# Data Engineering in the Cloud

## Real-time Data Processing with Azure Databricks

Xuemao Zhang

East Stroudsburg University

January 18, 2025

# Outline

- Azure Services Required
- Lab

# Azure Services Required

- Data Sources: Streaming data from IoT devices or social media feeds. W use simulated data in Event Hubs.
- Ingestion: Azure Event Hubs for capturing real-time data.
- Processing: Azure Databricks for stream processing using Structured Streaming.
- Storage: Processed data stored Azure Data Lake (Delta Format).
- The project is copied from Real Time Streaming with Azure Databricks and Event Hubs https://youtu.be/pwWIegHgNRw?si=SP1S9r0Z4AZPzGYA

# Azure Services Required

Compared to Azure Stream Analytics, Azure Databricks

- Highly scalable, can handle large volumes of data with Spark's distributed computing power.
- Interactive development environment with notebooks for collaborative data analysis and visualization.
  - Flexibility: Supports complex transformations and machine learning tasks using Python, Scala, SQL, and R.
- Advanced Analytics: Suitable for advanced analytics, machine learning, and data science use cases.

# Lab

**Step 1**. Create a datalake (storage account with hierarchical namespace enabled)
- To lower the cost, you may choose `Redundancy` as LRS - Create a container in
the storage account

## Lab

**Step 2:** Search for Event Hubs to create an event hub namespace to which the data will be streamed: - Event Hub acts as a message broker for the streaming data.

1. Enter the details shown in the image below and click on the "Review + create" button.

# Lab

2. Validate the configuration and click on the "Create" button.

# Lab

③ Once the deployment process is completed, click on the "Go to resource"
button.

# Lab

4. The Event Hubs dashboard will be as shown in the image below.

# Lab

⑤ click on the "+ Event Hub" button to create an event hub

- Add a name and select retention period as shown in the image below and click on the "Create" button. (Note down the event hub name for future reference)

# Lab

⑥ Click `Create`

# Lab

7. Go to the `Settings` of the Event Hubs Namespace and go to Shared access policies.

- Update the policy if necessay

# Lab

**Step 2**: Create an Azure Databricks workspace

# Lab

⑧ Click `Create`

Home > Azure Databricks >

## Create an Azure Databricks workspace ⋯

✓ Validation Succeeded

**Basics**

| | |
|---|---|
| Workspace name | xzhang2db |
| Subscription | AzureSubscription |
| Resource group | xzhang2 |
| Region | East US 2 |
| Pricing Tier | standard |
| Managed Resource Group name | |

**Networking**

| | |
|---|---|
| Deploy Azure Databricks workspace with Secure Cluster Connectivity (No Public IP) | No |
| Deploy Azure Databricks workspace in your own Virtual Network (VNet) | No |

**Encryption**

| | |
|---|---|
| Enable Infrastructure Encryption | No |
| Enable CMK for Managed Disks | No |
| Enable CMK for Managed Services | No |

**Security & compliance**

# Lab

9. Go to Azure Databricks resource and click on "Launch Workspace".

# Lab

**10** In the workspace, click on "New Cluster" to create a cluster.

# Lab

⑪ Settings of the cluster are shown below.

# Lab

12. Click on "Libraries" and then Click on "Install New"

# Lab

- 13 Go to Maven and click on "Search Packages"
- 14 Select `Maven Central` and search for "azure-eventhubs-spark" and Select "azure-eventhubs-spark_2.12"

# Lab

15. After selecting the package, click on "Install"

# Lab

16. Go to Azure Databricks workspace, and click on "Compute" and click on the cluster

# Lab

⑰ Click on "+ New" -> Notebook to create a new Notebook

# Lab

**18** Type the code in a cell and run the cell

```
#Importing the libraries
from pyspark.sql.functions import *
from pyspark.sql.types import *
```

# Lab

**19** Go to the Event Hubs Namespace and click on the event hub

# Lab

**㉑** Click `Shared access policies`

# Lab

㉑ Create the policy

# Lab

22. Click on the policy and Copy the connection string

# Lab

**23** Go back to the notebook in the databricks

```
# create a schema
spark.sql("create schema streaming;")
```

```
#Set up Azure Event hubs connection string
#Replace with your Event Hub namespace, name, and key
connectionString = "Endpoint=sb://xzhang2hubs.servicebus.windows.net/;
SharedAccessKeyName=databricks;
SharedAccessKey=yWtWfhtb9OhdAYGi9Wtrbzmw8WWYI3pNV+AEhKgRKdo=;
EntityPath=esuhub1"
eventHubName = "esuhub1"

ehConf = {
  'eventhubs.connectionString' : sc._jvm.org.apache.spark.eventhubs.EventHubsUtils.e
  'eventhubs.eventHubName': eventHubName
}
```

# Lab

- We can display the contents of a DBFS mount point or directory in Azure Databricks by running

```python
# Using shell command to list files
files = dbutils.fs.ls("/mnt")

# Display the files in a table format
display(files)
```

# Lab

```
         ✓ 2 minutes ago (2s)                    2
  # create a schema
  spark.sql("create schema streaming;")
```
```
DataFrame[]
```

```
  ▶ ✓   ✓ 1 minute ago (<1s)                    3                              Python  ✦  ⌗  ⋮

  #Set·up·Azure·Event·hubs·connection·string
  #Replace·with·your·Event·Hub·namespace,·name,·and·key
  connectionString·=·"Endpoint=sb://xzhang2hubs.servicebus.windows.net/;SharedAccessKeyName=RootManageSharedAccessKey;
  SharedAccessKey=jnBPua+J2mGHJuMmJyFyH01U8lZKjFm2o+AEhLTW7ws="
  eventHubName·=·"esuhub1"

  ehConf·=·{
  ··'eventhubs.connectionString'·:·sc._jvm.org.apache.spark.eventhubs.EventHubsUtils.encrypt(connectionString),
  ··'eventhubs.eventHubName':·eventHubName
  }
```

# Lab

24. Run the following cell

- Send data from the event hub while the stream is running

```
{
    "temperature": 20,
    "humidity": 60,
    "windSpeed": 10,
    "windDirection": "NW",
    "precipitation": 0,
    "conditions": "Partly Cloudy"
}
```

```
# Reading stream: Load data from Azure Event Hub into DataFrame 'df' using the previ
df = spark.readStream \
    .format("eventhubs") \
    .options(**ehConf) \
    .load() \

# Displaying stream: Show the incoming streaming data for visualization and debuggin
df.display()
```

# Lab

# Lab

25 Writing stream

```python
# Writing stream: Persist the streaming data to a Delta table 'stre
df.writeStream\
    .option("checkpointLocation", "/mnt/streaming/weather")\
    .outputMode("append")\
    .format("delta")\
    .toTable("streaming.weather")
```

# Lab

In Spark Structured Streaming, the outputMode specifies how the results of a streaming query should be written to the output sink. There are three output modes: append, complete, and update. Each mode has different behaviors and use cases:

- Append Mode: Only the new rows that are appended to the result table since the last trigger are written to the sink

- Complete Mode: The entire result table is written to the sink every time there is a trigger. This means all rows of the result are written out, not just the new rows.

- Update Mode: Only the rows that have changed since the last trigger are written to the sink. This includes both newly added rows and updated row

# Lab

26 Check the Catalog, you should be able to see the table weather

# Lab

**27** Defining the schema for the JSON object.

```python
# Defining the schema for the JSON object
json_schema = StructType([
    StructField("temperature", IntegerType()),
    StructField("humidity", IntegerType()),
    StructField("windSpeed", IntegerType()),
    StructField("windDirection", StringType()),
    StructField("precipitation", IntegerType()),
    StructField("conditions", StringType())
])
```

# Lab

**28** Reading and Transforming

```python
# Reading and Transforming: Load streaming data from the 'streaming.weather' Delta table, cast 'b
df = spark.readStream\
    .format("delta")\
    .table("streaming.weather")\
    .withColumn("body", col("body").cast("string"))\
    .withColumn("body",from_json(col("body"), json_schema))\
    .select("body.temperature", "body.humidity", "body.windSpeed",
            "body.windDirection", "body.precipitation",
            "body.conditions", col("enqueuedTime").alias('timestamp'))

# Displaying stream: Visualize the transformed data in the DataFrame for verification and analys
df.display()

# Writing stream: Save the transformed data to the 'streaming.weather2' Delta table in 'append' r
df.writeStream\
    .option("checkpointLocation", "/mnt/streaming/weather2")\
    .outputMode("append")\
    .format("delta")\
    .toTable("streaming.weather2")
```

# Lab

# Lab

29. Check the `Catalogy` again

# Lab

**30** Reading, aggregating and writing the stream

```
# Aggregating Stream: Read from 'streaming.silver.weather', apply watermarking and windowing, an
df = spark.readStream\
    .format("delta")\
    .table("streaming.weather2")\
    .withWatermark("timestamp", "1 minutes") \
    .groupBy(window("timestamp", "1 minutes")) \
    .agg(avg("temperature").alias('temperature'), avg("humidity").alias('humidity'),
     avg("windSpeed").alias('windSpeed'), avg("precipitation").alias('precipitation'))\
    .select('window.start', 'window.end', 'temperature', 'humidity', 'windSpeed', 'precipitation

# Displaying Aggregated Stream: Visualize aggregated data for insights into weather trends
df.display()

# Writing Aggregated Stream: Store the aggregated data in 'streaming.gold.weather_aggregated' wi
df.writeStream\
    .option("checkpointLocation", "/mnt/streaming/weather_summary")\
    .outputMode("append")\
    .format("delta")\
    .toTable("streaming.weather_summary")
```

# Lab

# Lab

**31** Check the catalog again

# Lab

**32** Summary of the three steps

```
json_schema = StructType([
    StructField("temperature", IntegerType()),
    StructField("humidity", IntegerType()),
    StructField("windSpeed", IntegerType()),
    StructField("windDirection", StringType()),
    StructField("precipitation", IntegerType()),
    StructField("conditions", StringType())
])

# Reading stream: Load data from Azure Event Hub into DataFrame 'df' using the previously config
df1 = spark.readStream \
    .format("eventhubs") \
    .options(**ehConf) \
    .load() \

# Writing stream: Persist the streaming data to a Delta table 'streaming.weather' in 'append' mo
df1.writeStream\
    .option("checkpointLocation", "/mnt/streaming/weather")\
    .outputMode("append")\
    .format("delta")\
    .toTable("streaming.weather")
```

# Lab

```python
# Reading and Transforming: Load streaming data from the 'streaming.weather' Delta table, cast 'b
df2 = spark.readStream\
    .format("delta")\
    .table("streaming.weather")\
    .withColumn("body", col("body").cast("string"))\
    .withColumn("body",from_json(col("body"), json_schema))\
    .select("body.temperature", "body.humidity", "body.windSpeed", "body.windDirection",
            "body.precipitation", "body.conditions", col("enqueuedTime").alias('timestamp'))

# Displaying stream: Visualize the transformed data in the DataFrame for verification and analys
df2.display()

# Writing stream: Save the transformed data to the 'streaming.weather2' Delta table in 'append' m
df2.writeStream\
    .option("checkpointLocation", "/mnt/streaming/weather2")\
    .outputMode("append")\
    .format("delta")\
    .toTable("streaming.weather2")
```

# Lab

```python
# Aggregating Stream: Read from 'streaming.silver.weather', apply watermarking and windowing, and
df3 = spark.readStream\
    .format("delta")\
    .table("streaming.weather2")\
    .withWatermark("timestamp", "2 minutes") \
    .groupBy(window("timestamp", "2 minutes")) \
    .agg(avg("temperature").alias('temperature'), avg("humidity").alias('humidity'),
        avg("windSpeed").alias('windSpeed'), avg("precipitation").alias('precipitation'))\
    .select('window.start', 'window.end', 'temperature', 'humidity', 'windSpeed', 'precipitation'

# Displaying Aggregated Stream: Visualize aggregated data for insights into weather trends
df3.display()

# Writing Aggregated Stream: Store the aggregated data in 'streaming.gold.weather_aggregated' wi
df3.writeStream\
    .option("checkpointLocation", "/mnt/streaming/weather_summary")\
    .outputMode("append")\
    .format("delta")\
    .toTable("streaming.weather_summary")
```

# Lab

# License