# Azure Data Pipeline Project Report

Nora Ayaz

October 10, 2024

# Contents

# 1   Initial Setup and Configuration

Before implementing the cloud-based data pipeline, several initial configurations were required on the Azure Portal. These steps were crucial for setting up the resources that enabled the pipeline's functionality. The following key components were created and configured:

## 1.1   Azure Student Account

An Azure Student account was created to gain access to a free tier of Azure services, which includes free credits and access to essential services like Azure Storage, Azure SQL Database, and Function Apps.

## 1.2   Resource Group

A Resource Group was created to organize and manage all related Azure resources under one grouping. This helps with administration, cost tracking, and resource management.

## 1.3   Azure Storage Account

An Azure Storage Account was provisioned to serve as the backbone for the queue service. Within the storage account:

- A **Queue Service** was created to hold sensor data messages sent by the data producer.

- The **Queue** was named `myqueue1090` to store and queue incoming data for processing.

## 1.4   Azure SQL Database and Server

An Azure SQL Database was set up to store the sensor data processed by the Azure Function. The following steps were taken:

- A **SQL Server** was created to manage the database. The server was configured with appropriate firewall rules to allow access from specific IP addresses.

- A **SQL Database** named `mydatabase1090` was created to store the processed sensor data (temperature, humidity, and timestamp).

- The database was set to the **serverless tier** to enable cost-effective scaling and automatic pausing during periods of inactivity.

## 1.5   Azure Function App

An **Azure Function App** was created to handle the processing of data messages from the queue. The function app automatically triggers whenever new data arrives in the queue, processes the message, and inserts the data into the Azure SQL Database.

## 1.6 Power BI and Grafana Setup

To visualize the stored data:

- **Power BI** was connected to the Azure SQL Database to create reports and visualizations based on the stored sensor data.

- **Grafana** was installed on a cloud-hosted environment and configured to connect to the Azure SQL Database, providing real-time data dashboards.

This initial setup allowed for the creation of a seamless data pipeline from data generation at the edge, through processing and storage in the cloud, to visualization using various tools.

# 2 Introduction

The goal of this project was to design and implement a serverless data pipeline using Azure services, demonstrating cloud computing skills in data generation, ingestion, storage, and consumption. The project simulates a real-world IoT-like system where data is generated at the edge, sent to the cloud for storage, and later consumed by data consumers for analysis and visualization.

The key components of this project include:

**Data Producer:** A Python script (`send_to_queue.py`) generates random sensor data (temperature, humidity, and timestamps) and sends it to an Azure Queue. This simulates the edge data generation and transmission to the cloud.

**Message Queue:** Azure Queue Storage is used as the queuing service to temporarily store the sensor data messages from the producer before processing.

**Data Ingestion:** An Azure Function (`function_app.py`) is triggered automatically when new data arrives in the Azure Queue. This serverless function processes the data and inserts it into an Azure SQL Database, ensuring that the data is ingested and stored efficiently.

**Data Storage:** The ingested data is stored in an Azure SQL Database (serverless tier), which provides a scalable and reliable cloud-based storage solution for the processed sensor data.

**Data Consumers:**

- **Grafana:** Connects to the Azure SQL Database to visualize the stored data in real-time dashboards, providing meaningful insights from the sensor data.

- **Power BI:** Retrieves data from the Azure SQL Database for advanced data analysis and reporting, enabling further business intelligence capabilities.

- **Python Data Consumer:** A custom Python script (`python_data_consumer.py`) connects to the Azure SQL Database to fetch the sensor data and visualizes it using `matplotlib`, providing an additional method of exploring the data.

This project demonstrates the integration of multiple Azure services to build a serverless, cost-efficient, and scalable cloud-based data pipeline. It showcases how serverless functions, cloud-based databases, and various data consumption tools (Grafana, Power BI, and a custom Python script) work together to process, store, and visualize data in near real-time, simulating a complete IoT data pipeline.

# 3   System Overview

This project implements a cloud-based data pipeline using Azure services. The system consists of a data producer that generates sensor data (temperature, humidity, and timestamps) and sends it to an Azure Queue. An Azure Function is automatically triggered by the arrival of new data in the queue. The Azure Function processes the data and inserts it into an Azure SQL Database.

Once the data is stored in the Azure SQL Database, it is visualized using various tools:

- **Grafana:** Visualizes real-time sensor data from the SQL Database in interactive dashboards, providing an intuitive way to monitor trends.

- **Power BI:** Retrieves data from the SQL Database to create detailed reports and conduct advanced data analysis.

- **Python Data Consumer:** A custom Python script fetches data from the SQL Database and visualizes it using `matplotlib`, offering a flexible data exploration method.

The entire system leverages serverless components such as Azure Functions and Azure SQL Database for scalability, cost-efficiency, and ease of management, ensuring secure and reliable data flow from generation to consumption. This architecture simulates a robust and scalable IoT data pipeline, enabling near real-time insights into sensor data.



Figure 1: Data Pipeline Overview.

# 4   Components

## 4.1   Data Producer

The data producer script (`send_to_queue.py`) generates random temperature, humidity, and timestamp data in JSON format and sends it to an Azure Queue. The script uses environment variables to securely store sensitive connection information such as the queue name and connection string. By continuously sending data, this component simulates an

IoT-like edge device in the data pipeline, providing real-time input for the system to process.



Figure 2: Data Producer Script.

## 4.2   Azure Queue

In this project, an Azure Storage Account was used to create a Queue service, allowing for reliable message queuing between the data producer and the next stages in the pipeline. The producer script sends JSON messages containing sensor data to this queue. The messages are then picked up by an Azure Function, which processes and inserts them into the Azure SQL Database. This ensures asynchronous and scalable communication between the components in the data pipeline.



Figure 3: Azure Queue Overview.

## 4.3 Azure Function and Data Ingestion

An Azure Function (`function_app.py`) is used to ingest data from the Azure Queue into the Azure SQL Database. The function is triggered automatically whenever new data arrives in the queue, ensuring real-time processing of sensor data. The function processes the JSON message, extracts the temperature, humidity, and timestamp, and inserts this data into the `SensorData` table in the Azure SQL Database. Environment variables are used to store sensitive information like database connection strings, ensuring security and scalability without relying on virtual machines.



```python
function_app.py > queue_trigger
18    DB_SERVER = os.getenv('DB_SERVER')  # The database server URL (e.g., mydbserver1090.database.windows.net)
19    DB_USER = os.getenv('DB_USER')  # The username for the SQL Server (e.g., adminuser)
20    DB_PASSWORD = os.getenv('DB_PASSWORD')  # The password for the SQL Server (ensure it's secure)
21    DB_NAME = os.getenv('DB_NAME')  # The name of the SQL database (e.g., mydatabase1090)
22    AZURE_QUEUE_NAME = os.getenv('AZURE_QUEUE_NAME')
23
24
25    # Define the function triggered by Azure Storage Queue events
26    @app.queue_trigger(arg_name="azqueue", queue_name=AZURE_QUEUE_NAME, connection="AzureWebJobsStorage")
27    def queue_trigger(azqueue: func.QueueMessage):
28        """
29        This function is triggered when a new message is added to the Azure Storage Queue 'AZURE_QUEUE_NAME'.
30        The message contains sensor data, which is parsed and inserted into the SQL database.
31        """
32
33        # Log the received queue message
34        logging.info('Queue trigger function processed a message: %s', azqueue.get_body().decode('utf-8'))
35
36        # Parse the JSON message from the queue
37        try:
38            # Decode the queue message and parse it as JSON
39            message = json.loads(azqueue.get_body().decode('utf-8'))
40            logging.info("Parsed message: %s", message)
41
42            # Establish a connection to the SQL database using pymssql
43            with pymssql.connect(DB_SERVER, DB_USER, DB_PASSWORD, DB_NAME) as conn:
44                with conn.cursor() as cursor:
45                    # Execute the SQL INSERT query to store sensor data into the database
46                    cursor.execute(
47                        """
48                        INSERT INTO SensorData (Temperature, Humidity, Timestamp)
49                        VALUES (%s, %s, %s)
50                        """,
51                        (
52                            message["temperature"],  # Insert the temperature value from the message
53                            message["humidity"],  # Insert the humidity value from the message
54                            datetime.utcfromtimestamp(message["timestamp"]).strftime('%Y-%m-%d %H:%M:%S')  # Conv
55                        )
56                    )
```

Figure 4: Azure Function for Data Ingestion.

## 4.4 Azure SQL Database

The Azure SQL Database acts as the central storage component in the data pipeline. It provides a serverless and scalable solution for storing the processed sensor data. The serverless tier was selected to minimize costs while maintaining the ability to scale automatically based on demand, and it pauses during inactivity.

The data from the Azure Queue is processed by the Azure Function and inserted into the `SensorData` table in the SQL Database, which stores:
**Temperature:** The sensor-generated temperature data.
**Humidity:** The humidity levels from the sensors.
**Timestamp:** The exact time when the data was generated.

The database offers a secure and high-availability storage solution, ensuring that the sensor data can be easily accessed by data consumers like Grafana, Power BI, and a custom Python Data Consumer for visualization and analysis. This component is critical for managing and storing data in the pipeline.
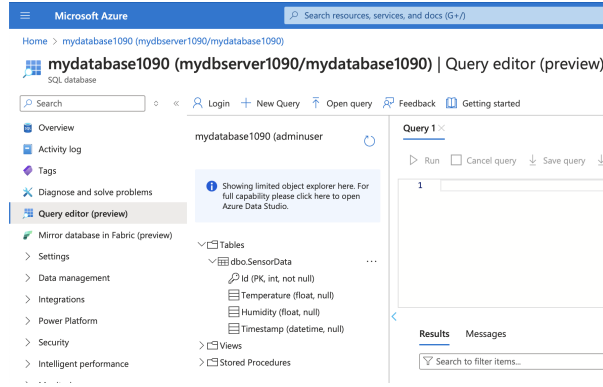
Figure 5: Azure SQL Database Overview.

## 4.5   Data Consumers

The stored data in the Azure SQL Database is visualized and analyzed using three different data consumers: **Grafana**, **Power BI**, and a **Custom Python Data Consumer**. These tools provide insight into the data, confirming that the pipeline is functioning as expected and allowing users to monitor the real-time flow of sensor data.

**Grafana:** Grafana was connected to the Azure SQL Database to create a real-time dashboard. It displays key metrics like temperature, humidity, and timestamp in an interactive interface, helping to monitor the incoming data. Grafana provides a clear, visual representation of the data pipeline's performance and ensures that the data flow is functioning correctly.



Figure 6: Grafana Dashboard Visualization.

**Power BI:** Power BI was used to visualize the same data stored in the SQL Database. Similar to Grafana, it retrieves and displays temperature, humidity, and timestamp data, providing advanced data analytics and reporting capabilities. Power BI's flexibility allows for more detailed reporting and further insights into the collected data.
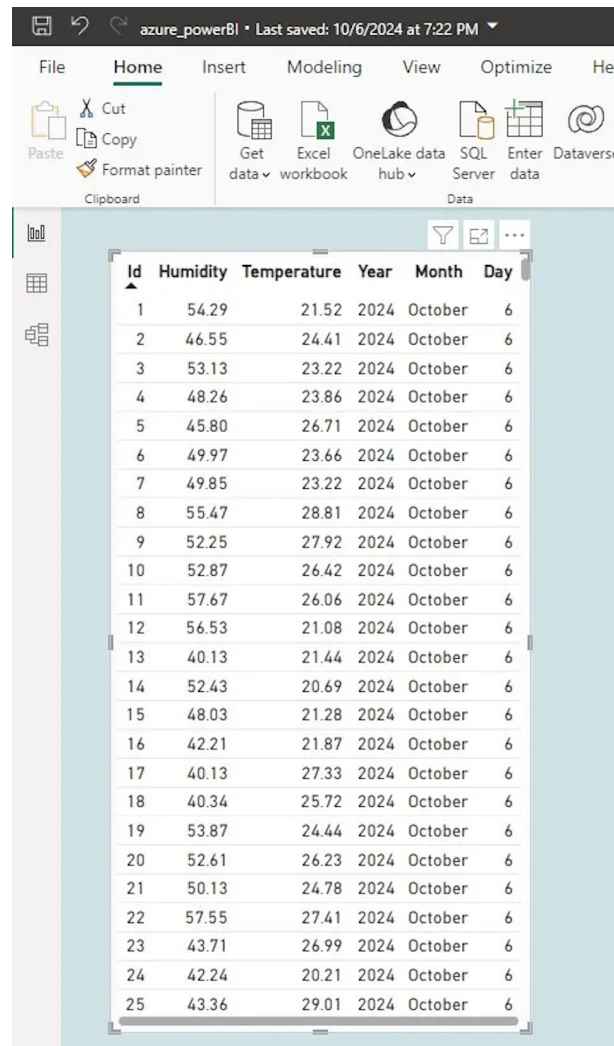
Figure 7: Power BI Dashboard Visualization.

**Python Data Consumer:** A custom Python script (`python_data_consumer.py`) was developed and executed on an Azure Virtual Machine (VM) to connect to the Azure SQL Database and retrieve sensor data. The script securely connects to the database using environment variables and fetches the latest 20 records of temperature, humidity, and timestamp data. The retrieved data is then visualized using `matplotlib`, which generates line plots showing trends of temperature and humidity over time. The generated plot is saved on the VM for further analysis or reporting.

```
import pyodbc
import matplotlib.pyplot as plt

# SQL bağlantısı ve verilerin çekilmesi
with pyodbc.connect(DB_CONNECTION_STRING) as conn:
    cursor = conn.cursor()
    cursor.execute("SELECT TOP 20 Temperature, Humidity, Timestamp FROM SensorData ORDER BY Timestamp DESC")
    rows = cursor.fetchall()

    if rows:
        # Her bir satırı yazdırma
        for row in rows:
            print(row)

        # Sıcaklık, nem ve zaman verilerini ayrı ayrı dizilere ayırma
        temperatures = [row[0] for row in rows]
        humidity = [row[1] for row in rows]
        timestamps = [row[2] for row in rows]

        # Zaman damgalarını insan tarafından okunabilir bir formata dönüştürme
        timestamps = [timestamp.strftime('%Y-%m-%d %H:%M:%S') for timestamp in timestamps]

        # Grafik oluşturma
        plt.plot(timestamps, temperatures, label='Temperature')
        plt.plot(timestamps, humidity, label='Humidity')
        plt.xlabel('Timestamp')
        plt.ylabel('Values')
        plt.xticks(rotation=45)
        plt.legend()
        plt.tight_layout()

        # Grafik dosyasına kaydetme
        plt.savefig('sensor_data_plot.png')  # Grafiği PNG dosyası olarak kaydeder
        print("Plot saved as 'sensor_data_plot.png'")
    else:
        print("No data found.")
ubuntu@my-ubuntu-vm:~$
```

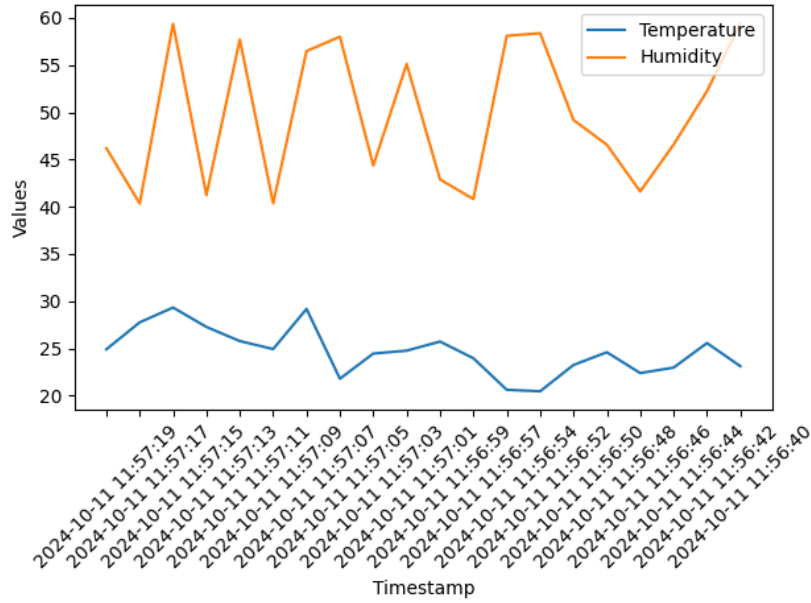Figure 8: Execution of the `python_data_consumer.py` script on the Azure VM.



Figure 9: Visualization of sensor data using `matplotlib`, showing trends in temperature and humidity.

Both Grafana, Power BI, and the Python script confirm the proper flow of data from generation through to storage and visualization, ensuring the pipeline's integrity and performance.
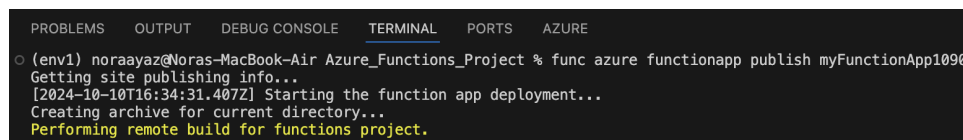
# 5   Results

The data pipeline was successfully implemented using Azure services, enabling the end-to-end flow of data from generation to storage and consumption. The Python producer script (`send_to_queue.py`) generated random sensor data and sent it to the Azure Queue.

An Azure Function (`function_app.py`) was automatically triggered by the arrival of new data in the queue, processing the messages and inserting the data into the Azure SQL Database.

Data consumers like Grafana, Power BI, and a custom Python Data Consumer were successfully configured to visualize the data stored in the database. Grafana provided real-time data insights using simple charts and graphs, while Power BI presented the data in a structured tabular format. The custom Python Data Consumer (`python_data_consumer.py`) retrieved the data from the Azure SQL Database and generated line plots using `matplotlib` to visualize temperature and humidity trends over time.

The system ran efficiently and without any major issues, with data being continuously inserted into the Azure SQL Database. The successful insertion of data was verified by querying the SQL Database and viewing it in the data consumer tools. This pipeline demonstrates the scalability and flexibility of a serverless architecture with Azure's integration, seamlessly handling data ingestion, storage, and visualization.
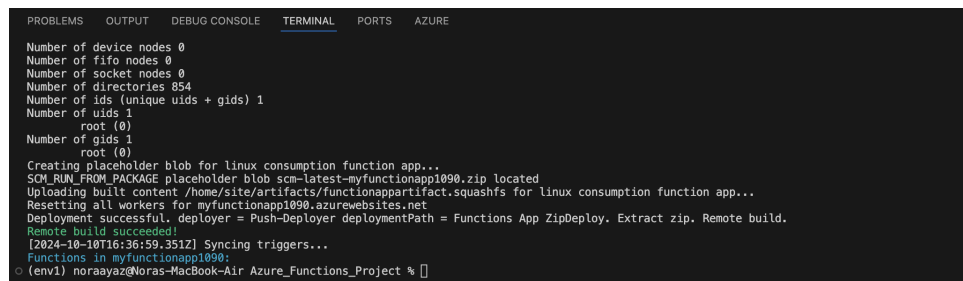


Figure 10: Deployment command initiated via Azure CLI.



Figure 11: Successful deployment of Azure Function via Azure CLI.



Figure 12: Logs from the local execution of the Azure Function, confirming successful data ingestion and processing during testing.

Figure 13: Azure Portal Log Stream displaying real-time execution logs of the deployed Azure Function.



Figure 14: Query results from the Azure SQL Database showing successfully ingested sensor data.

Screenshot showing the data successfully inserted into the Azure SQL Database by the Azure Function, confirming that the end-to-end data flow was successful.

# 6 Cost Analysis

For the cost analysis, I compared the pricing of a 10TB Azure SQL Database (serverless tier) against Amazon RDS for SQL server, focusing on monthly costs for both storage and compute resources.

**Azure SQL Database (Serverless):**
**4 vCores:**
Compute (Pay-as-you-go): $732/month$
Storage (10TB): $2,648/month$
Total: $3,385/month (Pay-as-you-go)$
With 1-year reservation: $527/month (Compute), reducing the total to $3,175/month.

**16 vCores:**
Compute (Pay-as-you-go): $2,648/month$
Storage (10TB): $2,648/month$
Total: $5,583/month (Pay-as-you-go)$
With 1-year reservation: $2,109/month (Compute), reducing the total to 4,757$/month.

**Amazon RDS for SQL server:**
**16 vCores with 64 GiB memory:**
Compute: $6,990/month$
Storage (10TB): $2,457/month$
Total: $9,448/month (Pay-as-you-go)$.

**Comparison and Conclusion:**
Azure SQL Database offers a significantly lower total cost compared to AWS RDS for similar compute and storage configurations. Using Azure's reserved pricing model further reduces costs, making Azure the more economical choice, especially for long-term use. Backup storage costs were minimal in both cases, with Azure offering a slightly cheaper option ($4.46/month for Azure vs 0.45$/month for AWS). Overall, Azure SQL Database provides a more cost-effective solution with better pricing for reserved compute instances, particularly for high storage scenarios like a 10TB database.



Figure 15: Azure SQL Database 16 vCores Estimated Cost



Figure 16: Azure SQL Database Storage Estimated Cost

Figure 17: Amazon RDS for SQL server Estimated Cost



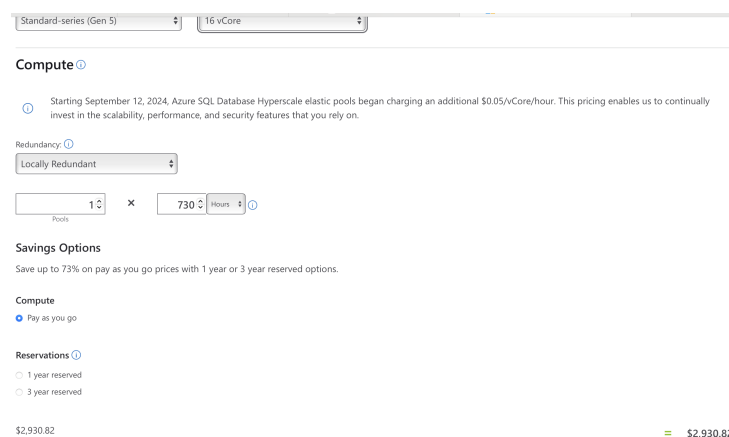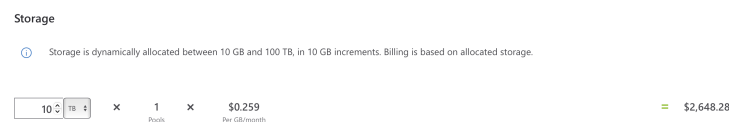Figure 18: Amazon RDS for SQL server storage Estimated Cost

# 7 Conclusion

In this project, I successfully implemented a cloud-based data pipeline using Azure services. The system efficiently handled data production, queuing, and storage, with real-time sensor data being ingested from a Python-based producer into an Azure Queue. This data was then processed by an Azure Function, which inserted it into an Azure SQL Database, enabling serverless, automated data ingestion. The stored data was visualized using Grafana and Power BI, providing meaningful insights through real-time dashboards and reports. Additionally, a custom Python data consumer was developed to visualize the data using `matplotlib`.

The project demonstrated the benefits of a serverless architecture for scalability, flexibility, and cost-effectiveness. By leveraging Azure Functions, Azure SQL Database, and data visualization tools, we ensured secure, reliable, and real-time data flow from generation to consumption. This implementation highlights how cloud-based solutions can simplify the management of data pipelines while ensuring efficient processing and visualization of data.

# 8 Challenges and Solutions

Throughout the implementation of this cloud-based data pipeline, I encountered several challenges, each contributing to valuable learning experiences. Below are the key difficulties I faced and the solutions I applied:

## 8.1 Data Generation and Queue Poison Issues

While developing the `send_to_queue.py` script to generate and send data to the Azure Queue, I initially faced a problem where the generated data was repeatedly sent to the **poison queue** instead of the target queue. This issue was quite challenging and time-consuming, as I struggled to understand the root cause. After extensive research, I came across a YouTube video that explained the problem was related to how Azure handles encoding for queue messages. The solution involved explicitly setting the message encoding and decoding policies with the following lines of code:

```
queue_client._message_encode_policy = BinaryBase64EncodePolicy()
queue_client._message_decode_policy = BinaryBase64DecodePolicy()
```

Once I implemented these policies, the data was successfully sent to the target queue.

## 8.2 Azure Function Deployment

Deploying the Azure Function was another significant challenge, as I had no prior experience with function deployment on the Azure platform. Initially, I struggled to understand how to deploy the function both locally and to the cloud. With the help of ChatGPT, I learned how to deploy the function using both the **Azure CLI** and **local deployment** methods. This guidance helped me navigate the Azure portal and correctly deploy my function, allowing it to process messages from the queue and insert them into the SQL database.

## 8.3 Database Connectivity Issues

One of the most frustrating challenges was related to transferring data from the Azure Function to the SQL database. Initially, I attempted to use the **ODBC driver** to connect to the database, but encountered compatibility issues. Since Azure Function Apps in Python are Linux-based, the ODBC driver was not supported in my environment, leading to repeated connection failures. I even attempted using Docker as an alternative, but that didn't resolve the issue either.

After more troubleshooting and assistance from ChatGPT, I was advised to switch from the ODBC driver to the **pymssql** library. After restructuring my code to use **pymssql**, I was finally able to connect to the Azure SQL Database successfully, resolving the connectivity issues.

## 8.4 Learning and Persistence

Each of these challenges required substantial research and persistence. From figuring out how to send messages to the correct queue to deploying the Azure Function and handling database connectivity, every step of this project taught me something new. While these obstacles consumed a significant amount of time, they also provided me with a deep understanding of how a data pipeline works in Azure, from data generation to ingestion, storage, and visualization. The experience of troubleshooting and learning from multiple sources, including ChatGPT, enhanced my technical knowledge and problem-solving skills.

# 9 Final Thoughts

This project provided a valuable hands-on experience in designing and implementing a serverless data pipeline using Azure cloud services. Each phase of the project, from setting up the environment to overcoming technical challenges, contributed to a deeper understanding of how cloud platforms can handle real-time data processing and visualization.

By working through the issues with queuing, database connectivity, and function deployment, I developed not only technical skills but also resilience and problem-solving abilities. The ability to research and troubleshoot effectively, coupled with leveraging tools like ChatGPT, allowed me to complete the pipeline and understand its full lifecycle from data generation to consumption.

This project reinforced the importance of serverless architecture for scalability, flexibility, and cost-efficiency, and it equipped me with practical knowledge of how to integrate multiple Azure components into a cohesive system. These skills will be invaluable in future cloud-based projects.