# Analytics dashboard using Azure Databricks

## Introduction:

In today's data-driven world, the ability to extract insights from streaming data in real-time is paramount for businesses to stay competitive and responsive. This project focuses on harnessing the power of Azure Databricks and PySpark to create a real-time analytics dashboard. By ingesting streaming data, processing it with PySparkSQL, and visualizing the results dynamically, organizations can gain actionable insights and make informed decisions instantaneously.

The real-time analytics dashboard created through this project not only enables organizations to monitor key metrics and KPIs as they evolve but also empowers stakeholders to respond proactively to changing conditions.
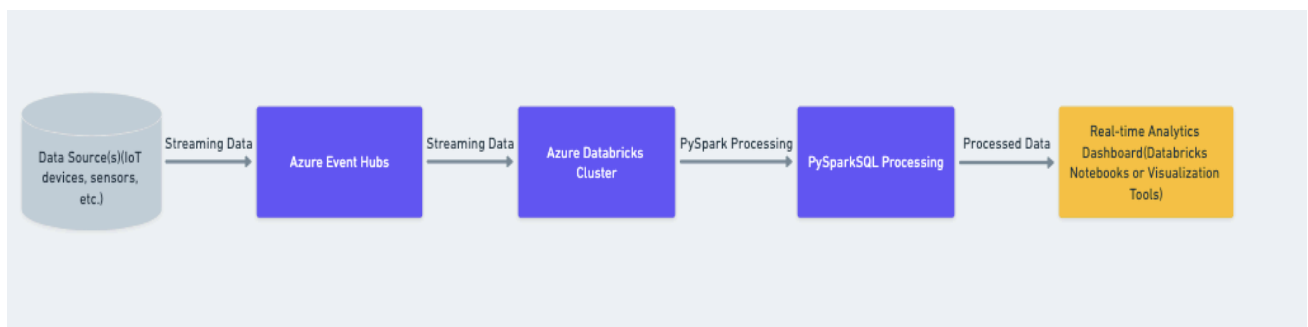
The aim of this project is to develop a real-time analytics dashboard leveraging Azure Databricks and PySpark. The objective is to ingest streaming data, process it using PySparkSQL, and visualize the results in real-time using Databricks notebooks or other visualization tools. This project aligns with the growing demand for real-time analytics solutions in various industries, providing insights into streaming data as it arrives.

## Project Overview

The primary goal of this project is to develop a real-time analytics dashboard using Azure Databricks and PySpark. The project aims to ingest streaming data, process it with PySparkSQL, and visualize the results in real-time. By achieving this objective, the project enables organizations to gain immediate insights from streaming data sources, facilitating quick decision-making and proactive responses to evolving situations.

# Architecture Diagram-



# Execution overview-

1. **Project Setup:**

   - Set up Azure Databricks environment, including creating a Databricks workspace and configuring a cluster with appropriate specifications.

   - Provision necessary Azure resources such as Azure Event Hubs for streaming data ingestion.

- Install required libraries and dependencies within the Databricks environment for PySpark development.

2. **Data Ingestion:**

- Configure Azure Event Hubs as the streaming data ingestion service.

- Set up event listeners to continuously ingest streaming data from the data source(s) into Azure Databricks.

3. **Data Processing with PySpark:**

- Define schemas for the streaming data to ensure proper parsing and transformation.

- Implement PySparkSQL queries and transformations to process the streaming data in real-time.

- Apply PySpark functions for tasks such as filtering, aggregation, and enrichment based on the analytical requirements.

4. **Real-time Analytics Dashboard Development:**

- Create Databricks notebooks or leverage compatible visualization tools for building the real-time analytics dashboard.

- Develop interactive charts, graphs, and widgets to visualize the insights derived from the streaming data.

- Implement live updates and streaming capabilities within the dashboard to reflect changes in the underlying data in real-time.

- Enhance user experience by providing interactive controls and filters for dynamic exploration of the data.

5. **Performance Optimization:**

- Optimize PySpark code and queries to improve processing performance and reduce latency in real-time analytics.

- Utilize caching, partitioning, and other optimization techniques to enhance the efficiency of data processing operations.

- Monitor cluster performance and resource utilization to ensure scalability and reliability under varying workload conditions.

6. **Testing and Validation:**

- Conduct thorough testing of the real-time analytics dashboard to ensure its functionality, performance, and reliability.

- Validate the accuracy of the insights generated by comparing them with expected results and business requirements.

- Address any issues or bugs identified during testing and make necessary refinements to the implementation.

7. **Deployment and Integration:**

- Deploy the real-time analytics dashboard to the production environment, making it accessible to stakeholders and end-users.

- Integrate the dashboard with other systems or data sources as needed to enrich the analysis and provide comprehensive insights.

- Ensure seamless integration with existing workflows and processes within the organization.

8. **Documentation and Knowledge Sharing:**

- Document the project architecture, implementation details, and deployment instructions for future reference.

- Conduct knowledge sharing sessions to disseminate insights gained from the project and promote best practices in real-time analytics and data visualization.

Flow diagram of the project

# Key components/Technologies Used

1. **Azure Databricks:**

   - Provides a cloud-based platform for big data processing and analytics.

   - Offers scalable clusters with Apache Spark capabilities for distributed data processing.

   - Enables seamless integration with other Azure services for data ingestion, storage, and visualization.

2. **PySpark:**

   - Python API for Apache Spark, facilitating data processing and analytics tasks.

   - Allows for the development of distributed data processing applications using the Spark framework.

   - Provides rich libraries and functions for data manipulation, querying, and analysis.

3. **Azure Event Hubs:**

   - Scalable event ingestion service for streaming data.

   - Facilitates the ingestion of high-volume, real-time data streams from various sources.

   - Integrates with Azure Databricks for seamless data ingestion into the analytics pipeline.

4. **Data Sources:**

   - IoT devices, sensors, application logs, or other streaming data sources.

   - Provide real-time data streams containing valuable insights for analysis.

5. **Databricks Notebooks:**

   - Interactive development environment for data exploration, analysis, and visualization.

   - Supports various programming languages including Python, SQL, and Scala.

   - Enables the creation of interactive dashboards and visualizations for real-time analytics.

6. **Visualization Tools:**

- Matplotlib, Plotly, Bokeh, or other Python libraries for creating interactive visualizations.

- Used to build real-time dashboards with charts, graphs, and widgets for data visualization.

7. **PySparkSQL:**

- Component of PySpark for executing SQL queries against Spark dataframes.

- Allows for querying and processing structured data in real-time.

- Enables the implementation of complex data transformations and aggregations for analytics.

# How it works-

1. ### Setting Up Azure Databricks Environment:

   - Sign in to the Azure portal and create an Azure Databricks workspace. Configure workspace settings, including pricing tier, region, and workspace name.

## 2. Developing PySpark Notebooks:

- Create a new PySpark notebook within the Databricks workspace. Begin writing PySpark code to perform ETL operations, data transformations, and other data processing tasks.

## 3. Create Cluster and Connecting to notebook

- The cluster is created with 4 working nodes and autoscaling is enabled which automatically adjust cluster size to accommodate changes in workload demand, allowing for seamless scalability without manual intervention.

## 4. Importing Necessary libraries and Loading the data



## 5. View a dataframe-

# 6. Run SQL queries



# 7. Visualize the dataframe

## 8. Create the dashboard—

# 9.Dashboard with widget enabled

# Final Output-

# Conclusion

In conclusion, the real-time analytics dashboard project utilizing Azure Databricks and PySpark offers a transformative solution for organizations seeking to harness the power of streaming data for actionable insights and informed decision-making. By ingesting, processing, and visualizing streaming data in real-time, the project enables stakeholders to monitor key metrics, detect trends, and respond proactively to changing conditions.

The project demonstrated the following key outcomes and benefits:

1. **Real-time Insights:** By processing streaming data in real-time, the dashboard provides immediate insights into operational performance, customer behaviour, and emerging trends.

2. **Agility and Scalability:** Leveraging Azure Databricks' cloud-based infrastructure, the solution offers agility and scalability to adapt to evolving data requirements and handle growing data volumes.

3. **Decision Support:** The interactive dashboards and visualizations empower stakeholders with the information needed to make data-driven decisions quickly and confidently.

4. **Efficiency and Performance:** Through optimization techniques and distributed computing capabilities of PySpark, the project ensures efficient data processing and low-latency analytics.

6. **Collaboration and Knowledge Sharing:** The project fosters collaboration among data engineers, data scientists, and analysts through shared notebooks, documentation, and knowledge-sharing sessions.