**Pete Champlin**
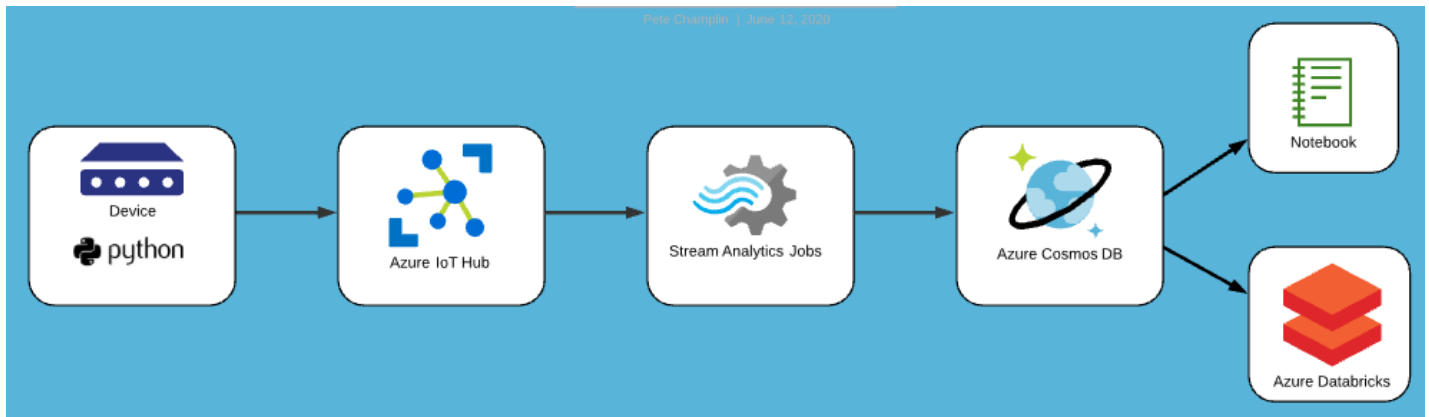
**Big Data 230A**

**Week 8 Assignment (Alternate)**
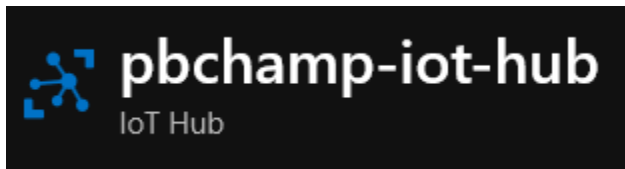
Using my own Azure subscription I followed this documentation - https://docs.microsoft.com/en-us/azure/iot-hub/quickstart-send-telemetry-python - to create a simulated IoT device (using Python), consume the output through Azure IoT Hub and a Stream Analytics job, and write the data to Cosmos DB. I was able to query the data in Data Explorer, an Azure Notebook, and a Databricks (Community Edition) Notebook.
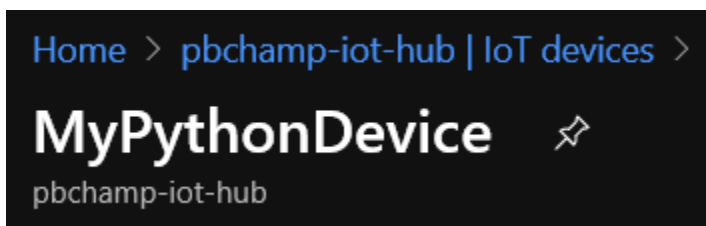
Here is the basic data flow:



Below are the steps followed:
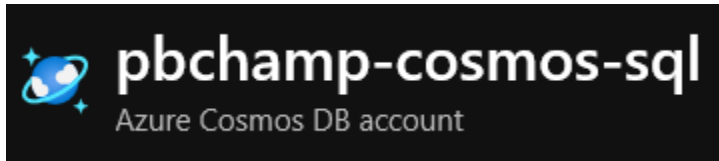
- Create an Azure IoT Hub:



- Associate simulated (Python) IoT device to IoT Hub using Azure CLI

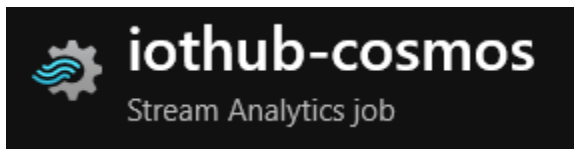  **az iot hub device-identity create --hub-name pbchamp-iot-hub --device-id MyPythonDevice**

- Create a Cosmos DB (SQL) account, container, and database



| Containers | | |
|---|---|---|
| ID | Database | Throughput (RU/s) |
| pbchamp-iot-data-container | pbchamp-iot-data | 400 (Shared) |

- Create a Stream Analytics job



- Create a Stream Analytics Input (Iot Hub)

**Input details**
iothub-input

🔌 Test   🗑 Delete

Input alias
iothub-input

⦿ Provide IoT Hub settings manually
◯ Select IoT Hub from your subscriptions

Subscription
Subscription information not needed

IoT Hub * ⓘ
pbchamp-iot-hub

- Create a Stream Analytics Output (Cosmos DB)

**Output details**
cosmos-output

🔌 Test   🗑 Delete

Output alias
cosmos-output

⦿ Provide Cosmos DB settings manually
◯ Select Cosmos DB from your subscriptions

Subscription
Subscription information not needed

Account id * ⓘ
pbchamp-cosmos-sql

- Create a Stream Analytics Query (Python device is running)
  - Note: Square brackets needed around input and output since they include a dash



- Query Cosmos DB using Data Explorer



- Query Cosmos DB using Azure Notebook



Obviously, much more complicated queries could be performed: https://docs.microsoft.com/en-us/azure/stream-analytics/stream-analytics-stream-analytics-query-patterns

- Query Cosmos DB using a Databricks/Spark (Community Edition) Notebook
  - Followed instructions here to set up install and configure the pyDocumentDB library (https://medium.com/big-data-and-cloud-a-z/connecting-databricks-to-cosmos-db-and-accessing-data-in-it-a1e1b9263944)

```
1  # Set query parameter
2  #querystr = "SELECT * from c"
3  querystr = "SELECT c.EventEnqueuedUtcTime, c.humidity, c.temperature FROM c order by c.EventEnqueuedUtcTime"
4
5  # Query for partitioned collections
6  query = client.QueryDocuments(collLink, querystr, options= { 'enableCrossPartitionQuery': True }, partition_key= True)
7  df = spark.createDataFrame(list(query))
8  display(df)
```

▶ (3) Spark Jobs

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [EventEnqueuedUtcTime: string, humidity: double ... 1 more fields]

| EventEnqueuedUtcTime | humidity | temperature |
|---|---|---|
| 2020-06-10T02:11:31.3040000Z | 73.99293013367489 | 27.6227960986376 |
| 2020-06-10T02:11:32.3350000Z | 76.10465538958384 | 23.4186172214242 |
| 2020-06-10T02:11:33.3980000Z | 63.702010936325486 | 20.4967233510979 |
| 2020-06-10T02:11:34.4460000Z | 78.6841953779532 | 30.5585858498586 |
| 2020-06-10T02:11:35.4770000Z | 67.48400840022711 | 25.6069807718323 |
| 2020-06-10T02:11:36.5090000Z | 71.58607486602933 | 25.8721849114242 |
| 2020-06-10T02:11:37.5570000Z | 62.18404913729657 | 26.6279205219442 |