

# Real-Time Big Data Processing

## Lab 4 – Aggregating Data Streams

### Overview

In this lab, you will extend Azure Stream Analytics solution that you created in lab 3 by aggregating the streamed data over a temporal window.

### What You'll Need

To complete the labs, you will need the following:

- A web browser
- A Microsoft account
- A Microsoft Azure subscription
- A Windows, Linux, or Mac OS X computer
- The lab files for this course
- The Azure resources created in the previous labs

**Important:** If you have not completed [Lab 1](#), [Lab 2](#), and [Lab 3](#), or you have deleted the event hub, storage account, IoT hub, and stream analytics jobs you created, complete the previous labs now.

### Exercise 1: Aggregating Events Over a Temporal Window

In most cases, data analysis involves aggregating individual data observations. When dealing with an unbounded stream of events, this aggregation is usually performed over temporal windows. For example, you find the average value of a field or the number of events within a specific time-period.

In this exercise, you will modify the query for the second Stream Analytics job to log alerts whenever a device experiences two or more readings greater than 0.5 over a 10 second period.

#### Modify a Stream Analytics Job

The second Stream Analytics job in your solution currently creates an alert for every reading greater than 0.5. This results in a high volume of alerts. You need to improve the solution so that an alert is only generated if a device experiences two or more high readings within a 10 second window.

1. In the Azure portal, browse to the second Stream Analytics job you created, which reads device data from an event hub input, and routes the processed results to a blob storage output.
2. In the blade for your Stream Analytics job, in the **Job Topology** section, click the **Query** tile.
3. Verify that the existing query looks like this:

```
SELECT  
*
```

```

INTO
    [DeviceAlerts]
FROM
    [HighReadings]

```

4. Modify the query as shown below, adding fields to show the start and end of each 10-second window in which a device experienced more than one high reading:

```

SELECT DateAdd(second,-10,System.Timestamp) AS WinStart,
       System.Timestamp AS WinEnd,
       DeviceName,
       AVG(CAST(reading as float)) AS AvgReading,
       COUNT(*) AS Alerts
INTO
    [DeviceAlerts]
FROM
    [HighReadings] TIMESTAMP BY EventProcessedUtcTime
GROUP BY DeviceName, SlidingWindow(second, 10)
HAVING COUNT(*) > 1

```

5. Save the query, and then close the query pane.

## Start the Jobs

Now you're ready to start both jobs and test the streaming topology.

1. Start both analytics jobs and wait for them to start – this can take a minute or so.
2. When the jobs have started, in the Node.JS console, in the **iotdevice** folder, enter the following command to run device simulation script and start submitting messages to the IoT hub:

```
node iotdevice.js
```

3. While the script is running, start Azure Storage Explorer, and if necessary, sign into your azure subscription using your Microsoft account.
4. Expand your storage account, and then expand **Blob Containers**.
5. Double-click the **device-readings** container, and then browse through the **alerts** folder, and the year, month, and date folder to view the most recent blob that has been generated by your job.
6. Download the blob to open it in a text editor or spreadsheet application, and verify that it contains average readings and counts for high readings in 10-second windows.
7. Close the downloaded file, and in the Node.JS console, press CTRL+C to stop the script.

## Stop the Jobs

When you want to stop processing events, you can stop the jobs.

1. In the Azure portal, stop both stream analytics jobs.