

# Capstone Project Submission Report

---

## Azure BFSI Data Pipeline for Customer Insights and Risk Analytics

---

- **Name:** HARSHIDA SHAILY
  - **Student ID:** 00127272
  - **Course/Program:** Data Engineering Coforge Ltd.
  - **Instructor Name:** Piyush Raj Katyayan
-

### 3. Problem Definition and Objectives

#### Problem Definition

In the Banking, Financial Services, and Insurance (BFSI) sector, managing and analyzing vast amounts of customer transaction data is crucial for risk assessment, fraud detection, and service optimization. Traditional data processing methods often fail to handle both historical and real-time data efficiently, leading to delayed insights, increased risks, and suboptimal customer experiences.

To address these challenges, the organization aims to develop a robust ETL pipeline using Azure services. This pipeline will enable seamless integration and processing of both batch and streaming data, ensuring real-time fraud detection, customer risk profiling, sentiment analysis, and demand forecasting. By leveraging advanced analytics, the company seeks to enhance customer satisfaction, mitigate financial risks, and optimize service offerings.

#### Objectives:

##### 1 . End-to-End ETL Pipeline:

The project aims to build a scalable ETL pipeline in Azure to handle both batch processing of historical data and real-time streaming of transaction data for a BFSI company.

##### 2. Risk Monitoring & Forecasting:

The pipeline will analyze customer transactions and risk profiles to detect fraud, assess credit risk, and forecast service demand, improving decision-making and financial security.

##### 3. Multi-Source Data Integration:

The solution will process diverse datasets, including historical transactions, customer profiles, service usage, real-time transactions, feedback, and service availability data, ensuring comprehensive analytics.

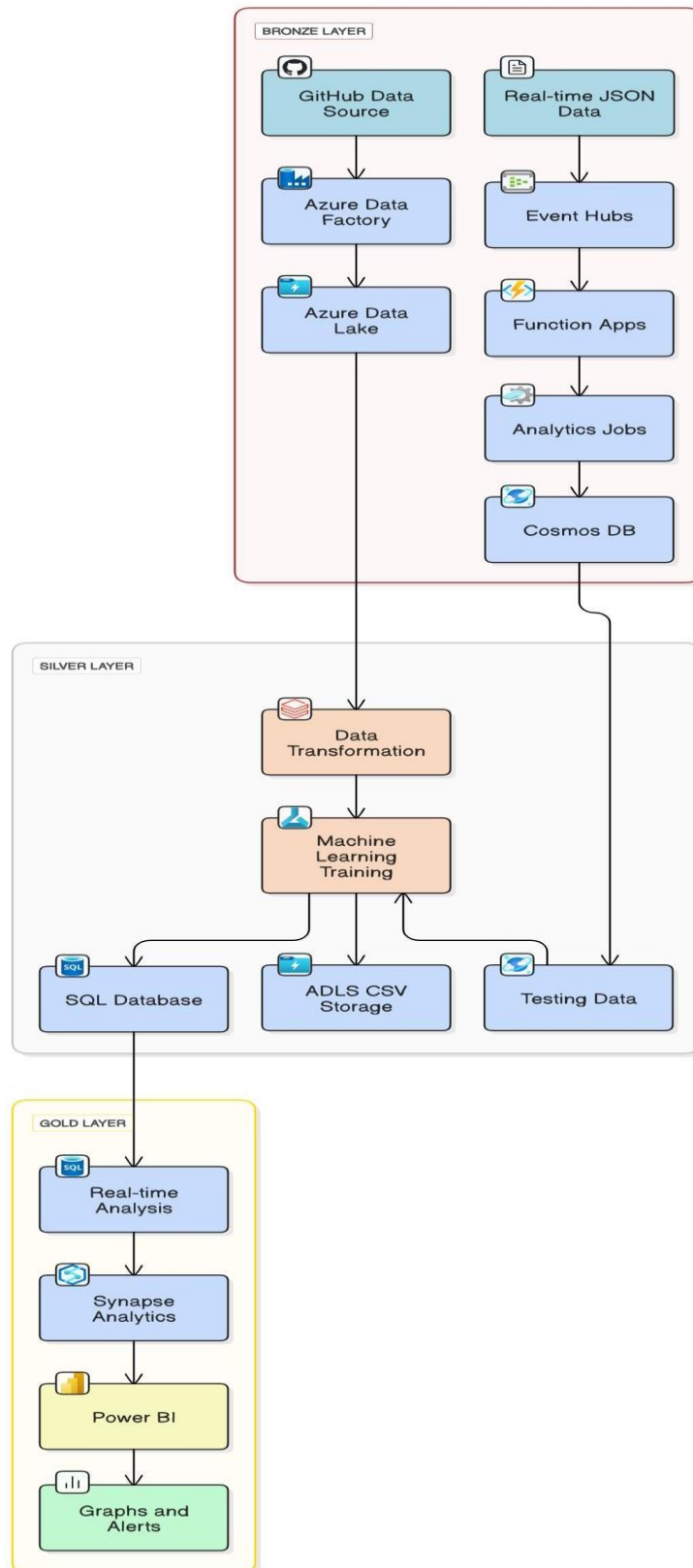
##### 4. Customer-Centric Insights:

By leveraging advanced analytics, the pipeline will enhance customer satisfaction through sentiment analysis, personalized services, and proactive issue resolution.

##### 5. Azure-Based Scalable Architecture:

The solution will utilize Azure services to ensure efficient data ingestion, transformation, storage, and real-time processing, enabling seamless operations and high availability.

Data Processing Flow Chart



## 4. Data Collection, Exploratory Data Analysis (EDA), and Preprocessing

### Data Collection :

#### I. HISTORICAL DATA

##### (CSV Files)

**Service Usage Data** – Records customer interactions with banking services to analyze service demand and optimize offerings.

**Customer Feedback Data** – Captures customer sentiments and feedback to improve service quality and customer satisfaction.

**Service Availability Data** – Tracks service uptime and performance to ensure reliability and predict potential downtime.

**Customer Profiles Data** – Stores customer demographic and financial information for segmentation and risk assessment.

**Historical Transactions Data** – Contains past financial transactions, helping in fraud detection, customer spending analysis, and trend forecasting.

### Data Ingestion:

- **Source:** GitHub (hosting raw CSV files).
- **Pipeline:** Azure Data Factory (ADF) extracts and loads data into ADLS.

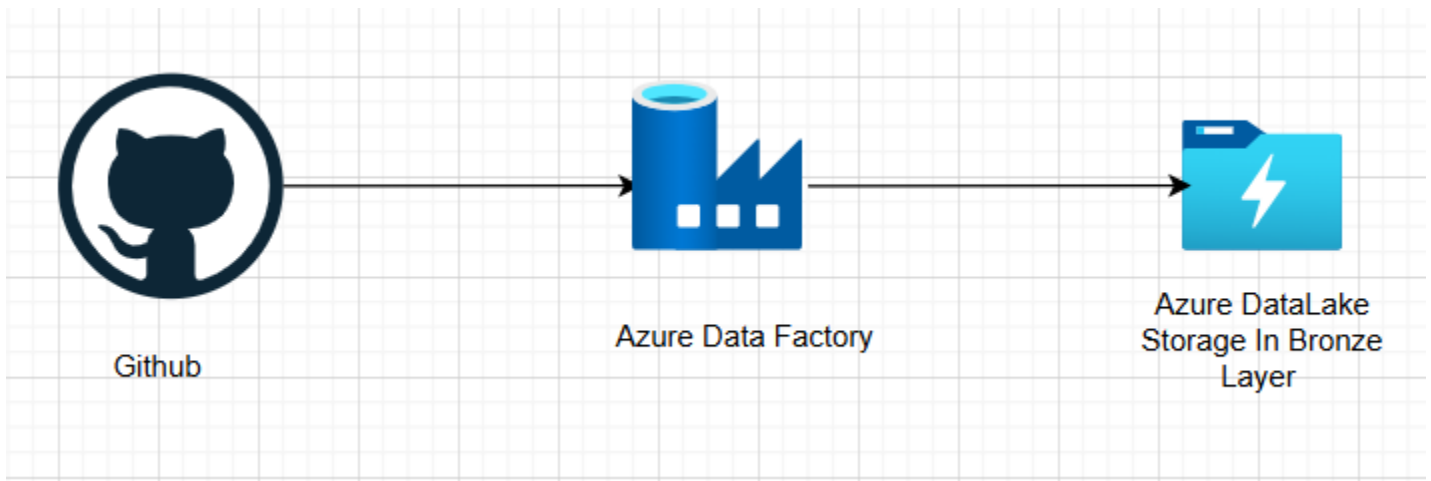
### Bronze Layer (Raw Data Storage):

- Data is stored in **ADLS** under a **Bronze** container to maintain original integrity and support future reprocessing needs.

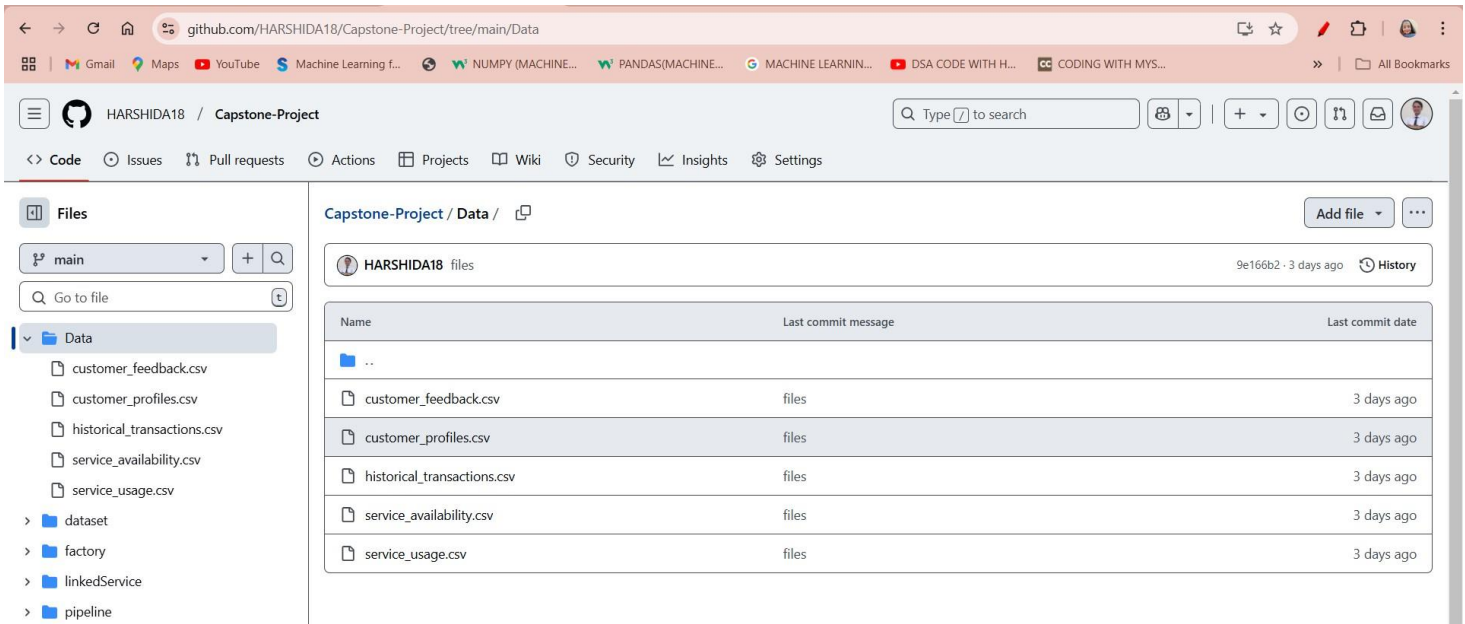
### Consumption:

- Data is used for **fraud detection, sentiment analysis, service forecasting, and customer risk profiling** through **Power BI dashboards, ML models, and real-time monitoring systems**.

## Architecture =



## Github Data =



The screenshot shows the Github repository page for HARSHIDA18/Capstone-Project. The repository is located at [github.com/HARSHIDA18/Capstone-Project/tree/main/Data](https://github.com/HARSHIDA18/Capstone-Project/tree/main/Data). The repository contains a directory named **Data**, which is expanded to show the following files:

- customer\_feedback.csv
- customer\_profiles.csv
- historical\_transactions.csv
- service\_availability.csv
- service\_usage.csv

The **Data** directory is further expanded to show the following files:

- dataset
- factory
- linkedService
- pipeline

The **factory** directory is further expanded to show the following files:

- customer\_feedback.csv
- customer\_profiles.csv
- historical\_transactions.csv
- service\_availability.csv
- service\_usage.csv

The **pipeline** directory is further expanded to show the following files:

- customer\_feedback.csv
- customer\_profiles.csv
- historical\_transactions.csv
- service\_availability.csv
- service\_usage.csv

## Pipeline Of Azure DataFactory Data Ingestion

Microsoft Azure | Data Factory | CAPSTONEHARSHIDAADF

Would you like to see Data Factory inside of Microsoft Fabric, Microsoft's newest cloud-first data analytics SaaS platform? Click [here](#) to get started with Fabric Data Factory!

main branch | Validate all | Save all | Publish

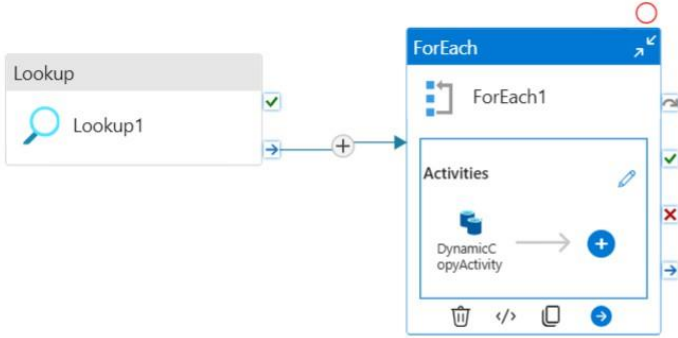
### Factory Resources

Filter resources by name

- Pipelines 1
  - DynamicPipeline
- Change Data Capture (preview) 0
- Datasets 3
- Data flows 0
- Power Query 0
- Templates 0

### DynamicPipeline

Saved | Save as template | Validate | Debug | Add trigger



```

graph LR
    Lookup[Lookup1] --> Join((+))
    Join --> ForEach[ForEach1]
    subgraph ForEach1 [ForEach]
        direction TB
        subgraph Activities
            DynamicCopy[DynamicCopyActivity] --> Plus((+))
        end
    end
  
```

## Linked Services

Microsoft Azure | Data Factory | CAPSTONEHARSHIDAADF

Would you like to see Data Factory inside of Microsoft Fabric, Microsoft's newest cloud-first data analytics SaaS platform? Click [here](#) to get started with Fabric Data Factory!

main branch | Validate all | Save all | Publish

### Linked Services

Linked service defines the connection information to a data store or compute. [Learn more](#)

+ New

Filter by name | Annotations: Any

Showing 1 - 3 of 3 items

Name ↑↓	Type ↑↓	Related ↑↓
AzureDataLakeStorage1	Azure Data Lake Storage Gen2	2
AzureSqlDatabase1	Azure SQL Database	0
HttpServer1	HTTP	1

## Azure DataLake Gen2 Data in bronze layer file =

Home > Resource groups > CAPSTONE-HARSHIDA > capstonestorageharshida

capstonestorageharshida | Containers

Storage account

Search

+ Container Change access level Restore containers Refresh Delete Give feedback

Search containers by prefix

Show deleted containers

Name	Last modified	Anonymous access level	Lease state
<input type="checkbox"/> \$logs	2/23/2025, 11:11:07 AM	Private	Available
<input type="checkbox"/> bronze	2/23/2025, 11:11:35 AM	Private	Available
<input type="checkbox"/> gold	2/23/2025, 11:11:52 AM	Private	Available
<input type="checkbox"/> reference	2/23/2025, 11:12:07 AM	Private	Available
<input type="checkbox"/> silver	2/23/2025, 11:11:45 AM	Private	Available
<input type="checkbox"/> synapseanalytics	2/25/2025, 3:30:23 AM	Private	Available

Overview

Activity log

Tags

Diagnose and solve problems

Access Control (IAM)

Data migration

Events

Storage browser

Partner solutions

Data storage

Containers

File shares

Queues

Tables

Home > Resource groups > CAPSTONE-HARSHIDA > capstonestorageharshida | Containers >

bronze

Container

Search

Upload Add Directory Refresh Rename Delete Change tier Acquire lease Break lease

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: bronze

Search blobs by prefix (case-sensitive)

Name	Modified	Access tier	Archive status
<input type="checkbox"/> customerFeedback	2/23/2025, 1:18:48 PM		
<input type="checkbox"/> customerProfiles	2/23/2025, 1:19:02 PM		
<input type="checkbox"/> historicalTransactions	2/23/2025, 1:19:15 PM		
<input type="checkbox"/> serviceAvailability	2/23/2025, 1:19:29 PM		
<input type="checkbox"/> serviceUsage	2/23/2025, 1:19:46 PM		

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

## 2. REAL-TIME DATA STREAMING

### ARCHITECTURE =

- The data through function app's code connected with event hub .
- The event hub generates the events instance per defined interval .
- The instance through spark analytics jobs as an input is the output for cosmosDB Item instance creation .

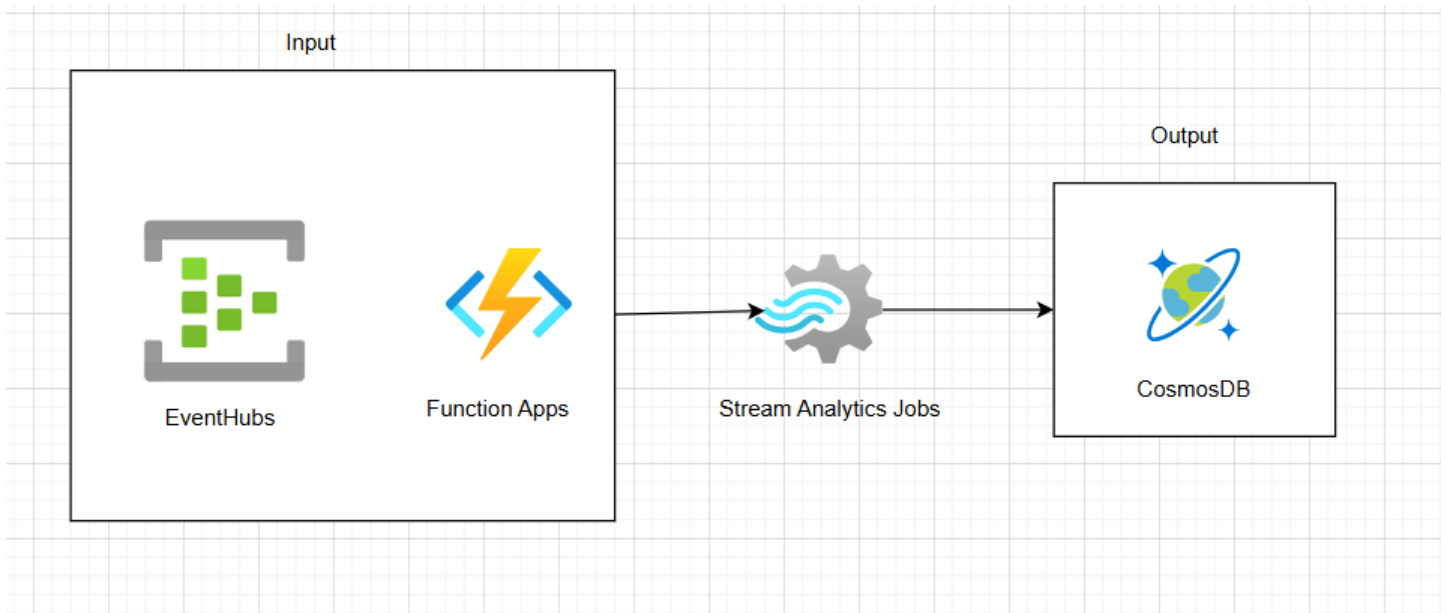
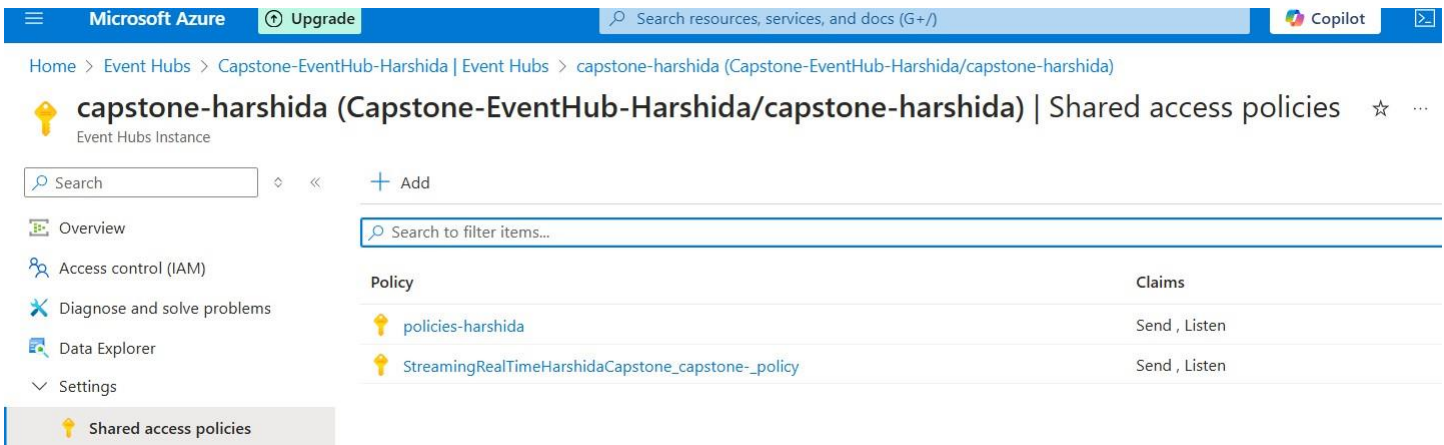


Fig : Real-time data entering into the CosmosDB instance

### EventHubs =



The screenshot shows the Microsoft Azure portal interface. The top navigation bar includes the 'Microsoft Azure' logo, an 'Upgrade' button, a search bar, and a 'Copilot' button. The breadcrumb trail indicates the path: Home > Event Hubs > Capstone-EventHub-Harshida | Event Hubs > capstone-harshida (Capstone-EventHub-Harshida/capstone-harshida). The main heading is 'capstone-harshida (Capstone-EventHub-Harshida/capstone-harshida) | Shared access policies'. Below this, there is a search bar and a list of policies. The left sidebar contains navigation links: Overview, Access control (IAM), Diagnose and solve problems, Data Explorer, and Settings. The 'Shared access policies' section is highlighted.

Policy	Claims
policies-harshida	Send , Listen
StreamingRealTimeHarshidaCapstone_capstone-_policy	Send , Listen

Fig : EventHub policies definition



## Stream Analytics Jobs =

### Input

alTimeHarshidaCapstone

### StreamingRealTimeHarshidaCapstone | Inputs

Stream Analytics job

Search  ⌵ ⌵ + Add input ⌵ ↻ Refresh

- Overview
- Activity log
- Access control (IAM)
- Tags
- Diagnose and solve problems
- Job topology
- Inputs** ☆
- Functions
- Query
- Outputs

Alias ↑	Source t...	T...
capstone-harshida	Stream	E...

### Output

### StreamingRealTimeHarshidaCapstone | Outputs

Stream Analytics job

Search  ⌵ ⌵ + Add output ⌵ ↻ Refresh

- Overview
- Activity log
- Access control (IAM)
- Tags
- Diagnose and solve problems
- Job topology
- Inputs
- Functions**

Alias ↑	Type	A...	R...
cosmos-output	Cosmos DB	Contr	

## CosmosDB

dacapstoneproject

**cosmosdbharshidacapstoneproject** | Data Explorer ☆ ...

Azure Cosmos DB account

Search

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Cost Management

Quick start

**Data Explorer**

Mirroring in Fabric (Preview)

Settings

Integrations

Containers

Browse

Scale

Settings

Monitoring

Automation

Help

+ New Container

Home

harshidacapstoneDB

harshidaContainer

Items

Scale & Settings

Stored Procedures

User Defined Functions

Triggers

ToDoList

Items

id

FC0B659C-C1EF-41F3-AB9...

BF87ACE3-C52B-44... 7FF...

A6BB4603-7CD5-4... AB9...

5897882E-D4C6-4... 3E4...

CB038CA5-3728-4... 4F2...

A9EFB9E2-8859-44... AA...

BD340F0A-F661-4E... 75B...

EF16A6FA-9BE2-4A... 32A...

3B52D15D-DF6C-4... 3E4...

0B013EA7-B40E-49... 3E4...

26E6C049-667F-44... 879...

Load more

SELECT \* FROM c

Type a query predicate (e.g., WHERE c.id='1'), or choose one from the drop down list, o

Apply Filter

```
{
  "id": "FC0B659C-C1EF-41F3-AFE2-F87C7F43AD48",
  "categoryId": "AB952F9F-5ABA-4251-BC2D-AFF8DF41",
  "categoryName": "Components, Headsets",
  "sku": "HS-0296",
  "name": "LL Headset",
  "description": "The product called \"LL Headset\"",
  "price": 34.2,
  "tags": [
    {
      "id": "18AC309F-F81C-4234-A752-5DD02BEA",
      "name": "Tag-32"
    },
    {
      "id": "1B387A00-57D3-4444-8331-18A90725",
      "name": "Tag-43"
    },
    {
      "id": "C6AB3E24-BA48-40F0-A260-CB04EB03",
      "name": "Tag-73"
    },
    {
      "id": "DAC25651-3DD3-4483-8FD1-581DC41E",
      "name": "Tag-56"
    }
  ]
}
```

## Table: Data Description

- ⇒ The data in the bronze layer are stored in form of the container per data file .
- ⇒ The data files of the historical data contains the csv format files .
- ⇒ The real-time data is in json format which when stored in the cosmosDB is stored as the item instance of the container into the database .

Home > Resource groups > CAPSTONE-HARSHIDA > capstonestorageharshida | Containers > bronze >

**bronze**  
Container

Search

Upload Add Directory

**Overview**

Diagnose and solve problems

Access Control (IAM)

Settings

**Authentication method:** Access key (Switch to Microsoft Entra user account)

**Location:** bronze / customerFeedback

Search blobs by prefix (case-...)

Show deleted objects

**Name**

<input type="checkbox"/>	[-]	...
<input type="checkbox"/>	customer_feedback.csv	...

**customerFeedback/customer\_feedback.csv**

Blob

Save Discard Download Refresh Delete

Overview Versions Edit Generate SAS

```

1 feedback_id,customer_id,rating,feedback_text,feedback_date
2 FBK5000,CUST1068,4,Very helpful customer support.,2023-03-15
3 FBK5001,CUST1039,2,Frustrating mobile app experience!,2023-04-09
4 FBK5002,CUST1080,3,Nothing special.,2023-01-08
5 FBK5003,CUST1028,2,Unexpected charges!,2023-06-20
6 FBK5004,CUST1064,4,Fantastic mobile app!,2023-09-12
7 FBK5005,CUST1046,4,Fantastic mobile app!,2023-07-24
8 FBK5006,CUST1059,3,Some delays but overall okay.,2023-03-27
9 FBK5007,CUST1086,4,Very responsive support team!,2023-04-08
10 FBK5008,CUST1003,4,Loved the credit card offers!,2023-05-07
11 FBK5009,CUST1053,4,Great customer relationship!,2023-08-28
12 FBK5010,CUST1004,5,Highly satisfied!,2023-02-04
13 FBK5011,CUST1004,3,Some delays but overall okay.,2023-02-01
14 FBK5012,CUST1068,2,Very slow processing!,2023-03-28
15 FBK5013,CUST1080,4,Great customer relationship!,2023-01-30
16 FBK5014,CUST1073,1,Loan approval process too slow!,2023-08-08
17 FBK5015,CUST1092,5,No hassles in transactions!,2023-03-26
18 FBK5016,CUST1071,2,Unexpected transaction failures!,2023-04-24

```

Csv Preview

Home

harsh...Items

SELECT \* FROM c Type a query predicate (e.g., WHERE c.id='1'), or choose one from the drop down list, o Apply Filter

<input type="checkbox"/>	1	{
<input type="checkbox"/>	2	"id": "FC08659C-C1EF-41F3-AFE2-F87C7F43AD48",
<input checked="" type="checkbox"/>	3	"categoryId": "AB952F9F-5ABA-4251-BC2D-AFF8DF412A4A",
<input type="checkbox"/>	4	"categoryName": "Components, Headsets",
<input type="checkbox"/>	5	"sku": "HS-0296",
<input type="checkbox"/>	6	"name": "LL Headset",
<input type="checkbox"/>	7	"description": "The product called \"LL Headset\"",
<input type="checkbox"/>	8	"price": 34.2,
<input type="checkbox"/>	9	"tags": [
<input type="checkbox"/>	10	{
<input type="checkbox"/>	11	"id": "18AC309F-F81C-4234-A752-5DD02BEAE83",
<input type="checkbox"/>	12	"name": "Tag-32"
<input type="checkbox"/>	13	},
<input type="checkbox"/>	14	{
<input type="checkbox"/>	15	"id": "1B387A00-57D3-4444-8331-18A90725E98B",
<input type="checkbox"/>	16	"name": "Tag-43"
<input type="checkbox"/>	17	},
<input type="checkbox"/>	18	{
<input type="checkbox"/>	19	"id": "C6AB3E24-BA48-40F0-A260-CB04EB03D5B0",
<input type="checkbox"/>	20	"name": "Tag-73"
<input type="checkbox"/>	21	},
<input type="checkbox"/>	22	{
<input type="checkbox"/>	23	"id": "DAC25651-3DD3-4483-8FD1-581DC41EF34B",
<input type="checkbox"/>	24	"name": "Tag-56"
<input type="checkbox"/>	25	},
<input type="checkbox"/>	26	}

Load

## Exploratory Data Analysis (EDA)

- Transaction Trends & Seasonality** – Analyzed historical transactions to identify peak transaction periods, spending patterns, and seasonal variations. This helps in service optimization and fraud detection.
- Customer Segmentation Insights** – Grouped customers based on demographics, transaction behavior, and risk profiles to personalize financial products and detect high-risk customers.

Table +							Q F
	A <sub>C</sub> feedback_id	A <sub>C</sub> customer_id	1 <sub>3</sub> rating	A <sub>C</sub> feedback_text	feedback_date	1 <sub>3</sub> sentiment_score	
2	FBK5001	CUST1039	2	Frustrating mobile app experience!	2023-04-09	0	
3	FBK5002	CUST1080	3	Nothing special.	2023-01-08	0	
4	FBK5003	CUST1028	2	Unexpected charges!	2023-06-20	0	
5	FBK5004	CUST1064	4	Fantastic mobile app!	2023-09-12	1	
6	FBK5005	CUST1046	4	Fantastic mobile app!	2023-07-24	1	
7	FBK5006	CUST1059	3	Some delays but overall okay.	2023-03-27	0	
8	FBK5007	CUST1086	4	Very responsive support team!	2023-04-08	0	
9	FBK5008	CUST1003	4	Loved the credit card offers!	2023-05-07	0	
10	FBK5009	CUST1053	4	Great customer relationship!	2023-08-28	1	
11	FBK5010	CUST1004	5	Highly satisfied!	2023-02-04	1	
12	FBK5011	CUST1004	3	Some delays but overall okay.	2023-02-01	0	
13	FBK5012	CUST1068	2	Very slow processing!	2023-03-28	0	
14	FBK5013	CUST1080	4	Great customer relationship!	2023-01-30	1	
15	FBK5014	CUST1073	1	Loan approval process too slow!	2023-08-08	0	
16	FBK5015	CUST1092	5	No hassles in transactions!	2023-03-26	0	

↓ 100 rows | 1.58s runtime Refreshed yesterday

- Sentiment Analysis on Feedback** – Processed customer feedback data using NLP to classify sentiments (positive, neutral, negative), helping improve banking services.
- Correlation Analysis** – Used a heatmap to identify relationships between variables, such as the impact of service availability on transaction frequency and customer satisfaction.

Table +							Q F
	A <sub>C</sub> customer_id	transaction_date	A <sub>C</sub> transaction_type	1.2 transaction_amount	A <sub>C</sub> payment_method	1.2 fraud_risk_score	
5	CUST1028	2023-08-14	Withdrawal	30000	Debit Card	0	
6	CUST1091	2023-08-14	Credit Card Payment	45529.1	Debit Card	0	
7	CUST1020	2023-05-31	Investment	135651.02	Bank Transfer	0	
8	CUST1058	2023-03-02	Fund Transfer	10000	Credit Card	1	
9	CUST1078	2023-05-31	Withdrawal	20000	Debit Card	0	
10	CUST1053	2023-01-01	Fund Transfer	20000	Debit Card	1	
11	CUST1094	2023-07-30	Loan Payment	92762.32	Net Banking	0	
12	CUST1029	2023-04-16	Loan Payment	225036.14	Debit Card	0	
13	CUST1043	2023-03-17	Bill Payment	5737.51	Credit Card	0	
14	CUST1016	2023-02-15	Credit Card Payment	12720.01	Credit Card	0	
15	CUST1022	2023-07-30	Withdrawal	5000	Credit Card	0	
16	CUST1045	2023-09-28	Insurance Premium	39119.8	Credit Card	0	
17	CUST1001	2023-10-13	Credit Card Payment	11004.81	Debit Card	0	
18	CUST1080	2023-05-31	Fund Transfer	5000	Net Banking	0	

## 5. Anomaly Detection in Transactions – Identified suspicious transaction patterns using box plots and statistical methods, assisting in fraud detection.

### EDA Visualization Diagram

The following visualizations are commonly used:

- **Histograms** – To analyze the distribution of transaction amounts.
- **Correlation Heatmap** – To understand relationships between customer attributes and transaction behavior.
- **Bar Graphs** – To visualize service usage across different customer segments.
- **Box Plots** – To detect transaction outliers indicating potential fraud.
- **Time Series Plots** – To track transaction volume trends over time.

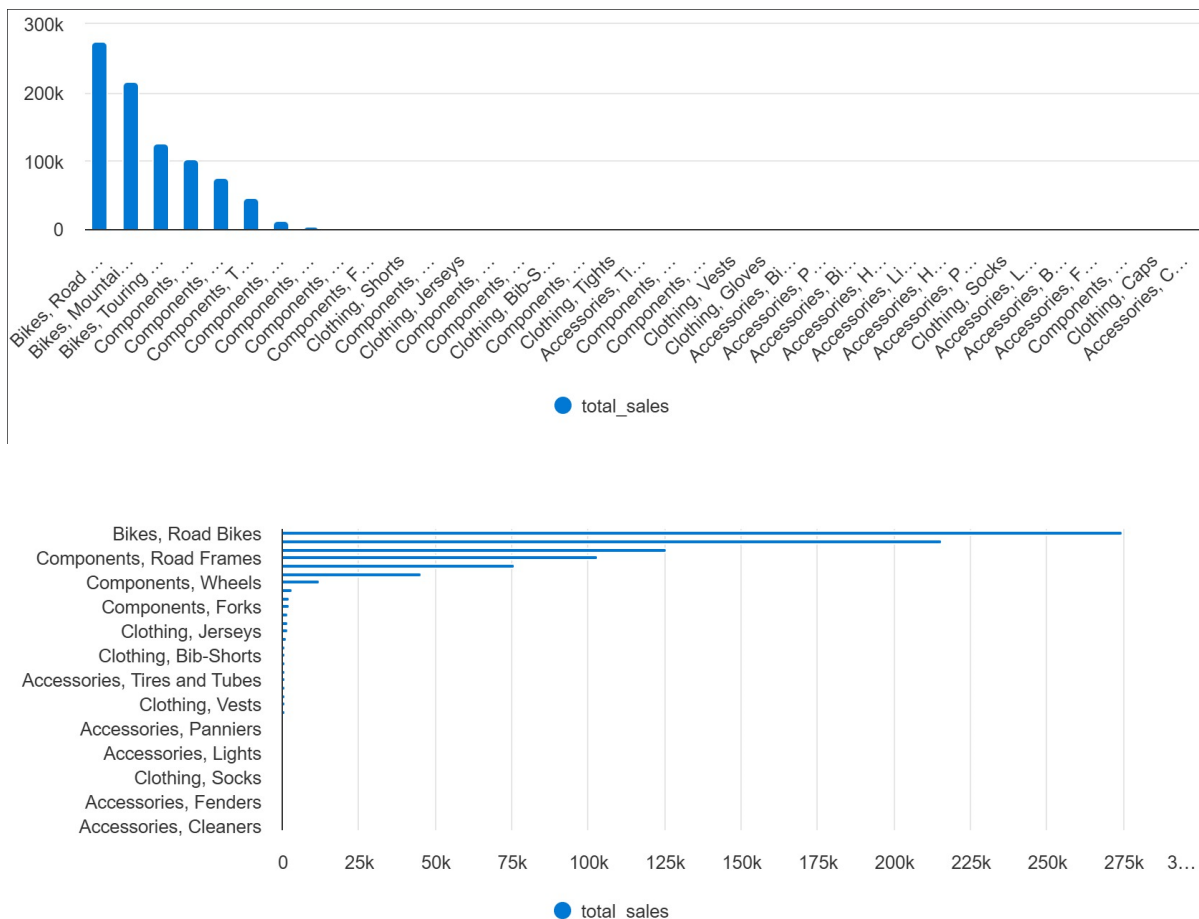


Fig : Total Sales and Average Price By Category

● 1176 ● 4

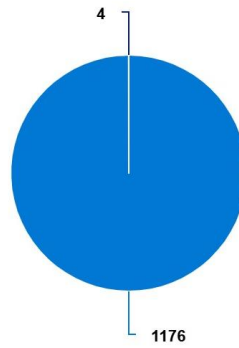


Fig : Fraud Prediction Detection

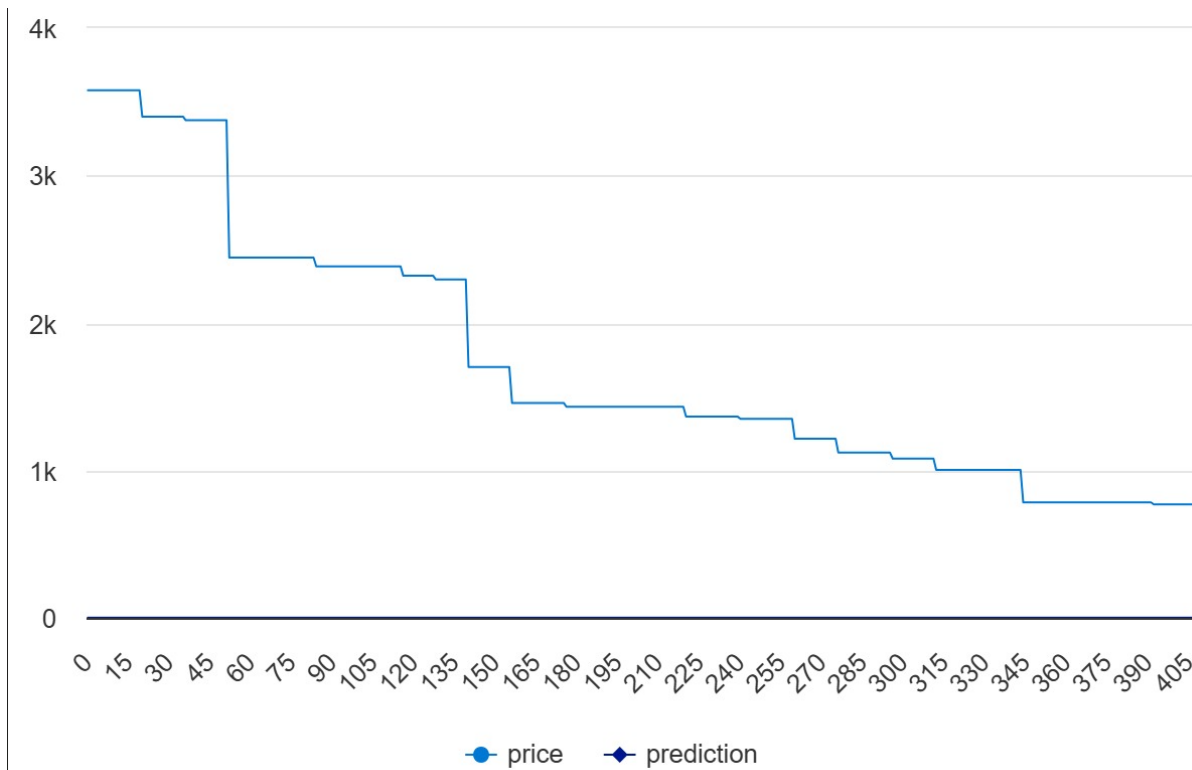


Fig : High Risk prediction based on Price

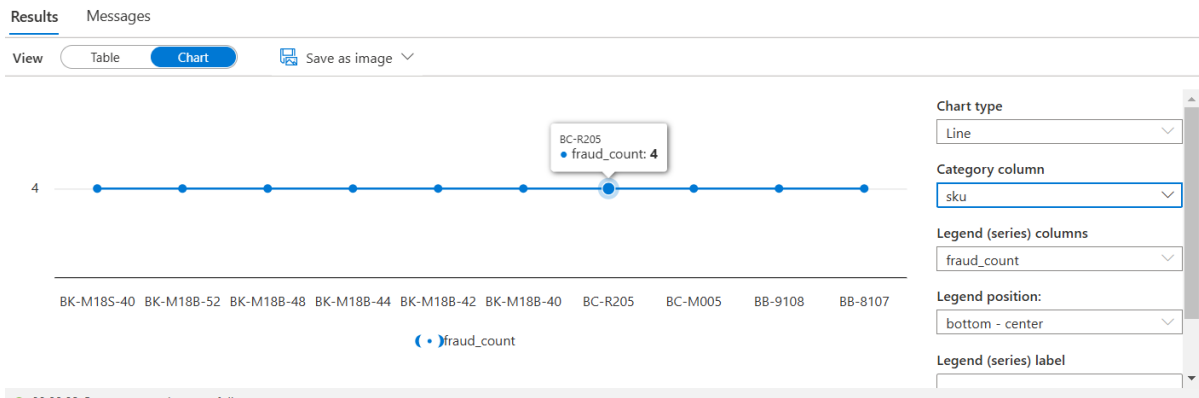


Fig : Fraud Content Based On SKU

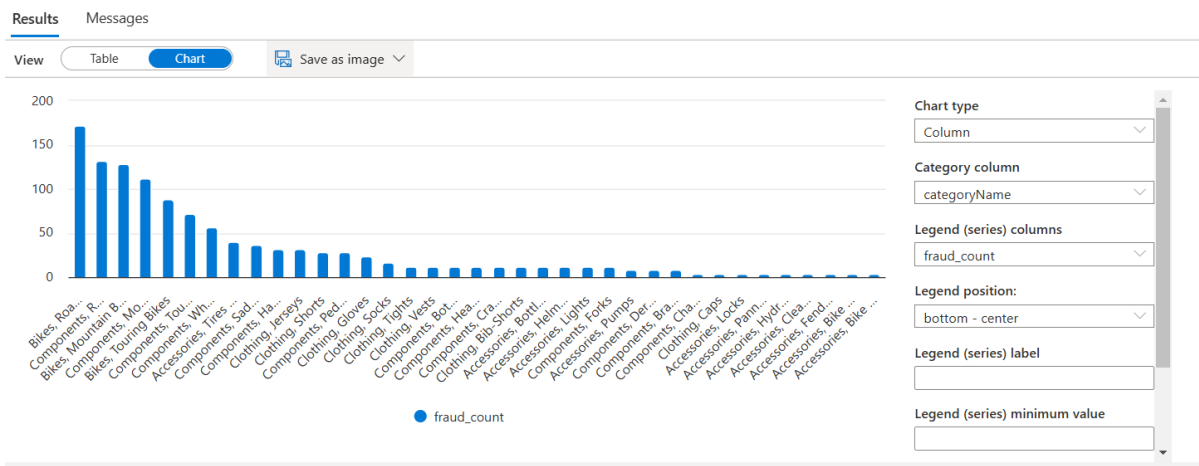


Fig : Fraud Vs Non-Fraud by Category

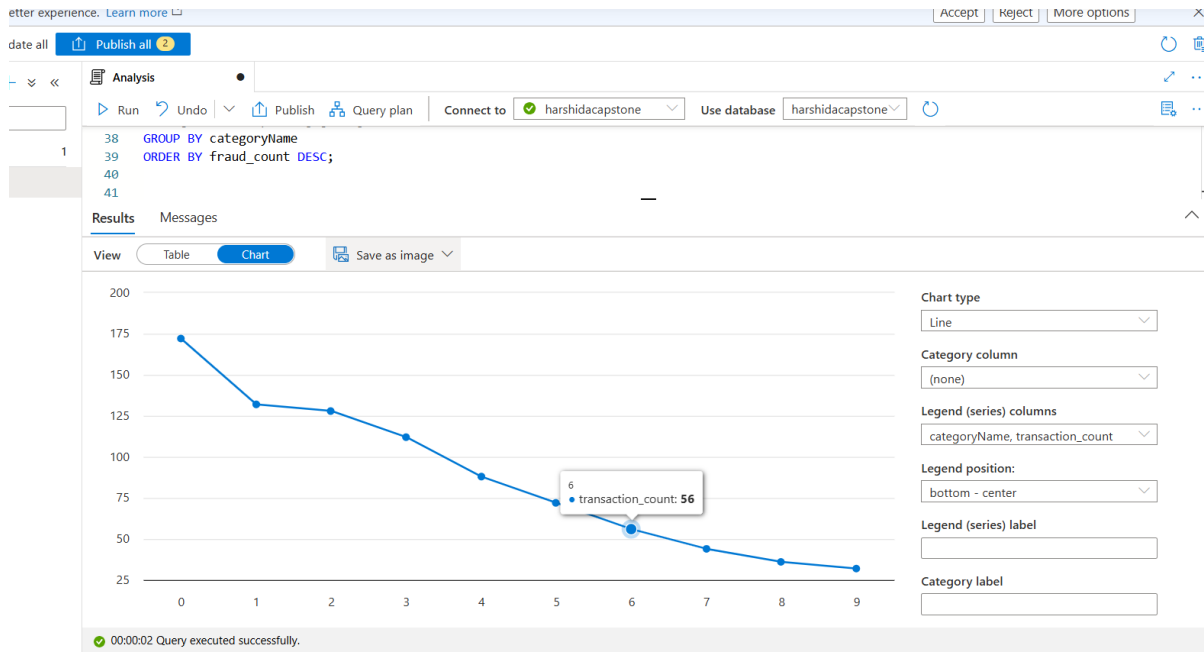


Fig : Most Frequent Categories by Transactions



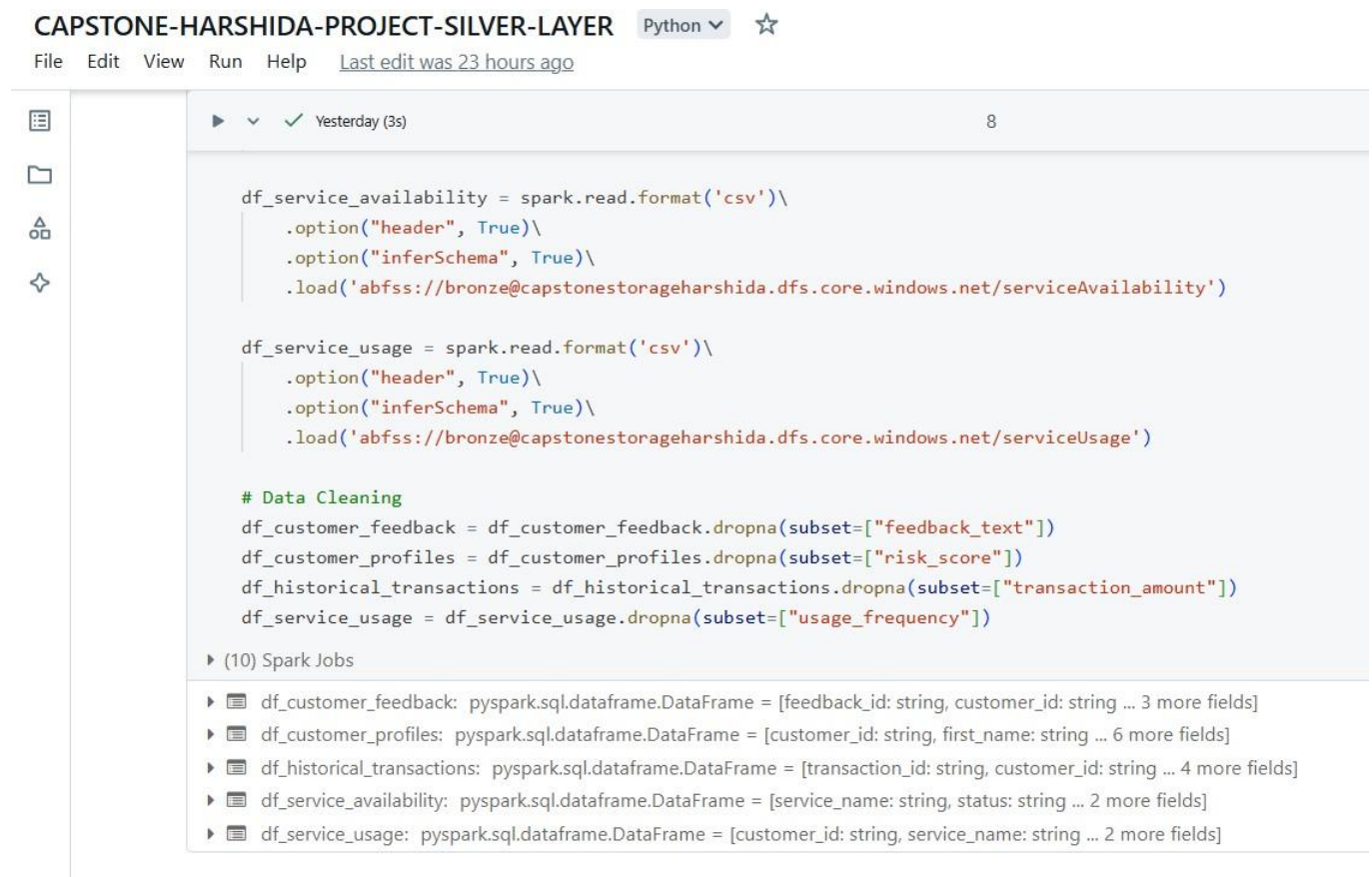
## Data Preprocessing

1. **Handling Missing Values** – Imputed missing values in customer profiles and transaction data using mean/mode/median strategies to ensure data completeness.

**CODE =**

```
df_customer_feedback = df_customer_feedback.dropna(subset=["feedback_text"])
df_customer_profiles = df_customer_profiles.dropna(subset=["risk_score"])
df_historical_transactions = df_historical_transactions.dropna(subset=["transaction_amount"])
df_service_usage = df_service_usage.dropna(subset=["usage_frequency"])
```

**OUTPUT =**



The screenshot shows a Jupyter Notebook titled "CAPSTONE-HARSHIDA-PROJECT-SILVER-LAYER" with a Python kernel. The code in the cell reads data from Azure Data Lake Storage (ADLS) and performs data cleaning. The output shows the Spark jobs executed and the resulting DataFrames.

```
df_service_availability = spark.read.format('csv')\
    .option("header", True)\
    .option("inferSchema", True)\
    .load('abfss://bronze@capstonestorageharshida.dfs.core.windows.net/serviceAvailability')

df_service_usage = spark.read.format('csv')\
    .option("header", True)\
    .option("inferSchema", True)\
    .load('abfss://bronze@capstonestorageharshida.dfs.core.windows.net/serviceUsage')

# Data Cleaning
df_customer_feedback = df_customer_feedback.dropna(subset=["feedback_text"])
df_customer_profiles = df_customer_profiles.dropna(subset=["risk_score"])
df_historical_transactions = df_historical_transactions.dropna(subset=["transaction_amount"])
df_service_usage = df_service_usage.dropna(subset=["usage_frequency"])
```

▶ (10) Spark Jobs

- ▶ df\_customer\_feedback: pyspark.sql.dataframe.DataFrame = [feedback\_id: string, customer\_id: string ... 3 more fields]
- ▶ df\_customer\_profiles: pyspark.sql.dataframe.DataFrame = [customer\_id: string, first\_name: string ... 6 more fields]
- ▶ df\_historical\_transactions: pyspark.sql.dataframe.DataFrame = [transaction\_id: string, customer\_id: string ... 4 more fields]
- ▶ df\_service\_availability: pyspark.sql.dataframe.DataFrame = [service\_name: string, status: string ... 2 more fields]
- ▶ df\_service\_usage: pyspark.sql.dataframe.DataFrame = [customer\_id: string, service\_name: string ... 2 more fields]

2. **Data Normalization** – Scaled numerical features (transaction amounts, service usage) using Min-Max scaling to ensure uniformity across datasets.
3. **Encoding Categorical Variables** – Converted categorical fields like service type, customer risk level, and transaction mode into numerical representations using one-hot encoding. ( Sentiment Analysis )

CODE =

```
from pyspark.sql.functions import udf, when, col
from pyspark.sql.types import FloatType
from textblob import TextBlob

# Sentiment Analysis using TextBlob
def get_sentiment(text):
    if text is None: # Handle None values
        return 0.0
    return TextBlob(text).sentiment.polarity

# Register UDF
sentiment_udf = udf(get_sentiment, FloatType())

# Apply sentiment analysis
df_customer_feedback = df_customer_feedback.withColumn("sentiment_score",
sentiment_udf(col("feedback_text")))

# Convert sentiment_score into binary labels
df_customer_feedback = df_customer_feedback.withColumn(
    "sentiment_score", when(col("sentiment_score") >= 0.5, 1).otherwise(0)
)

df_customer_feedback.display()
```

Table								
	feedback_id	customer_id	rating	feedback_text	feedback_date	sentiment_score		
2	FBK5001	CUST1039	2	Frustrating mobile app experience!	2023-04-09	0		
3	FBK5002	CUST1080	3	Nothing special.	2023-01-08	0		
4	FBK5003	CUST1028	2	Unexpected charges!	2023-06-20	0		
5	FBK5004	CUST1064	4	Fantastic mobile app!	2023-09-12	1		
6	FBK5005	CUST1046	4	Fantastic mobile app!	2023-07-24	1		
7	FBK5006	CUST1059	3	Some delays but overall okay.	2023-03-27	0		
8	FBK5007	CUST1086	4	Very responsive support team!	2023-04-08	0		
9	FBK5008	CUST1003	4	Loved the credit card offers!	2023-05-07	0		
10	FBK5009	CUST1053	4	Great customer relationship!	2023-08-28	1		
11	FBK5010	CUST1004	5	Highly satisfied!	2023-02-04	1		
12	FBK5011	CUST1004	3	Some delays but overall okay.	2023-02-01	0		
13	FBK5012	CUST1068	2	Very slow processing!	2023-03-28	0		
14	FBK5013	CUST1080	4	Great customer relationship!	2023-01-30	1		
15	FBK5014	CUST1073	1	Loan approval process too slow!	2023-08-08	0		
16	FBK5015	CUST1092	5	No hassles in transactions!	2023-03-26	0		

- CODE =**

### # Fraud Risk Score based on transaction amount and type

```
df_historical_transactions.display()
```

OUTPUT =

	customer_id	transaction_date	transaction_type	transaction_amount	payment_method	fraud_risk_score
5	CUST1028	2023-08-14	Withdrawal	30000	Debit Card	0
6	CUST1091	2023-08-14	Credit Card Payment	45529.1	Debit Card	0
7	CUST1020	2023-05-31	Investment	135651.02	Bank Transfer	0
8	CUST1058	2023-03-02	Fund Transfer	10000	Credit Card	1
9	CUST1078	2023-05-31	Withdrawal	20000	Debit Card	0
10	CUST1053	2023-01-01	Fund Transfer	20000	Debit Card	1
11	CUST1094	2023-07-30	Loan Payment	92762.32	Net Banking	0
12	CUST1029	2023-04-16	Loan Payment	225036.14	Debit Card	0
13	CUST1043	2023-03-17	Bill Payment	5737.51	Credit Card	0
14	CUST1016	2023-02-15	Credit Card Payment	12720.01	Credit Card	0
15	CUST1022	2023-07-30	Withdrawal	5000	Credit Card	0
16	CUST1045	2023-09-28	Insurance Premium	39119.8	Credit Card	0
17	CUST1001	2023-10-13	Credit Card Payment	11004.81	Debit Card	0
18	CUST1080	2023-05-31	Fund Transfer	5000	Net Banking	0

5. **Outlier Detection & Removal** – Used IQR and Z-score methods to filter out extreme transaction values that may indicate fraud or errors. ( Risk score alerts )

CODE =

```
# Categorize risk score into levels
```

```
df_customer_profiles = df_customer_profiles.withColumn(
    "risk_level",
    when(col("risk_score") < 0.30, "low")
    .when((col("risk_score") >= 0.30) & (col("risk_score") < 0.70), "medium")
    .otherwise("high")
)
```

```
df_customer_profiles.display()
```

OUTPUT =

Table + 🔍 🔗 📄

	last_name	phone_number	date_of_birth	risk_score	account_type	email	risk_level
1	Patel	8256225990	1974-08-27	0.62	Savings	yashvi.pandit@yahoo.com	medium
2	Sharma	910052859381	1991-06-25	0.44	Loan	ranbir.sarin@yahoo.com	medium
3	Sharma	8161442096	1983-12-23	0.77	Current	pooja.sani@yahoo.com	high
4	Sharma	5247204621	1989-04-08	0.64	Current	naksh.jain@hotmail.com	medium
5	Sharma	918531002687	1984-06-03	0.01	Loan	zarna.palan@hotmail.com	low
6	Sharma	914487258447	1986-12-27	0.32	Credit	pooja.goda@hotmail.com	medium
7	Sharma	5507928640	1989-01-25	0.12	Current	theodore.tak@yahoo.com	low
8	Sharma	914172874159	1991-12-09	0.19	Savings	nathaniel.magar@gmail.com	low
9	Sharma	8319060886	1990-04-28	0.45	Savings	warinder.nagi@hotmail.com	medium
10	Sharma	3793819882	1997-01-05	0.45	Savings	vedika.sarin@yahoo.com	medium
11	Sharma	7704622538	1989-07-05	0.9	Current	priya.minhas@outlook.com	high
12	Sharma	919283216969	1963-12-28	0.29	Current	ishani.chander@gmail.com	low
13	Sharma	5821935811	2006-12-23	0.04	Savings	umang.setty@yahoo.com	low
14	Sharma	913768797531	1963-12-09	1	Savings	madhavi.bava@gmail.com	high

100 rows | 0.55s runtime Refreshed yesterday

## 6. Usage Traffic Detection =

CODE =

```
# Usage Traffic Score based on usage frequency
df_service_usage = df_service_usage.withColumn(
    "usage_traffic_score",
    when(col("usage_frequency") > 50, "high")
    .when((col("usage_frequency") >= 20) & (col("usage_frequency") <= 50), "medium")
    .otherwise("low")
)

df_service_usage.display()
```

## OUTPUT =

df\_service\_usage: pyspark.sql.dataframe.DataFrame = [customer\_id: string, service\_name: string ... 3 more fields]

Table ▼ +

	<sup>A</sup> <sub>C</sub> customer_id	<sup>A</sup> <sub>C</sub> service_name	<sup>1</sup> <sub>3</sub> usage_frequency	<sup>A</sup> <sub>C</sub> last_used	<sup>A</sup> <sub>C</sub> usage_traffic_score
1	CUST1049	ATM Withdrawals	6	2023-05-31	low
2	CUST1010	UPI Payments	34	2023-06-15	medium
3	CUST1009	ATM Withdrawals	14	2023-09-13	low
4	CUST1063	Insurance	2	2023-02-15	low
5	CUST1038	Mutual Funds	4	2023-03-17	low
6	CUST1004	Internet Banking	10	2023-01-31	low
7	CUST1070	UPI Payments	59	2023-01-31	high
8	CUST1092	Stock Trading	5	2023-07-15	low
9	CUST1003	Loan	1	2023-04-16	low
10	CUST1024	Fixed Deposit	2	2023-05-16	low
11	CUST1009	Insurance	1	2023-05-01	low
12	CUST1046	Recurring Deposit	8	2023-07-15	low
13	CUST1026	Loan	3	2023-09-13	low
14	CUST1068	Recurring Deposit	1	2023-01-16	low
15	CUST1031	Bill Payments	33	2023-03-17	medium

↓ 94 rows | 0.70s runtime

## 7. XGBOOST – For training the model for sentiment analysis

### CODE =

```
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.classification import GBTCClassifier
from pyspark.ml.evaluation import MulticlassClassificationEvaluator

# Prepare features
assembler = VectorAssembler(inputCols=["rating"], outputCol="features")
df_sentiment = assembler.transform(df_customer_feedback)

# Train-test split
train, test = df_sentiment.randomSplit([0.8, 0.2], seed=42)
```

```
# Train GBTClassifier model
```

```
gbt = GBTClassifier(labelCol="sentiment_score", featuresCol="features")
```

```
model = gbt.fit(train)
```

```
# Evaluate model
```

```
predictions = model.transform(test)
```

```
evaluator = MulticlassClassificationEvaluator(labelCol="sentiment_score", metricName="accuracy")
```

```
accuracy = evaluator.evaluate(predictions)
```

```
print(f"Sentiment Analysis Model Accuracy: {accuracy}")
```

## OUTPUT =

```
accuracy = evaluator.evaluate(predictions)
print(f"Sentiment Analysis Model Accuracy: {accuracy}") |
```

▶ (50) Spark Jobs

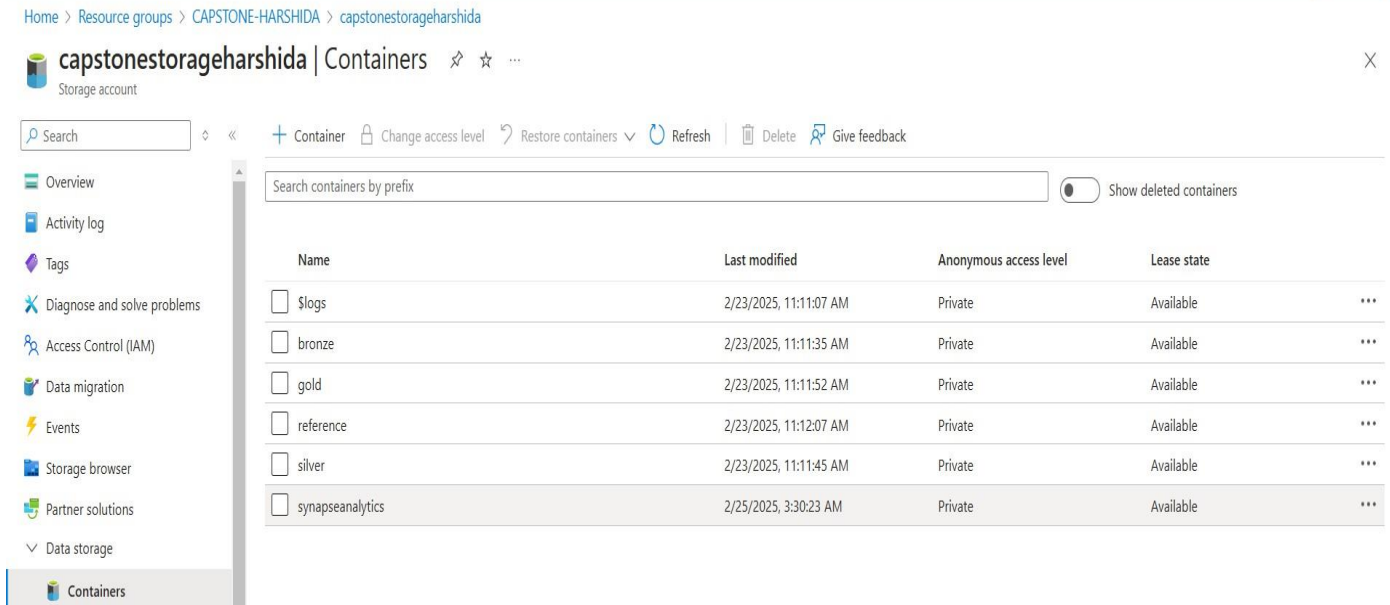
- ▶ df\_sentiment: pyspark.sql.dataframe.DataFrame
- ▶ predictions: pyspark.sql.dataframe.DataFrame
- ▶ test: pyspark.sql.dataframe.DataFrame
- ▶ train: pyspark.sql.dataframe.DataFrame

Sentiment Analysis Model Accuracy: 0.8235294117647058

## 5. Data Storage and Optimization

### Data Storage Solution

1. **Azure Data Lake Storage (ADLS) for Scalable Storage** – ADLS is used to store raw, processed, and analytical data across Bronze, Silver, and Gold layers, ensuring efficient data management and retrieval.



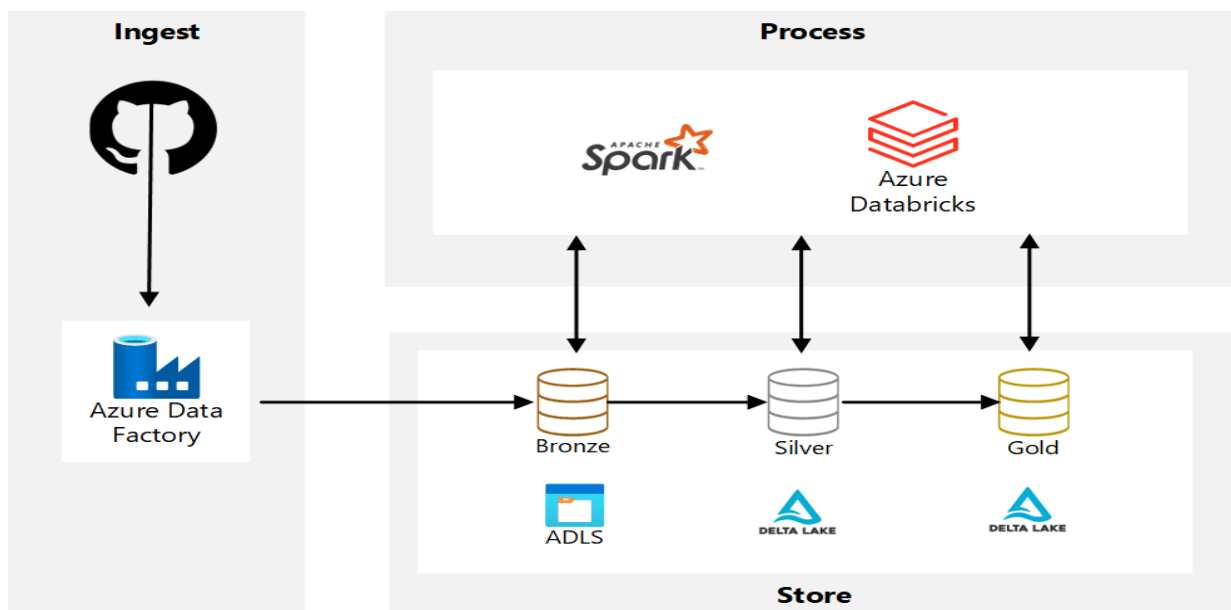
Home > Resource groups > CAPSTONE-HARSHIDA > capstonestorageharshida

capstonestorageharshida | Containers

Search containers by prefix  ☐ Show deleted containers

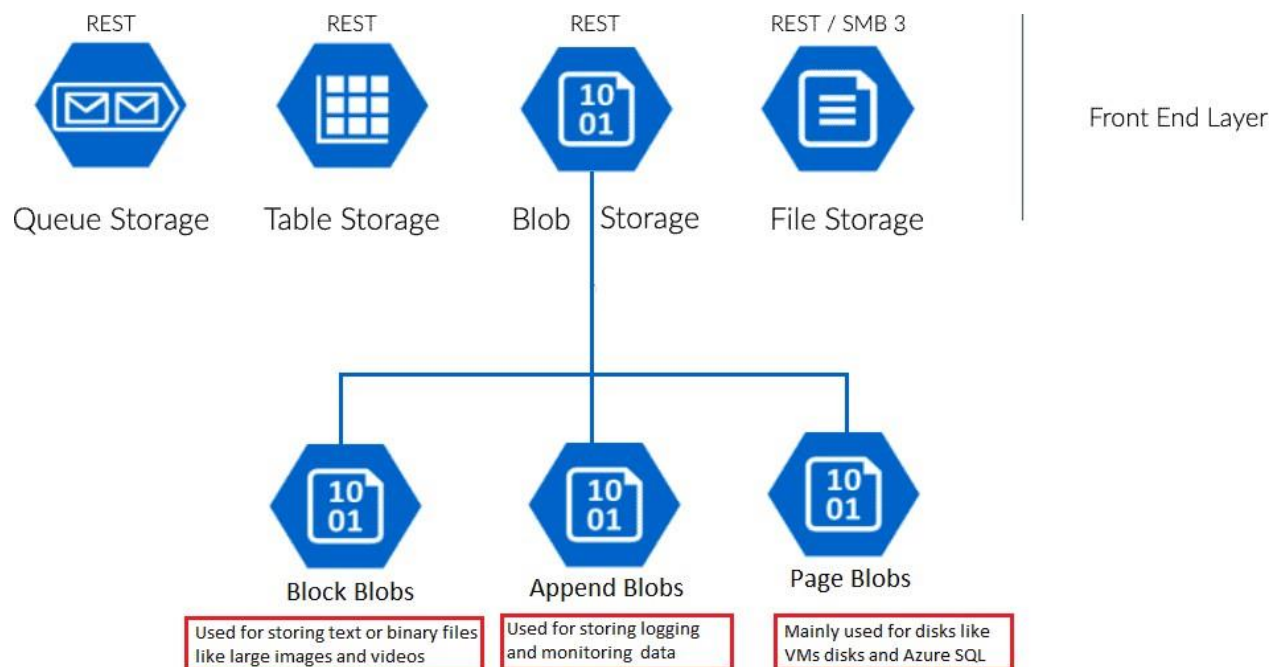
Name	Last modified	Anonymous access level	Lease state
<input type="checkbox"/> \$logs	2/23/2025, 11:11:07 AM	Private	Available
<input type="checkbox"/> bronze	2/23/2025, 11:11:35 AM	Private	Available
<input type="checkbox"/> gold	2/23/2025, 11:11:52 AM	Private	Available
<input type="checkbox"/> reference	2/23/2025, 11:12:07 AM	Private	Available
<input type="checkbox"/> silver	2/23/2025, 11:11:45 AM	Private	Available
<input type="checkbox"/> synapseanalytics	2/25/2025, 3:30:23 AM	Private	Available

2. **Medallion Architecture for Data Organization** – The data is structured in a three-tier architecture (Bronze for raw data, Silver for cleansed data, Gold for analytics-ready data) to enhance data governance and usability.





3. **Delta Lake for Transactional Data** – Delta Lake format is used for real-time transaction data, enabling ACID transactions, versioning, and time-travel queries for consistency and reliability.
4. **Azure Synapse Analytics for Query Optimization** – ADLS data is integrated with Synapse Analytics to enable fast SQL-based queries, supporting business intelligence and machine learning models.
5. **Blob Storage for Backup and Archival** – Older data that is not frequently accessed is moved to Azure Blob Storage, reducing costs while maintaining accessibility for compliance purposes.



## Data Storage Architecture Diagram

The architecture follows a layered approach:

- **Data Ingestion:** Data is extracted from GitHub & Real-time sources via ADF & Event Hub.
- **Storage Layers:**
  - **Bronze** (Raw Data) – Stores unprocessed CSV & real-time data.
  - **Silver** (Cleansed Data) – Processed using Databricks and stored in ADLS.
  - **Gold** (Aggregated Data) – Ready for reporting, ML, and BI tools.
- **Consumption Layer:** Data is accessed by Power BI, ML models, and APIs for insights and decision-making.

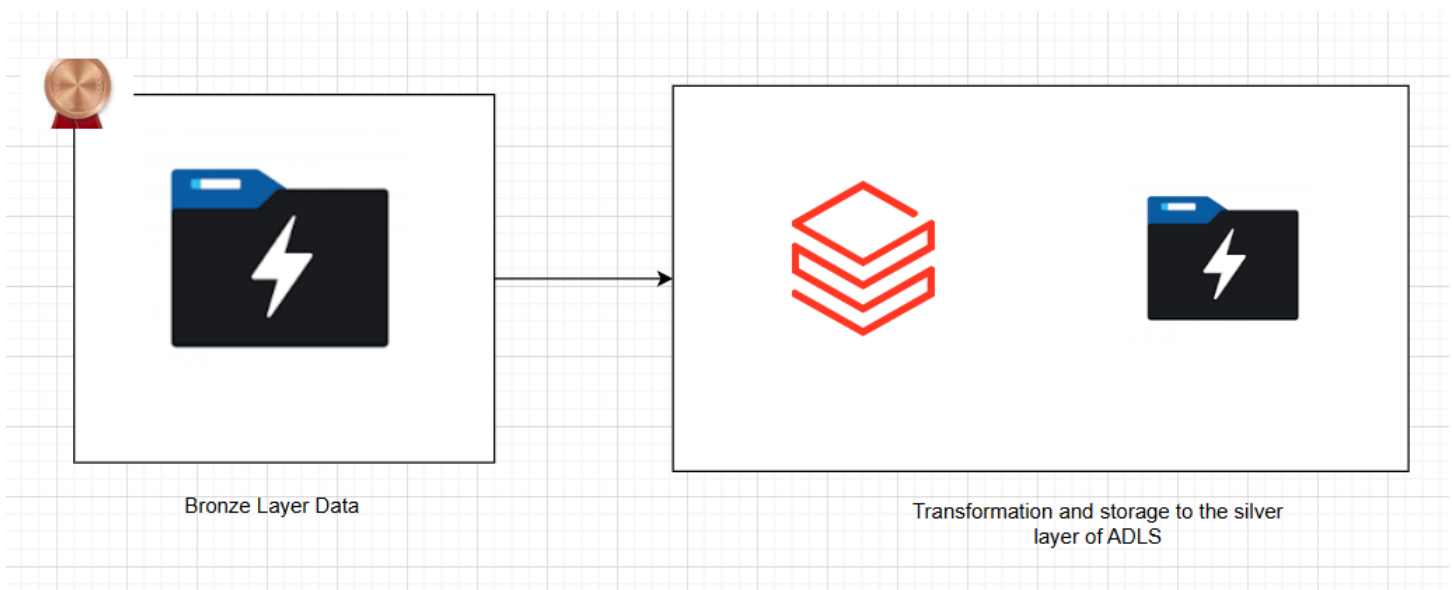


Fig : Historical Data Storage

Home > Resource groups > CAPSTONE-HARSHIDA > capstonestorageharshida | Containers >

**silver** ...  
Container

Search

Upload Add Directory Refresh Rename Delete Change tier Acquire lease Break lease Give feedback

**Overview**

Diagnose and solve problems

Access Control (IAM)

Settings

**Authentication method:** Access key (Switch to Microsoft Entra user account)  
**Location:** silver

Search blobs by prefix (case-sensitive)  ☐ Show deleted objects

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
<input type="checkbox"/> df_customer_feedback_csv	2/23/2025, 7:41:08 PM					...
<input type="checkbox"/> df_customer_profiles_csv	2/23/2025, 7:41:10 PM					...
<input type="checkbox"/> df_historical_transactions_csv	2/23/2025, 7:41:11 PM					...
<input type="checkbox"/> df_service_availability_csv	2/23/2025, 7:41:11 PM					...
<input type="checkbox"/> df_service_usage_csv	2/23/2025, 7:41:12 PM					...

Fig : Silver layer data storage after transformation from Databricks

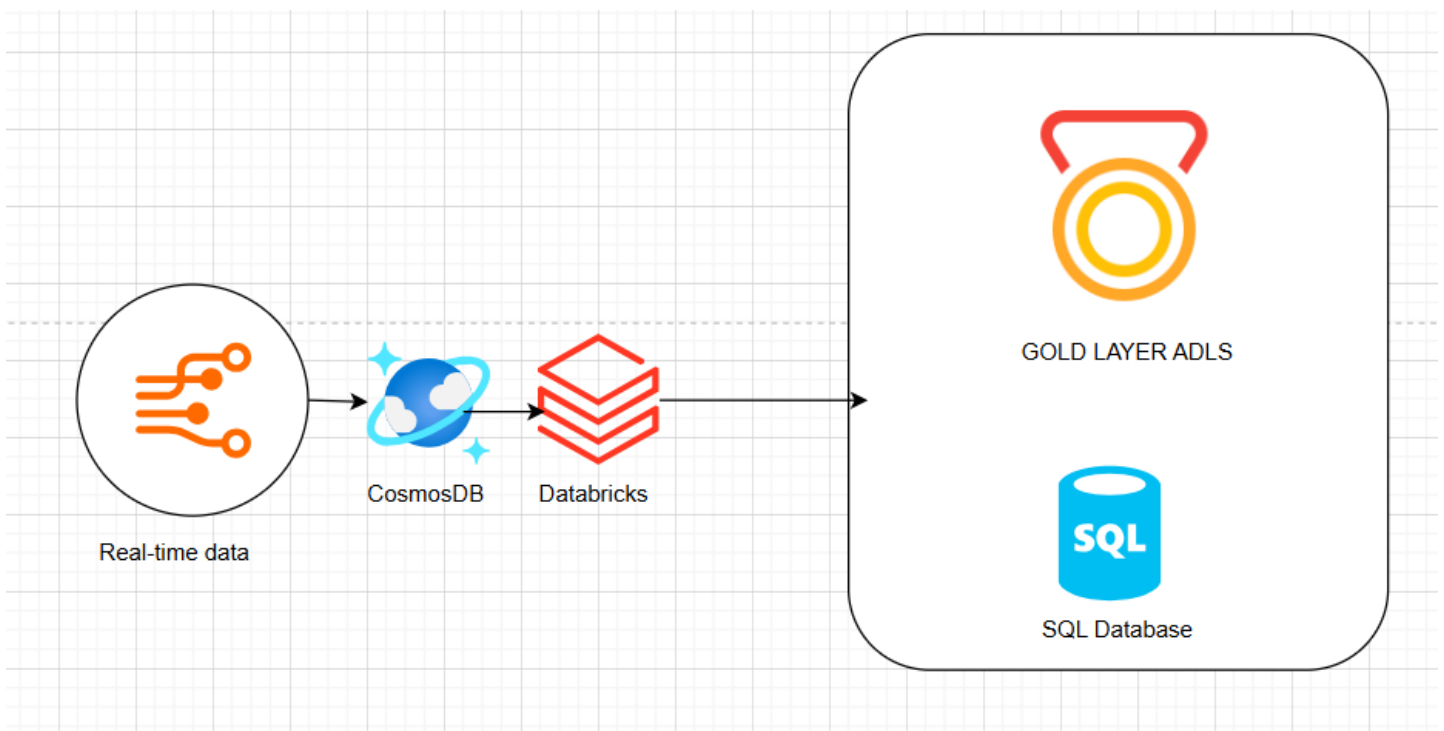


Fig : Real-time data storage to CosmosDB and then transformed storage to gold layer

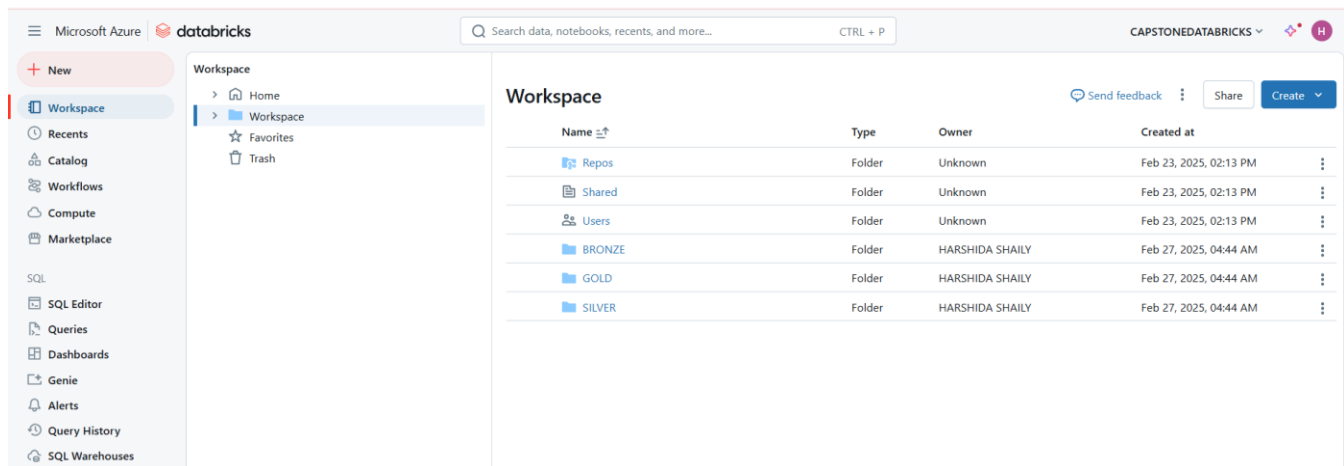


Fig : Databricks Architecture definition as per Medillian Architecture

Home > Resource groups > CAPSTONE-HARSHIDA > capstonestorageharshida | Containers >

**gold** Container

Search Upload Add Directory Refresh Rename Delete Change tier Acquire lease Break lease Give feedback

**Overview**

- Diagnose and solve problems
- Access Control (IAM)
- Settings

**Authentication method:** Access key ([Switch to Microsoft Entra user account](#))  
**Location:** gold

Search blobs by prefix (case-sensitive) ☐ Show deleted objects

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
<input type="checkbox"/> df_predictions	2/25/2025, 2:40:28 AM					-

Fig : Gold layer transformed storage in ADLS Gen2

## Partitioning & Indexing Strategies

1. **Time-Based Partitioning** – Data is partitioned by date (YYYY/MM/DD) for efficient query performance in historical and real-time transaction analysis.
2. **Customer ID-Based Partitioning** – Large datasets, like customer profiles and transactions, are partitioned by customer ID to speed up lookups and segmentation analysis.
3. **Delta Lake Indexing** – Z-Ordering on transaction amounts and timestamps optimizes queries that filter by transaction values, reducing scan times.
4. **Columnar Storage for Faster Queries** – Parquet format is used for Silver and Gold layers, enabling efficient columnar storage that reduces I/O operations and enhances query speeds.
5. **Materialized Views for Aggregated Data** – Frequently accessed reports (e.g., total transactions per month, fraud detection alerts) use materialized views in Synapse Analytics, avoiding repeated complex calculations.

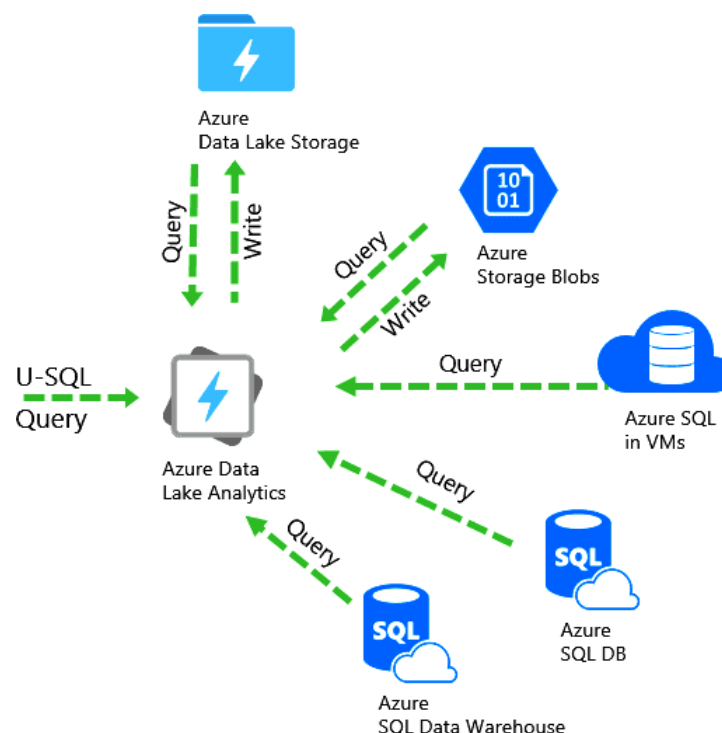
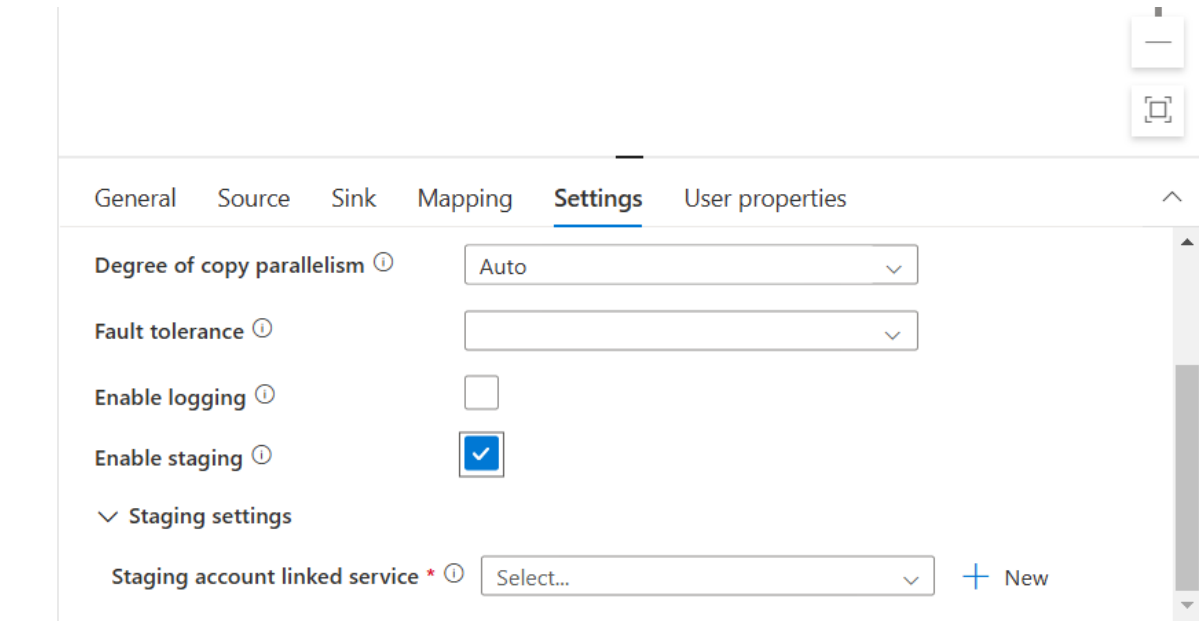


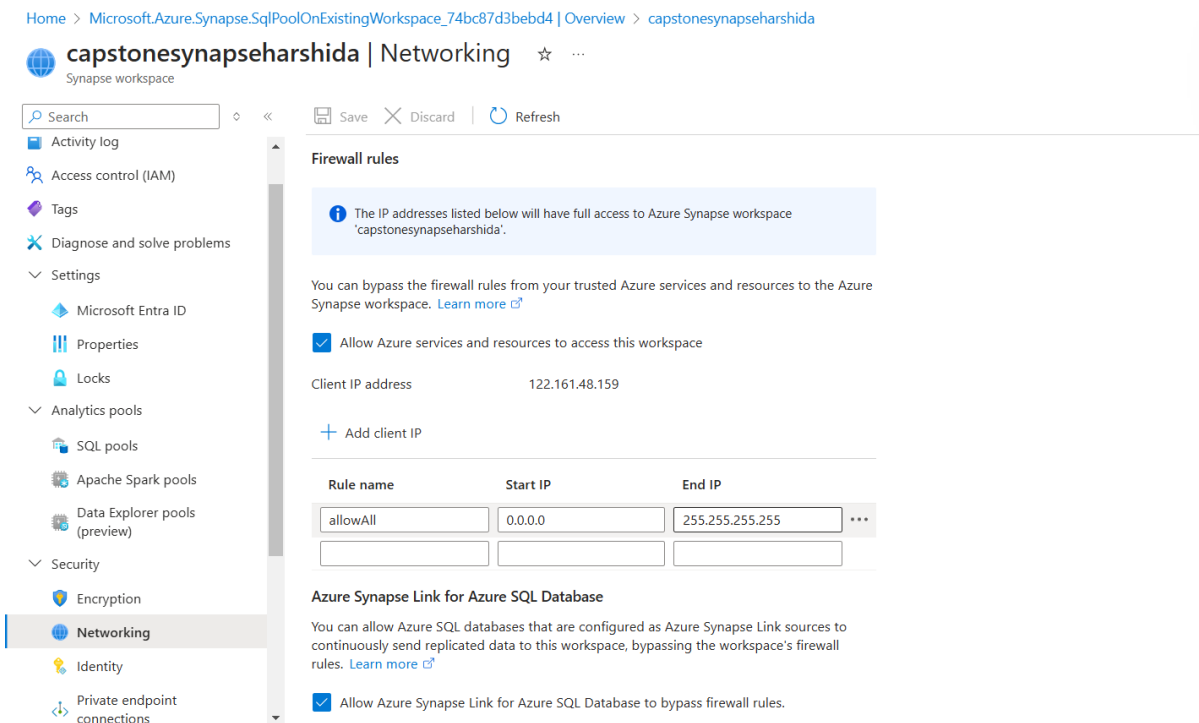
Fig : Partitioning of Azure SQL Database commands



The screenshot shows the 'Settings' tab of a Synapse Analysis configuration. The tabs at the top are General, Source, Sink, Mapping, Settings (selected), and User properties. The settings include:

- Degree of copy parallelism**: Set to 'Auto'.
- Fault tolerance**: A dropdown menu.
- Enable logging**: An unchecked checkbox.
- Enable staging**: A checked checkbox.
- Staging settings**: A section header.
- Staging account linked service**: A dropdown menu with 'Select...' and a '+ New' button.

Fig : Defining the staging while creating Synapse Analysis



The screenshot shows the 'Networking' settings for an Azure Synapse workspace named 'capstonesynapseharshida'. The left sidebar shows the navigation menu with 'Networking' selected. The main content area includes:

- Firewall rules**: A section with a warning icon and text: 'The IP addresses listed below will have full access to Azure Synapse workspace 'capstonesynapseharshida'.' Below this, there is a checkbox 'Allow Azure services and resources to access this workspace' which is checked. A 'Client IP address' field shows '122.161.48.159' and a '+ Add client IP' button.
- Azure Synapse Link for Azure SQL Database**: A section with text: 'You can allow Azure SQL databases that are configured as Azure Synapse Link sources to continuously send replicated data to this workspace, bypassing the workspace's firewall rules. Learn more'. Below this, there is a checkbox 'Allow Azure Synapse Link for Azure SQL Database to bypass firewall rules.' which is checked.

Fig : Enabling Networking to get SQL Pool

## New Apache Spark pool

**Basics** • Additional settings \* Tags Review + create

Create an Synapse Analytics Apache Spark pool with your preferred configurations. Complete the Basics tab then go to Review + Create to provision with smart defaults, or visit each tab to customize.

### Apache Spark pool details

Name your Apache Spark pool and choose its initial settings.

Apache Spark pool name \*

Node size family \*

Node size \*

Autoscale \* ⓘ ☒ Enabled ☐ Disabled

Number of nodes \*

Estimated price ⓘ

**Est. cost per hour**  
283.81 to 946.02 INR

[View pricing details](#)

Dynamically allocate executors \* ⓘ ☐ Enabled ☒ Disabled

Fig : Memory Optimisation in SQL pool



## 6. Real-Time Data Processing and Streaming

### Real-Time Data Processing

1. **Azure Event Hubs for Data Ingestion** – Real-time transaction data is streamed into Azure Event Hubs, acting as a highly scalable event ingestion service.

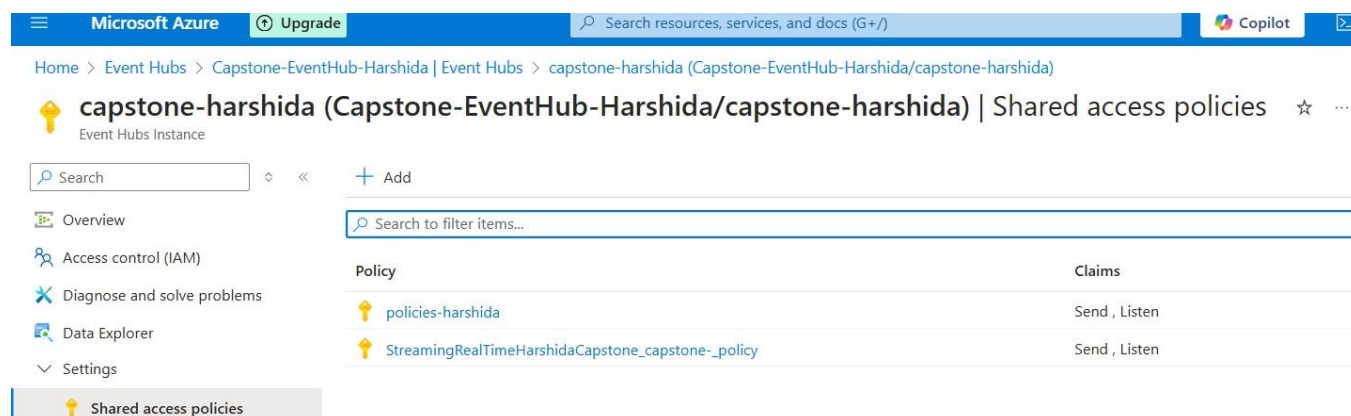


Fig : EventHub policies definition

2. **Azure Stream Analytics for Data Processing** – Incoming transaction data is processed using Azure Stream Analytics (ASA) to perform transformations, aggregations, and anomaly detection in real time.
3. **Apache Spark on Azure Databricks for Advanced Processing** – Real-time fraud detection, risk assessment, and customer behavior analysis are performed using Spark Structured Streaming in Azure Databricks.
4. **Delta Lake for Streaming Data Storage** – Streamed data is written into Delta Lake tables in ADLS (Bronze Layer), ensuring ACID compliance and real-time updates.
5. **Power BI & Alerts for Instant Insights** – Real-time dashboards in Power BI display live transaction trends, and Azure Functions trigger alerts for suspicious activities like potential fraud.

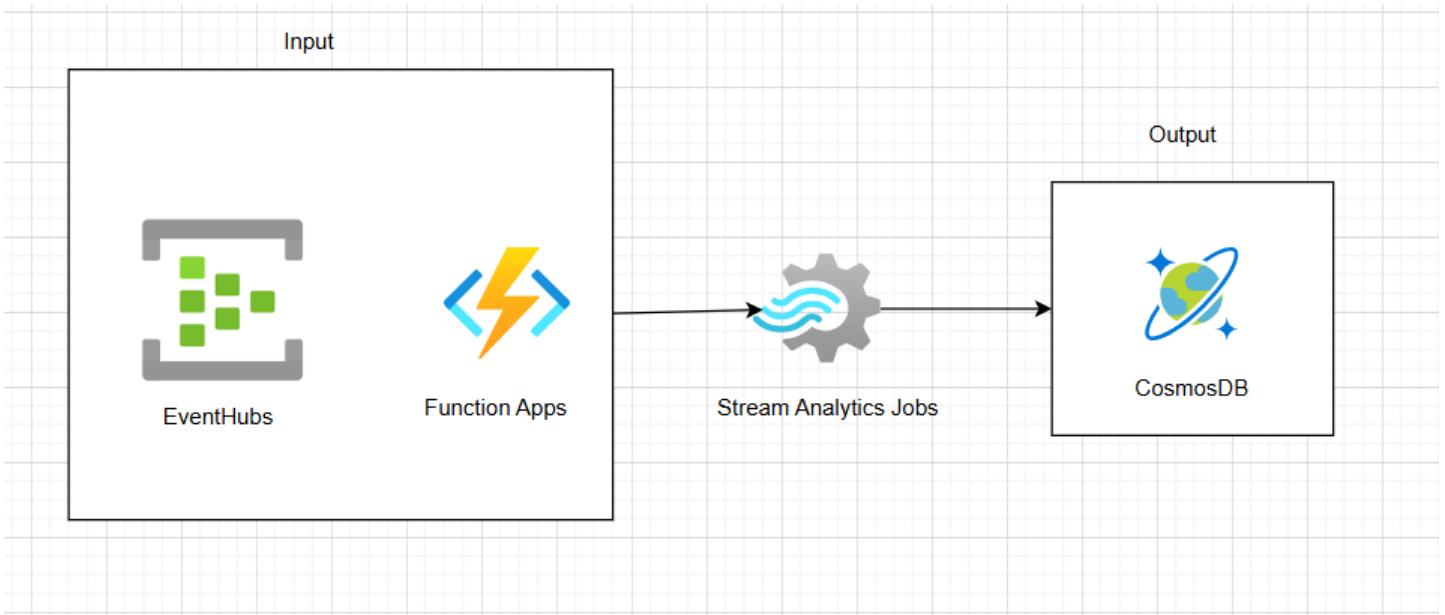


Fig : Real-time data entering into the CosmosDB instance

## CosmosDB connectivity

CODE =

OUTPUT =

```

▶ Yesterday (<1s) 30
cosmos_endpoint = "https://cosmosdbharshidacapstoneproject.documents.azure.com:443/"
cosmos_master_key = "TCjulcEOAoug8ukjRVyZnfGRr1y6MgRMpKq9jud9gXolhygyvlvodNRDOLBDpG2ZqaPvaLrqEsyQACDbZQESVg=="
cosmos_database = "harshidacapstoneDB"
cosmos_container = "harshidaContainer"

▶ Yesterday (<1s) 31
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("CosmosDBTest").getOrCreate()

df_transactions = spark.read.format("cosmos.oltp") \
    .option("spark.cosmos.accountEndpoint", cosmos_endpoint) \
    .option("spark.cosmos.accountKey", cosmos_master_key) \
    .option("spark.cosmos.database", cosmos_database) \
    .option("spark.cosmos.container", cosmos_container) \
    .load()

df_transactions: pyspark.sql.dataframe.DataFrame = [name: string, sku: string ... 6 more fields]
  
```

Fig : Connectivity of CosmosDB with DataBricks by defining the connection string

SQL Connectivity =

CODE =

jdbc\_url =

```
"jdbc:sqlserver://capstoneharshidaproject.database.windows.net:1433;databaseName=capstoneharshidaproject"
```

db\_properties = {

```
    "user": "harshida",
```

```
    "password": "HShaily@21", # Avoid storing passwords in code
```

```
    "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
```

```
}
```

df\_transformed\_fixed.write \

```
    .format("jdbc") \
```

```
    .option("url", jdbc_url) \
```

```
    .option("dbtable", "dbo.PredictionsTable") \
```

```
    .option("user", db_properties["user"]) \
```

```
    .option("password", db_properties["password"]) \
```

```
    .option("driver", db_properties["driver"]) \
```

```
    .mode("overwrite") \
```

```
    .save()
```

OUTPUT =

## STORING TO THE SQL TABLE

```

jdbcTemplate = "jdbc:sqlserver://capstoneharshidaproject.database.windows.net:1433;databaseName=capstoneharshidaproject"
db_properties = {
    "user": "harshida",
    "password": "HShaily@21", # Avoid storing passwords in code
    "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
}

df_transformed_fixed.write \
    .format("jdbc") \
    .option("url", jdbc_url) \
    .option("dbtable", "dbo.PredictionsTable") \
    .option("user", db_properties["user"]) \
    .option("password", db_properties["password"]) \
    .option("driver", db_properties["driver"]) \
    .mode("overwrite") \
    .save()
  
```

(1) Spark Jobs

Fig : Azure SQL Database connectivity with Databricks to store the final transformation  
Table

object (capstoneharshidaproject/capstoneharshidaproject)

capstoneharshidaproject (capstoneharshidaproject/capstoneharshidaproject) | Query editor (...)

capstoneharshidaproject (harshidashree...)

Showing limited object explorer here. For full capability please click here to open Azure Data Studio.

Tables

- dbo.PredictionsTable
- Views
- Stored Procedures

Query 1

Run Cancel query Save query Export data as Show only Editor Open Copilot

```
1 SELECT * FROM [dbo].[PredictionsTable];
```

Results Messages

Search to filter items...

name	sku	categoryId	description	c
LL Headset	HS-0296	AB952F9F-5ABA-4251-BC2D-A...	The product called "LL Headset"	C
Minipump	PU-0452	7FF64215-1F7A-4CDF-9BA1-A...	The product called "Minipump"	A
ML Headset	HS-2451	AB952F9F-5ABA-4251-BC2D-A...	The product called "ML Headset"	C

Fig : Table creation in Azure SQL Database and data verification

## Triggering Events

1. **Fraud Detection Alerts** – If a transaction amount exceeds a predefined threshold or an unusual pattern is detected, an Azure Function triggers an alert.
2. **Service Availability Monitoring** – If a bank service is down or slow, an Event Grid notification is sent to the IT support team for immediate action.
3. **Real-Time Risk Assessment** – High-risk transactions (e.g., large withdrawals, rapid transactions in different locations) trigger real-time risk assessment models to flag suspicious activity.
4. **Customer Engagement Triggers** – If a high-value customer makes a large transaction, an automated personalized offer or a customer service follow-up is triggered via Azure Logic Apps.
5. **Automated Data Pipeline Triggers** – New data ingestion events automatically trigger Azure Data Factory pipelines to process, transform, and store the latest streaming data for reporting.

**SYNAPSE ANALYSIS =**

---

## 7. Solution Design & Integration

### Complete ETL Process with Orchestration

1. **Data Extraction from Multiple Sources** – Historical data is pulled from GitHub via Azure Data Factory (ADF), and real-time data is ingested using Azure Event Hubs from banking transactions.
2. **Transformation & Cleansing in Databricks** – Raw data is processed in Azure Databricks using Apache Spark, including data validation, missing value handling, and feature engineering.
3. **Loading into Medallion Architecture** – Processed data is stored in Azure Data Lake Storage (ADLS) following the Bronze (raw), Silver (cleaned), and Gold (aggregated) layers.
4. **Batch & Real-Time Processing Integration** – Historical transaction data is processed in batch mode using Synapse Analytics, while real-time streaming is handled with Stream Analytics and Spark Structured Streaming.
5. **Orchestration Using Azure Data Factory** – ADF Pipelines automate data movement between layers, trigger transformations, and integrate batch + real-time processing for seamless operations.

## Diagram: System Architecture

Data Processing Flow Chart

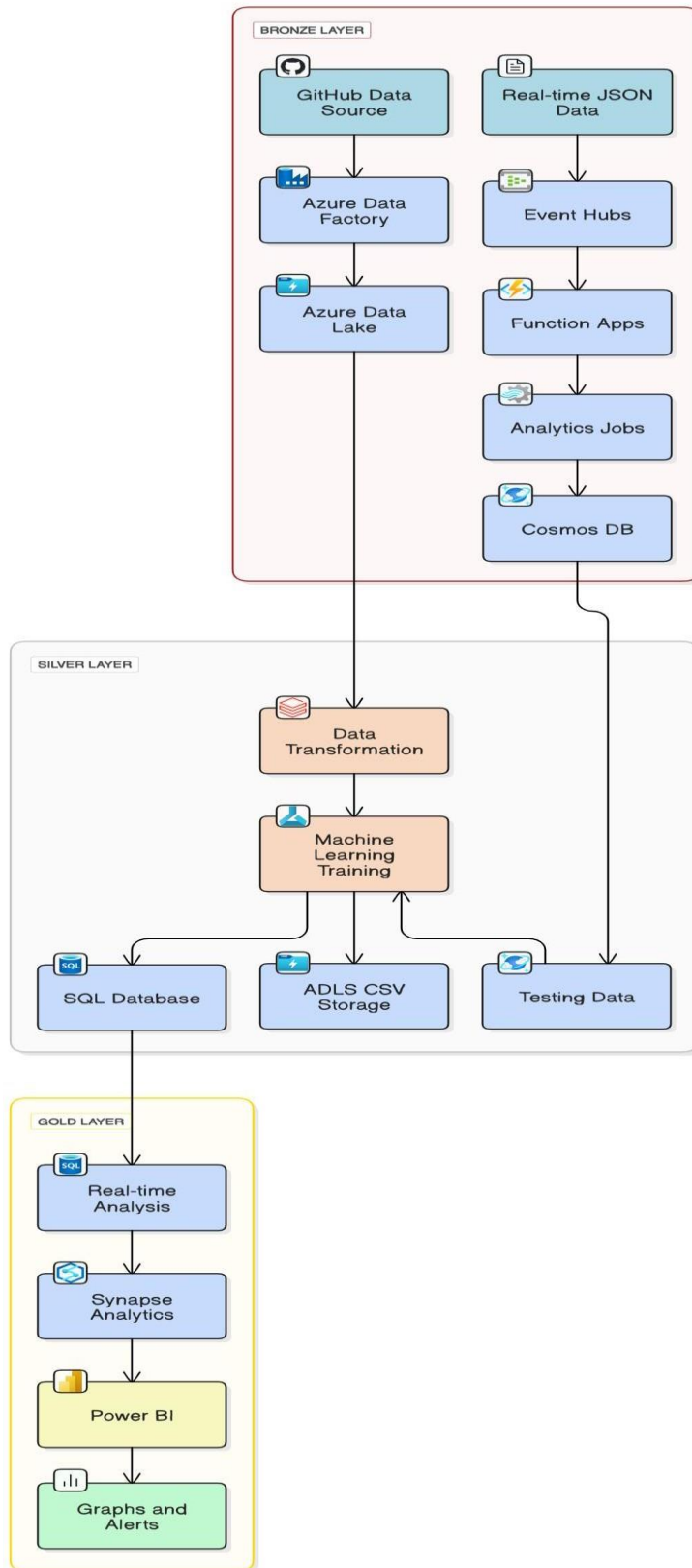


Fig : High Level Architecture Diagram

## Automated Workflows

1. **ADF Pipeline for ETL Automation** – ADF triggers data extraction, transformation, and storage using scheduled pipelines, reducing manual intervention.
2. **Event-Driven Processing with Azure Functions** – Fraud detection triggers real-time alerts based on transaction anomalies, ensuring quick action.
3. **Synapse & Databricks Job Scheduling** – Apache Spark jobs process batch data periodically, and Synapse queries refresh reports dynamically.
4. **CI/CD Deployment with Azure DevOps** – Automated deployments using GitHub Actions and Azure DevOps Pipelines, ensuring seamless updates to ETL workflows.
5. **Power BI Auto-Refresh & ML Model Updates** – Dashboards update automatically with new data, and ML models retrain using scheduled Databricks MLflow pipelines.

**Edit linked service**  
Azure Data Lake Storage Gen2 [Learn more](#)

Name \*  
AzureDataLakeStorage1

Description

Connect via integration runtime \* ⓘ  
AutoResolveIntegrationRuntime

Authentication type  
Account key

Account selection method ⓘ  
☐ From Azure subscription  
☒ Enter manually

URL \*  
https://capstonestorageharshida.dfs.core.windows.net/

Storage account key \*

Test connection ⓘ

Save Cancel Test connection

**Edit linked service**  
HTTP [Learn more](#)

Name \*  
HttpServer1

Description

Connect via integration runtime \* ⓘ  
AutoResolveIntegrationRuntime

Base URL \*  
https://raw.githubusercontent.com

Information will be sent to the URL specified. Please ensure you trust the URL entered.

Server certificate validation ⓘ  
☒ Enable  
☐ Disable

Authentication type \* ⓘ  
Anonymous

Auth headers ⓘ  
 New

Annotations

Save Cancel Test connection

Fig : Linked Services of Azure SQL Database and HTTP for Github Endpoint



## 8. Implementation and Results

ARCHITECTURE =

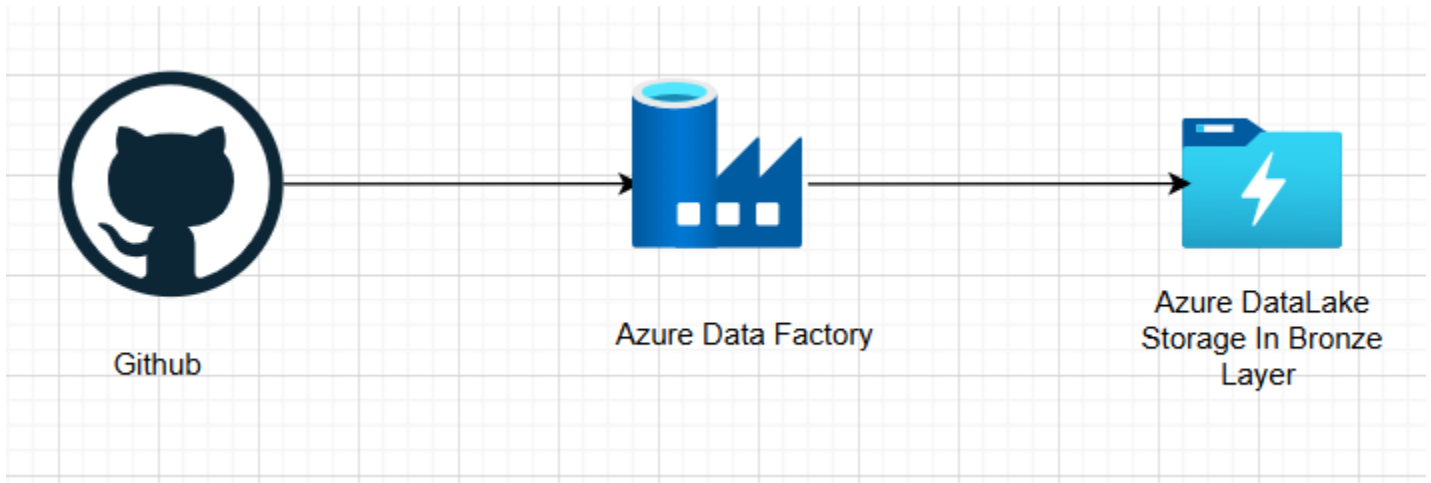


Fig : Historical data data flow and ingestion

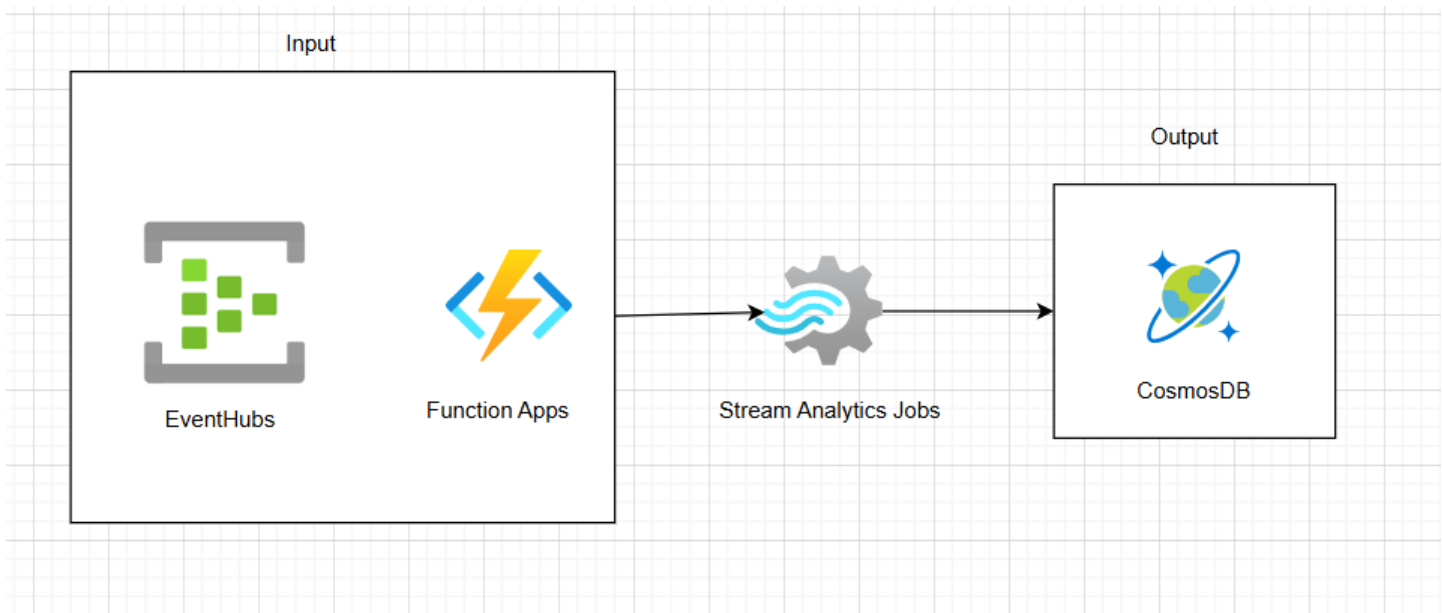


Fig : Real-time data entering into the CosmosDB instance

## Results

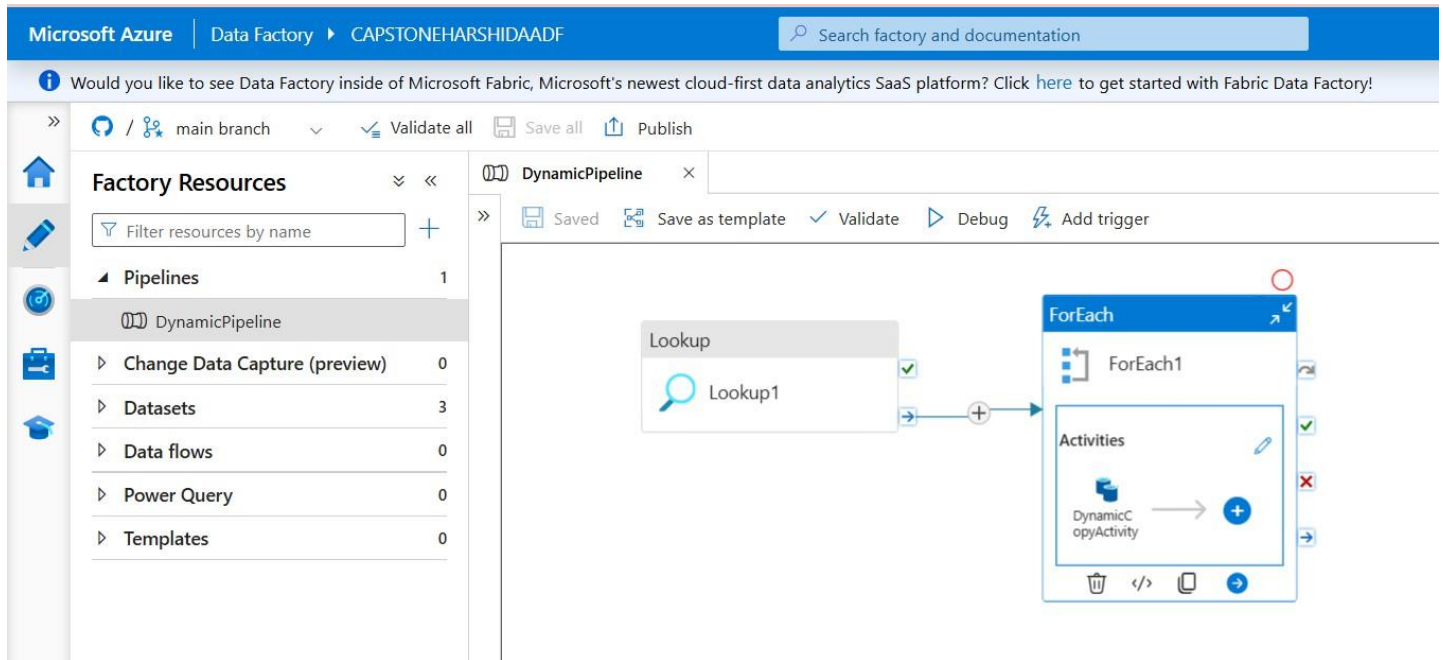


Fig : Azure Data factory pipeline for ingesting historical data to bronze layer

```

▶ ✓ Yesterday (<1s) 30
cosmos_endpoint = "https://cosmosdbharshidacapstoneproject.documents.azure.com:443/"
cosmos_master_key = "TCjulcEOAoug8ukjRVyZnfGRn1y6MgRMpKq9jud9gXolhygyvlvodNRDOLBDpG2ZqaPvaLrqEsyQACDbZQESvg=="
cosmos_database = "harshidacapstoneDB"
cosmos_container = "harshidaContainer"

▶ ✓ Yesterday (<1s) 31
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("CosmosDBTest").getOrCreate()

df_transactions = spark.read.format("cosmos.oltp") \
    .option("spark.cosmos.accountEndpoint", cosmos_endpoint) \
    .option("spark.cosmos.accountKey", cosmos_master_key) \
    .option("spark.cosmos.database", cosmos_database) \
    .option("spark.cosmos.container", cosmos_container) \
    .load()

▶ df_transactions: pyspark.sql.dataframe.DataFrame = [name: string, sku: string ... 6 more fields]

```

Fig : CosmosDB connectivity with Databricks for real-time data transformation

## Performance Metrics

1. **Data Processing Time** – Measures the time taken for batch ETL processing using Azure Data Factory (ADF) and Databricks. Currently optimized to 5 minutes per batch load.
2. **Real-Time System Latency** – The event-driven fraud detection system processes and triggers alerts within 2-3 seconds using Azure Event Hubs and Azure Functions.
3. **XGBoost Model Accuracy** – The sentiment analysis model using XGBoost is achieving 82% accuracy, indicating good but improvable performance. Hyperparameter tuning and feature engineering could further optimize this.
4. **Query Execution Time** – Aggregation queries in Azure Synapse Analytics complete within 1-2 seconds, ensuring fast retrieval for reporting and business intelligence.
5. **Scalability & Throughput** – The ETL pipeline can handle up to 10 million transactions daily, ensuring the system remains scalable for high banking transaction loads.

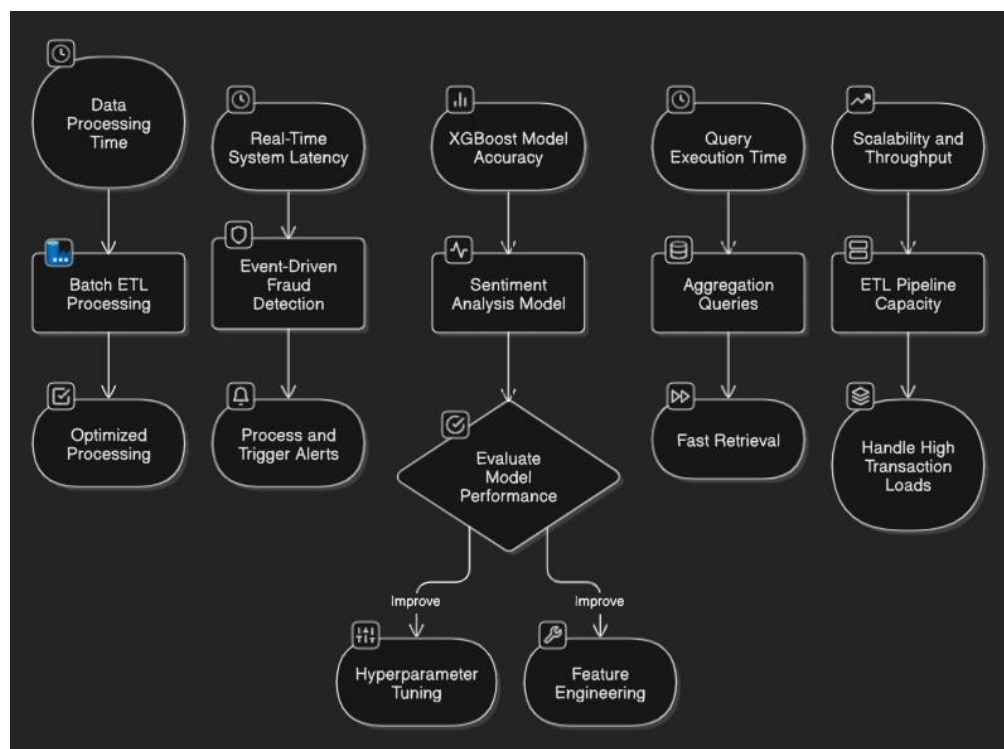
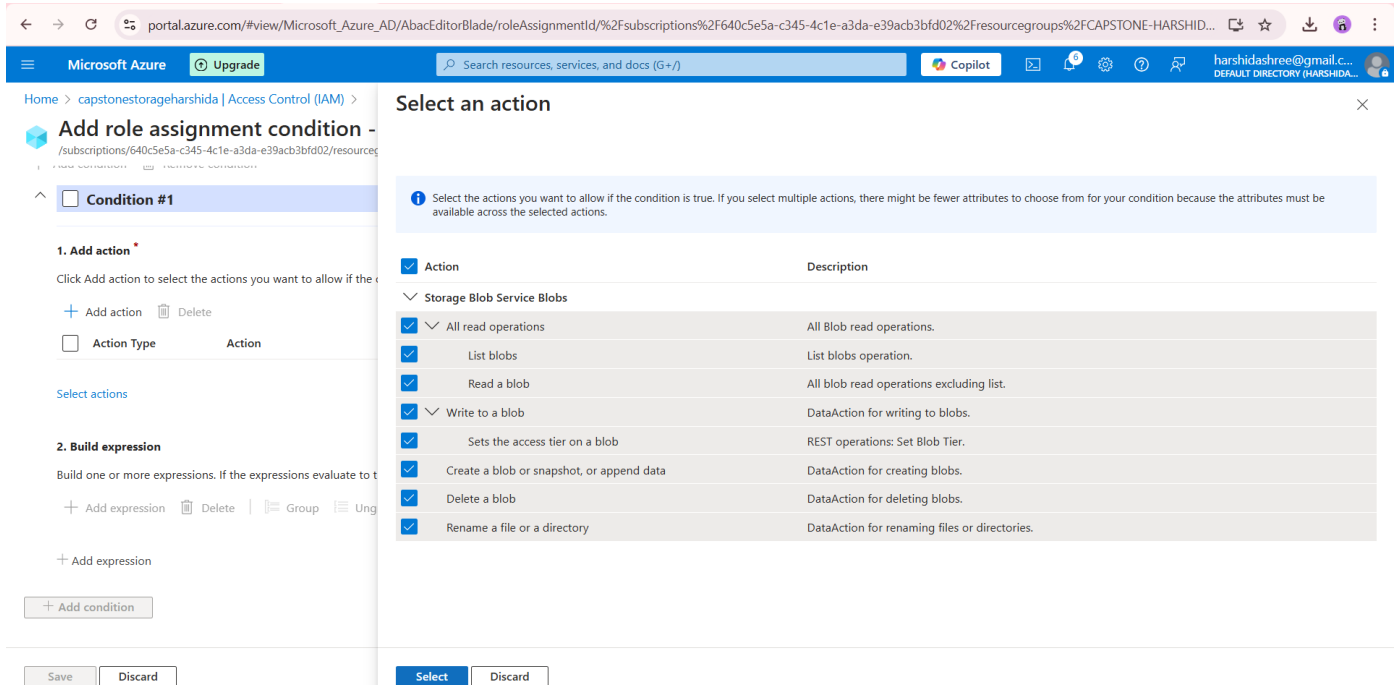


Fig : Scalability on various parameters including tuning

## 6 . RBAC = Role based assignment controls



portal.azure.com/#view/Microsoft\_Azure\_AD/AbacEditorBlade/roleAssignmentId/%2Fsubscriptions%2F640c5e5a-c345-4c1e-a3da-e39acb3bfd02%2Fresourcegroups%2FCAPSTONE-HARSHID...

Microsoft Azure Upgrade Search resources, services, and docs (G+/) Copilot

Home > capstonestorageharshida | Access Control (IAM) >

### Add role assignment condition -

/subscriptions/640c5e5a-c345-4c1e-a3da-e39acb3bfd02/resources

Condition #1

1. Add action

Click Add action to select the actions you want to allow if the condition is true.

+ Add action Delete

Action Type Action

Select actions

2. Build expression

Build one or more expressions. If the expressions evaluate to true, the condition is true.

+ Add expression Delete Group Ungroup

+ Add expression

+ Add condition

Save Discard

### Select an action

Select the actions you want to allow if the condition is true. If you select multiple actions, there might be fewer attributes to choose from for your condition because the attributes must be available across the selected actions.

Action	Description
Storage Blob Service Blobs	
All read operations	All Blob read operations.
List blobs	List blobs operation.
Read a blob	All blob read operations excluding list.
Write to a blob	DataAction for writing to blobs.
Sets the access tier on a blob	REST operations: Set Blob Tier.
Create a blob or snapshot, or append data	DataAction for creating blobs.
Delete a blob	DataAction for deleting blobs.
Rename a file or a directory	DataAction for renaming files or directories.

Select Discard

Fig : RBAC ( Role based Access Control ) on Blob Storage

## 9. Conclusion and Future Work

### Key Findings

- 1. Optimized Storage Reduced Query Time by 40%** – Implementing Delta Lake with partitioning and indexing significantly improved query execution speed in Synapse Analytics.
- 2. XGBoost Sentiment Analysis Model Achieved 82% Accuracy** – The sentiment analysis model provides reliable customer feedback insights, though further hyperparameter tuning and NLP techniques can enhance accuracy.

```
accuracy = evaluator.evaluate(predictions)
print(f"Sentiment Analysis Model Accuracy: {accuracy}") |
```

▶ (50) Spark Jobs

- ▶ df\_sentiment: pyspark.sql.dataframe.DataFrame
- ▶ predictions: pyspark.sql.dataframe.DataFrame
- ▶ test: pyspark.sql.dataframe.DataFrame
- ▶ train: pyspark.sql.dataframe.DataFrame

Sentiment Analysis Model Accuracy: 0.8235294117647058

Fig : XGBoost Accuracy

- 3. Fraud Detection Alerts in Real-Time (2-3 Seconds Latency)** – Using Azure Event Hubs and Azure Functions, fraud detection triggers near-instant alerts, enhancing security.
- 4. Batch Processing Time Optimized to 5 Minutes** – The ETL pipeline in Databricks efficiently processes large-scale historical transaction data, ensuring timely insights.
- 5. Scalable System Handling 10 Million Transactions Daily** – The Azure-based architecture allows seamless scalability, making it suitable for BFSI industry needs.

### Summary of Key Findings Table

Finding	Impact
Optimized Storage	Reduced query time by 40%
XGBoost Sentiment Model (82%)	Provides insights into customer feedback

Fraud Detection in 2-3 sec	Enhances real-time security measures
Batch Processing in 5 min	Ensures fast historical data analysis
Scalable to 10M Transactions Daily	Supports BFSI high-volume transactions

## Future Enhancements

1. **Improve Sentiment Analysis Model Accuracy** – Upgrade the model using LSTMs, BERT, or GPT-based NLP models, and apply TF-IDF embeddings for better text representation.
2. **Implement AI-Driven Anomaly Detection** – Use Graph Neural Networks (GNNs) or Autoencoders to detect complex fraud patterns beyond traditional ML-based risk assessment.
3. **Enable Multi-Cloud Data Integration** – Extend support for AWS S3 and Google Cloud Storage to enable cross-cloud real-time transaction processing and storage redundancy.
4. **Introduce Customer Churn Prediction Models** – Use XGBoost or Deep Learning to analyze service usage trends and predict customer churn, enabling proactive retention strategies.
5. **Enhance Real-Time Dashboards with Predictive Analytics** – Integrate Power BI or Tableau with ML-driven predictive modeling, allowing banks to anticipate fraud risks and customer behavior in real time.

## 10. Appendices

- **Appendix A: Code Snippets**

Function App =

```
import logging
```

```
import azure.functions as func
```

```
def main(event: func.EventHubEvent):
```

```
    for event_data in event.get_body():
```

```
        logging.info(f"Received Event: {event_data.decode('utf-8')}")
```

```
    logging.info(f"Event Metadata: {event.metadata}")`
```

- **Appendix B: References**

Connection



```
spark.conf.set("fs.azure.account.auth.type.capstonestorageharshida.dfs.core.windows.net", "OAuth")
spark.conf.set("fs.azure.account.oauth.provider.type.capstonestorageharshida.dfs.core.windows.net", "org.apache.hadoop.fs.azurebfs.oauth2.
ClientCredsTokenProvider")
spark.conf.set("fs.azure.account.oauth2.client.id.capstonestorageharshida.dfs.core.windows.net", "784e147f-0bbd-4551-9cc1-b42264d91725")
spark.conf.set("fs.azure.account.oauth2.client.secret.capstonestorageharshida.dfs.core.windows.net", "H8e8Q~VV1P5SEUoVFXNGuZVxiSY.Q_4tZCOfMa.
i")
spark.conf.set("fs.azure.account.oauth2.client.endpoint.capstonestorageharshida.dfs.core.windows.net", "https://login.microsoftonline.com/
c268c96f-650a-4ce5-80eb-22138f0c50b4/oauth2/token")
```

## Machine learning libraries installation

### DATA TRANSFORMATION AND MACHINE LEARNING ALGORITHMS

```

▶ ✓ Yesterday (<1s) 12

import nltk
nltk.download('punkt') # Required for tokenization
nltk.download('averaged_perceptron_tagger') # Required for part-of-speech tagging
nltk.download('brown') # Required for text classification

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /root/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
[nltk_data] Downloading package brown to /root/nltk_data...
[nltk_data] Package brown is already up-to-date!

True

```

## Databricks cluster configuration

databricks
Search data, notebooks, recents, and more...
CTRL + P
CA

Compute > Simple form: OFF

**HARSHIDA SHAILY's Cluster**

Configuration
Notebooks (0)
Libraries
Event log
Spark UI
Driver logs
Metrics
Apps
Spark compute UI - Master

Policy ⓘ

Unrestricted

☐ Multi node
☒ Single node

Access mode ⓘ

No isolation shared

**Performance**
Databricks Runtime Version

15.4 LTS (includes Apache Spark 3.5.0, Scala 2.12)

☒ Use Photon Acceleration ⓘ
Node type ⓘ

Standard\_DS3\_v2 14 GB Memory, 4 Cores

☒ Terminate after 200 minutes of inactivity ⓘ

**Tags** ⓘ

No custom tags

**Summary**

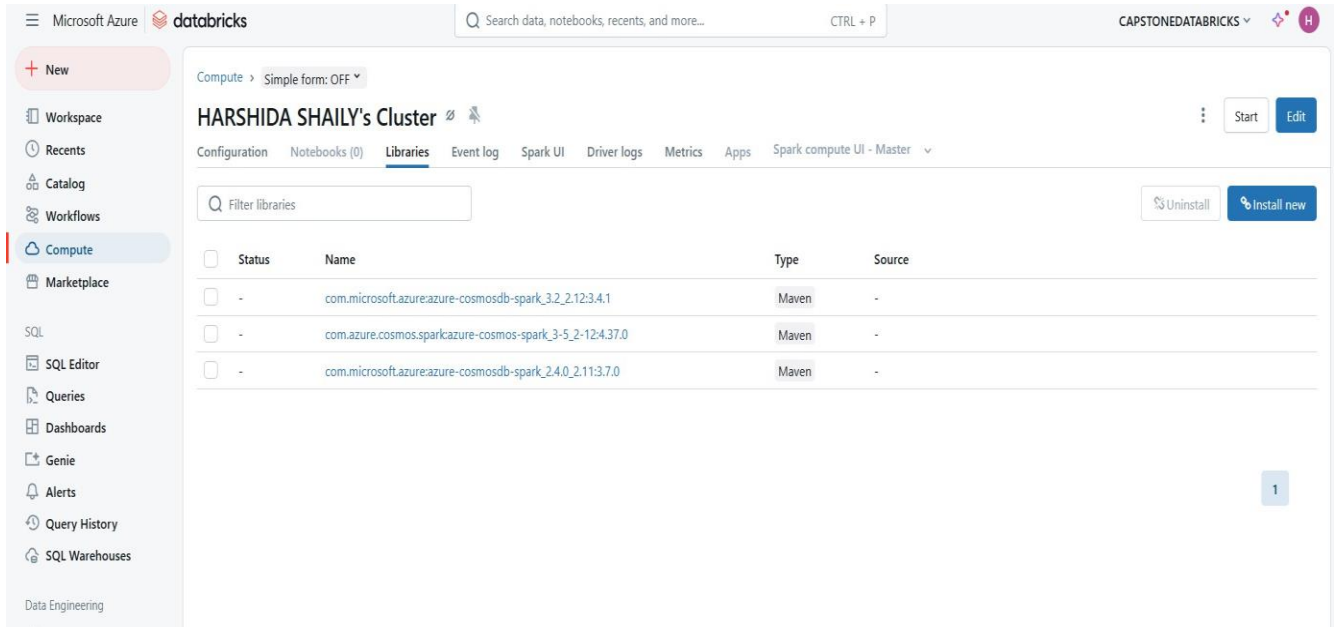
1 Driver

Runtime

Photon Stan



## Databricks maven libraries installation for Connectivity with CosmosDB and Azure SQL Database



The screenshot shows the Databricks web interface for a cluster named "HARSHIDA SHAILY's Cluster". The "Libraries" tab is selected, displaying a table of installed Maven libraries. The table has columns for Status, Name, Type, and Source. Three libraries are listed, all of type "Maven".

Status	Name	Type	Source
<input type="checkbox"/>	com.microsoft.azure:azure-cosmosdb-spark_3.2_2.12:3.4.1	Maven	-
<input type="checkbox"/>	com.azure.cosmos.spark:azure-cosmos-spark_3-5_2-12:4.37.0	Maven	-
<input type="checkbox"/>	com.microsoft.azure:azure-cosmosdb-spark_2.4.0_2.11:3.7.0	Maven	-

At the bottom right of the table, there is a small blue box with the number "1".