# BIG DATA MANAGEMENT SYSTEMS: PROJECT #4 – AZURE STREAM ANALYTICS

**July 5, 2019**

Zoe Kotti and Chryssa Nampouri

*Department of Management Science and Technology*

*Athens University of Economics and Business*

Athens, Greece

{t8150062, t8150096}@aueb.gr

Supervisor: Prof. Damianos Chatziantoniou

# Contents

# 1 PROJECT DESCRIPTION

In the context of this assignment a set of Reference Data is given that consists of three JSON files:

1. **Customer.json:** Personal data about customers who have made transactions. Each customer example is described by the following attributes:

   - *card_ number* (integer)
   - *first_ name* (string)
   - *last_ name* (string)
   - *age* (integer)
   - *gender* (string)
   - *area_ code* (integer)

2. **Atm.json:** Describes ATMs as:

   - *atm_ code* (integer)
   - *area_ code* (integer)

3. **Area.json:** Describes areas as:

   - *area_ code* (integer)
   - *area_ country* (string)
   - *area_ city* (string)

The purpose of this project was to process a data stream of ATM transactions and answer stream queries. The schema of the stream is the following: (ATMCode, CardNumber, Type, Amount).

# 2  Azure Configuration

In order to execute the above process on Azure Stream Analytics Platform the following steps are required:

1. Create an Azure account.

2. Setup an Event Hub.

3. Generate a Security Access Signature [1].

4. Edit Generator.html and update the CONFIG variables with your security access signature.

5. Feed the Event Hub with streaming data by using the Generator.html.

6. Setup a Storage account.

7. Upload the Reference Data files to your storage account.

8. Setup a Stream Analytics Job.

9. Use the Event Hub and the Reference Data Files as Input.

10. Create a Blob Storage Output.

11. Run queries.

---

[1]https://github.com/sandrinodimattia/RedDog/releases
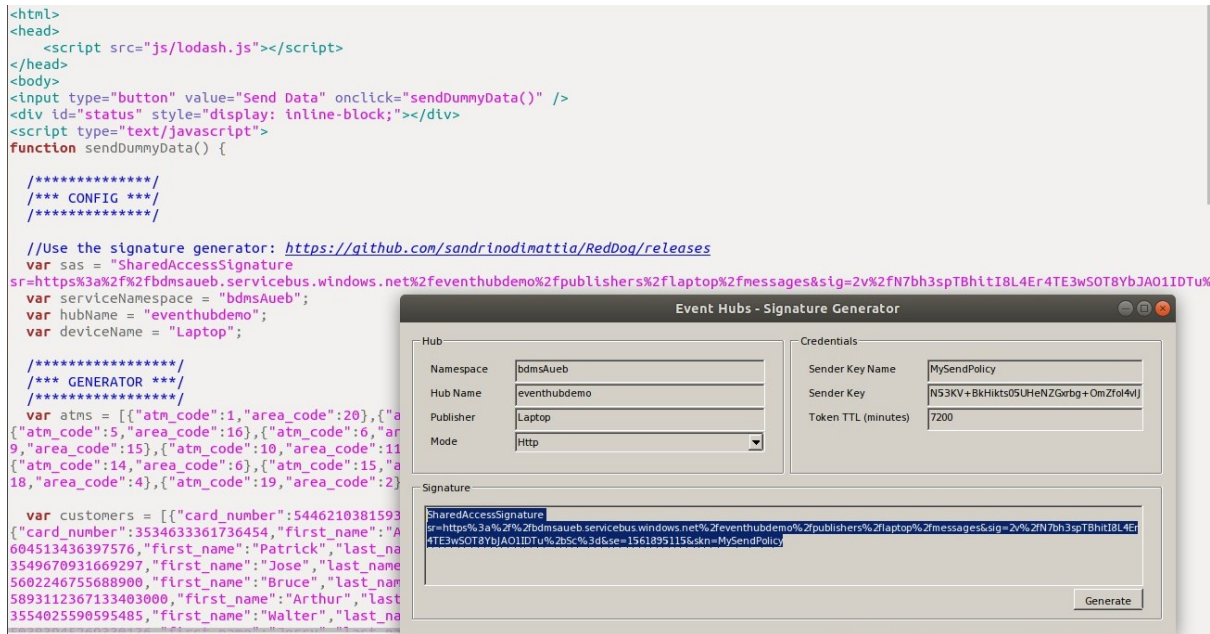
# 3 Azure Configuration: Screenshots



**Figure 1:** Security Access Signature (*Step 3,4*).



**Figure 2:** Random Data Generator (*Step 5*).


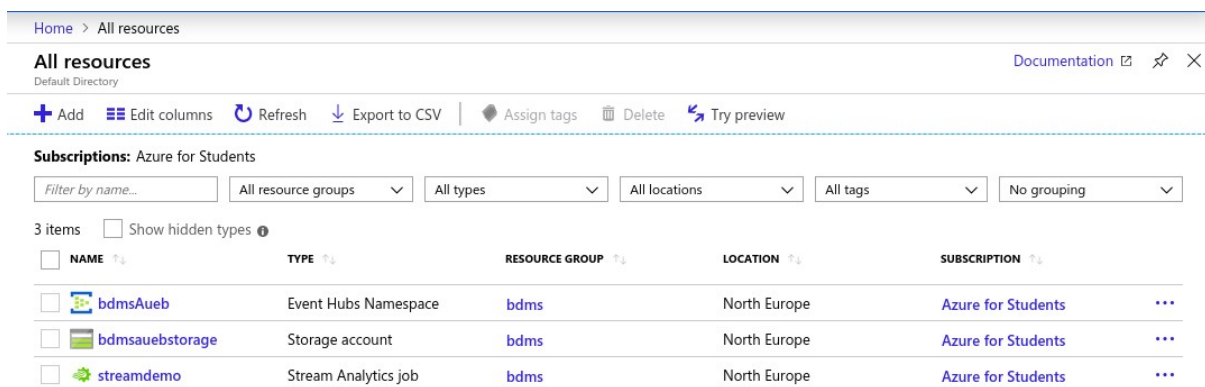
**Figure 3:** All Resources (*Step 2, 6, 8*).

**Figure 4:** Storage Account (*Step 7).*



**Figure 5:** Job Inputs (*Step 9).*



**Figure 6:** Job Output (*Step 10).*

# 4 AZURE STREAM ANALYTICS: QUERIES

After specifying the input source of the streaming data and the output sink for the results, the following queries were expressed in a *SQL-like query language (T-SQL)*, in order to be executed on the Azure Platform:

1. Show the total 'Amount' of 'Type = 0' transactions at 'ATM Code = 21' of the last 10 minutes. Repeat as new events keep flowing in (use a sliding window).

2. Show the total 'Amount' of 'Type = 1' transactions at 'ATM Code = 21' of the last hour. Repeat once every hour (use a tumbling window).

3. Show the total 'Amount' of 'Type = 1' transactions at 'ATM Code = 21' of the last hour. Repeat once every 30 minutes (use a hopping window).

4. Show the total 'Amount' of 'Type = 1' transactions per 'ATM Code' of the last one hour (use a sliding window).

5. Show the total 'Amount' of 'Type = 1' transactions per 'Area Code' of the last hour. Repeat once every hour (use a tumbling window).

6. Show the total 'Amount' per ATM's 'City' and Customer's 'Gender' of the last hour. Repeat once every hour (use a tumbling window).

7. Alert (SELECT '1') if a Customer has performed two transactions of 'Type = 1' in a window of an hour (use a sliding window).

8. Alert (SELECT '1') if the 'Area Code' of the ATM of the transaction is not the same as the 'Area Code' of the 'Card Number' (Customer's Area Code) - (use a sliding window)

For each of the above queries we ran a stream analytics job and we present some indicative results in the report below. In some cases, we reduced the time limit of the windows, in order to precipitate the progress of the data transformation.

**Query 1:**

```
1  /*
2      Q1: Show the total 'Amount' of 'Type = 0' transactions at
3      'ATM Code = 21' of the last 10 minutes. Repeat as new events keep
4      flowing in (use a sliding window).
5  */
6
7  SELECT
8      sum(CAST([input].[Amount] AS BIGINT)) as Total_Amount,
9      System.Timestamp() AS Event_Time
10 INTO [output]
11 FROM [input]
12 WHERE [input].[Type] = 0 and [input].[ATMCode] = 21
13 GROUP BY SlidingWindow(minute, 10)
```

**Output_1 (SlidingWindow(second, 10)):**

```
1  [
2    {
3       "Total_Amount":21,
4       "Event_Time":"2019-07-03T00:18:29.2580000Z"
5    },
6    {
7       "Total_Amount":38,
8       "Event_Time":"2019-07-03T00:18:39.2430000Z"
9    },
10   {
11      "Total_Amount":17,
12      "Event_Time":"2019-07-03T00:18:39.2580000Z"
13   },
14   {
15      "Total_Amount":105,
16      "Event_Time":"2019-07-03T00:18:40.2430000Z"
17   },
18   {
19      "Total_Amount":88,
20      "Event_Time":"2019-07-03T00:18:49.2430000Z"
21   },
22   {
23      "Total_Amount":119,
24      "Event_Time":"2019-07-03T00:18:49.2590000Z"
25   }
26 ]
```

**Query 2:**

```
1  /*
2      Q2: Show the total 'Amount' of 'Type = 1' transactions at
3      'ATM Code = 21' of the last hour. Repeat once every hour
4      (use a tumbling window).
5  */
6
7  SELECT
8      sum(CAST([input].[Amount] AS BIGINT)) as Total_Amount,
9      System.Timestamp() AS Time
10 INTO [output]
11 FROM [input]
12 WHERE [input].[Type] = 1 and [input].[ATMCode] = 21
13 GROUP BY TumblingWindow(hour, 1)
```

**Output_2 (TumblingWindow(minute, 1)):**

```
1  [
2    {
3      "Total_Amount": 335,
4      "Event_Time": "2019-06-29T17:36:00.0000000Z"
5    },
6    {
7      "Total_Amount": 603,
8      "Event_Time": "2019-06-29T17:37:00.0000000Z"
9    },
10   {
11     "Total_Amount": 465,
12     "Event_Time": "2019-06-29T17:38:00.0000000Z"
13   }
14 ]
```

**Query 3:**

```
1  /*
2      Q3: Show the total 'Amount' of 'Type = 1' transactions at
3      'ATM Code = 21' of the last hour. Repeat once every 30 minutes
4      (use a hopping window).
5  */
6
7  SELECT
8      sum(CAST([input].[Amount] AS BIGINT)) as Total_Amount,
9      System.Timestamp() AS Time
10 INTO [output]
11 FROM [input]
12 WHERE [input].[Type] = 1 and [input].[ATMCode] = 21
13 GROUP BY HoppingWindow(minute, 60, 30)
```

**Output_3 (HoppingWindow(second, 60, 30)):**

```
1  [
2    {
3      "Total_Amount": 406,
4      "Event_Time": "2019−06−29T17:48:00.0000000Z"
5    },
6    {
7      "Total_Amount": 651,
8      "Event_Time": "2019−06−29T17:48:30.0000000Z"
9    },
10   {
11     "Total_Amount": 728,
12     "Event_Time": "2019−06−29T17:49:00.0000000Z"
13   }
14 ]
```

**Query 4:**

```
/*
    Q4: Show the total 'Amount' of 'Type = 1' transactions per 'ATM Code'
    of the last one hour (use a sliding window).
*/

SELECT
    [input].[ATMCode],
    sum(CAST([input].[Amount] AS BIGINT)) as Total_Amount,
    System.Timestamp() AS Time
INTO [output]
FROM [input]
WHERE Type = 1
GROUP BY [input].[ATMCode],
         SlidingWindow(hour, 1)
```

**Output_4 (SlidingWindow(minute, 1)):**

```
[
  {
    "ATMCode":15
    "Total_Amount": 2519,
    "Event_Time": "2019-06-29T17:56:55.2870000Z"
  },
  {
    "ATMCode":18
    "Total_Amount": 1624,
    "Event_Time": "2019-06-29T17:56:55.2920000Z"
  },
  {
    "ATMCode":18
    "Total_Amount": 1660,
    "Event_Time": "2019-06-29T17:56:55.3020000Z"
  },
  {
    "ATMCode":19
    "Total_Amount": 2031,
    "Event_Time": "2019-06-29T17:56:55.3020000Z"
  }
]
```

**Query 5:**

```
/*
    Q5: Show the total 'Amount' of 'Type = 1' transactions per 'Area Code'
    of the last hour. Repeat once every hour (use a tumbling window).
*/

SELECT
    [inputAtm].[area_code],
    sum(CAST([input].[Amount] AS BIGINT)) as Total_Amount,
    System.Timestamp() AS Time
INTO [output]
FROM [input]
INNER JOIN [inputAtm]
    ON [input].[ATMCode] = [inputAtm].[atm_code]
WHERE [input].[Type] = 1
GROUP BY [inputAtm].[area_code],
         TumblingWindow(hour, 1)
```

**Output_5 (TumblingWindow(minute, 1)):**

```
[
  {
      "area_code":10,
      "Total_Amount":670,
      "Time":"2019-06-30T09:40:00.0000000Z"
  },
  {
      "area_code":14,
      "Total_Amount":43,
      "Time":"2019-06-30T09:40:00.0000000Z"
  },
  {
      "area_code":12,
      "Total_Amount":50,
      "Time":"2019-06-30T09:40:00.0000000Z"
  },
  {
      "area_code":2,
      "Total_Amount":1813,
      "Time":"2019-06-30T09:40:00.0000000Z"
  }
]
```

**Query 6:**

```
1  /*
2      Q6: Show the total 'Amount' per ATM's 'City' and Customer's 'Gender'
3      of the last hour. Repeat once every hour (use a tumbling window).
4  */
5
6  SELECT
7      [inputArea].[area_city],
8      [inputCustomers].[gender],
9      sum(CAST([input].[Amount] AS BIGINT)) as Total_Amount,
10     System.Timestamp() AS Time
11  INTO [output]
12  FROM [input]
13  INNER JOIN [inputAtm]
14      ON [input].[ATMCode] = [inputAtm].[atm_code]
15  INNER JOIN [inputArea]
16      ON [inputAtm].[area_code] = [inputArea].[area_code]
17  INNER JOIN [inputCustomers]
18      ON [input].[CardNumber] = [inputCustomers].[card_number]
19  GROUP BY [inputArea].[area_city],
20           [inputCustomers].[gender],
21           TumblingWindow(hour, 1)
```

**Output_6 (TumblingWindow(minute, 1)):**

```
1  [
2    {
3      "area_city":"Springfield",
4      "gender":"Male",
5      "Total_Amount":1066,
6      "Time":"2019-06-30T20:12:00.0000000Z"
7    },
8    {
9      "area_city":"Baltimore",
10     "gender":"Male",
11     "Total_Amount":384,
12     "Time":"2019-06-30T20:12:00.0000000Z"
13   }
14  ]
```

**Query 7:**

```
1  /*
2      Q7: Alert (SELECT '1') if a Customer has performed two transactions
3      of 'Type = 1' in a window of an hour (use a sliding window).
4  */
5
6  SELECT
7      [inputCustomers].[first_name],
8      [inputCustomers].[last_name],
9      [input].[CardNumber] AS Card_Number,
10     COUNT (*) AS Transactions,
11     System.Timestamp AS Time
12 INTO [output]
13 FROM [input]
14 INNER JOIN [inputCustomers]
15     ON [inputCustomers].[card_number] = [input].[CardNumber]
16 WHERE [input].[Type] = 1
17 GROUP BY [inputCustomers].[first_name],
18          [inputCustomers].[last_name],
19          [input].[CardNumber],
20          SlidingWindow(hour, 1)
21 HAVING Transactions = 2
```

**Output_7 (SlidingWindow(hour, 1)):**

```
1  [
2    {
3        "first_name":"Brenda",
4        "last_name":"Carroll",
5        "Card_Number":560222217915598000,
6        "Transactions":2,
7        "Time":"2019-06-30T12:49:59.7440000Z"
8    },
9    {
10       "first_name":"Bruce",
11       "last_name":"Morrison",
12       "Card_Number":5602246755688900,
13       "Transactions":2,
14       "Time":"2019-06-30T12:50:02.6970000Z"
15   }
16 ]
```

**Query 8:**

```
/*
    Q8: Alert (SELECT '1') if the 'Area Code' of the ATM of the transaction
    is not the same as the 'Area Code' of the 'Card Number'
    (Customer's Area Code) - (use a sliding window).
*/

SELECT
    [inputAtm].[area_code] AS Atm_Area_Code,
    [inputCustomers].[area_code] AS Customer_Area_Code,
    COUNT (*),
    System.Timestamp AS Time
INTO [output]
FROM [input]
INNER JOIN [inputCustomers]
    ON [inputCustomers].[card_number] = [input].[CardNumber]
INNER JOIN [inputAtm]
    ON [inputAtm].[atm_code] = [input].[ATMCode]
WHERE [inputAtm].[area_code] != [inputCustomers].[area_code]
GROUP BY [inputAtm].[area_code],
         [inputCustomers].[area_code],
         SlidingWindow(hour, 1)
```

**Output_8 (SlidingWindow(hour, 1)):**

```
[
  {
    "Atm_Area_Code":5,
    "Customer_Area_Code":7,
    "COUNT":523,
    "Time":"2019-06-30T20:18:13.5610000Z"
  },
  {
    "Atm_Area_Code":2,
    "Customer_Area_Code":1,
    "COUNT":333,
    "Time":"2019-06-30T20:18:13.5760000Z"
  }
]
```

# Bibliography

[1] D. Chatziantoniou and S. Safras, *Course Notes on Big Data Management Systems.*

[2] "Microsoft Azure: Stream Analytics Documentation." https://docs.microsoft.com/en-us/azure/stream-analytics/.