_binary"), 3).alias("DiabetesRate"),
        round(avg("Obese"), 3).alias("ObesityRate"),
        round(avg("PhysActivity"), 3).alias("LowActivityRate"),
        round(avg("AlcoholRisk"), 3).alias("AlcoholRiskRate"),
        round(avg("HighMentalDistress"), 3).alias("MentalDistressRate"),
        round(avg("HighBP"), 3).alias("HighRiskRate")
    )
```

```
df_silver = df_silver \
    .withColumn("Obese", when(col("BMI") >= 30, 1.0).otherwise(0.0)) \
    .withColumn("HighMentalDistress", when(col("MentHlth") > 15, 1.0).otherwise(0.0))
    .withColumn("HighPhysicalDistress", when(col("PhysHlth") > 15, 1.0).otherwise(0.0
```

```
26

# Save the gold data in Delta format, partitioned by AgeGroup
df_gold.write.format("delta").mode("overwrite").partitionBy('AgeGroup').save("dbfs:/mnt/diabetes/gold/aggregated_results")
```
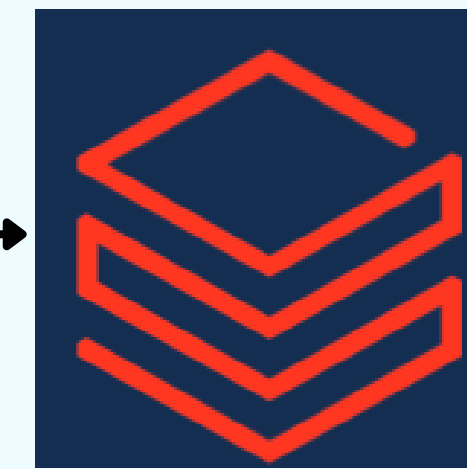
# ETL Orchestration

**Tool: Azure Data Factory**

1. COPY RAW DATA TO BRONZE

2. TRANSFORM BRONZE ➜ SILVER USING `BRONZE_TO_SILVER` NOTEBOOK

3. TRANSFORM SILVER ➜ GOLD USING `SILVER_TO_GOLD` NOTEBOOK

1. GENERATE VISUAL INSIGHTS USING `GOLD_TO_DASHBOARD` NOTEBOOK

# Key Findings

## DATA ANALYSIS

- Obesity strongly linked to risk
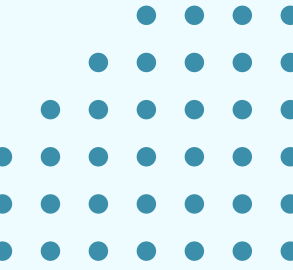- Males aged 40–60 at higher risk

## CHALLENGE FACED

Extracting data Directly from kaggle

## SOLUTION

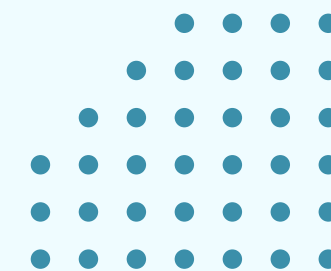Downloaded the data and Uploaded in My Github

# Conclusion

## Project Summary

- Full ETL pipeline from ingestion to insights
- Used Medallion Architecture (Bronze → Gold)
- Learned orchestration, storage formats, and dashboards

# Thank You