# Project – Designing Different Data Solutions

**Objective**: To build a data solutions using different Azure Technologies for different data realated challenges including Cosmos DB, Dedicated SQL Pool, Spark Pool, Event Hubs & Stream Analytics.

**Project outline:**

- Technologies used -
    - Azure Cosmos DB – NoSQL database for real-time data storage.
    - Azure Dedicated SQL Pool – Data warehousing for structured data.
    - Azure Synapse Spark Pool – Large-scale data processing with Spark.
    - Azure Event Hub – Real-time data ingestion.
    - Azure Stream Analytics – Stream processing for real-time analytics.

- Prerequisites – Setting up technologies for each and using the created ADLS storage and container 'project' for the source data – business_employment.csv
- Data Ingestion & Preparation – Used Event Hubs to ingest the data in real time
- Data Transformation & Cleaning – Used Stream jobs, Spark pool, Dedicated SQL Pool
- Conclusion

**Dedicated SQL Pool:** Perform large-scale batch processing and analytics on structured data.

- Set up the Dedicated SQL Pool for scalable data storage.
- Develop SQL queries for data extraction and transformation.
- Implement keys (alternate and surrogate) and support Slowly Changing Dimensions (SCD).

Created an External Data Source - 'project_src' & created a file format for our 'CSV'

```
--- Created External Data Source---
CREATE EXTERNAL DATA SOURCE project_src
WITH(
LOCATION = 'abfss://project@synapsestorageadls12.dfs.core.windows.net/Data'
)

---Created a CSV File Format----
IF NOT EXISTS (SELECT * FROM sys.external_file_formats WHERE name ='csv_file_format')
CREATE EXTERNAL FILE FORMAT csv_file_format

WITH (

FORMAT_TYPE = DELIMITEDTEXT,
FORMAT_OPTIONS (

FIELD_TERMINATOR = ','

, STRING_DELIMITER = '"'
, First_Row = 2
, USE_TYPE_DEFAULT = FALSE
, Encoding = 'UTF8'
)
```

Created Schema Name 'Staging' & created an external table with the column names from the dataset 'business_employment.csv' present in ADLS storage

```
---Created Schema----
CREATE SCHEMA Staging
GO
---Created an External Table----
CREATE EXTERNAL TABLE Staging.externaltable
    (    Series_reference NVARCHAR(100),
         Period NVARCHAR(50),
         Data_value FLOAT,
         Suppressed NVARCHAR(10),
         STATUS NVARCHAR(50),
         UNITS NVARCHAR(50),
         Magnitude FLOAT,
         Subject NVARCHAR(100),
         [Group] NVARCHAR(100),
         Series_title_1 NVARCHAR(100),
         Series_title_2 NVARCHAR(100),
         Series_title_3 NVARCHAR(100),
         Series_title_4 NVARCHAR(100),
         Series_title_5 NVARCHAR(100)
        )
    WITH (
            LOCATION = 'business_employment.csv',
            DATA_SOURCE = project_src,
            FILE_FORMAT = csv_file_format
    );
```

This our external table view after running Select* FROM External Table

Created a table 'mytable' and added alternate key

```sql
CREATE TABLE dbo.mytable
(    Series_reference NVARCHAR(100) NOT NULL,
     Period NVARCHAR(50),
     Data_value FLOAT,
     Suppressed NVARCHAR(10),
     STATUS NVARCHAR(50),
     UNITS NVARCHAR(50),
     Magnitude FLOAT,
     Subject NVARCHAR(100),
     [Group] NVARCHAR(100),
     Series_title_1 NVARCHAR(100),
     Series_title_2 NVARCHAR(100),
     Series_title_3 NVARCHAR(100),
     Series_title_4 NVARCHAR(100),
     Series_title_5 NVARCHAR(100)
     )
WITH
(
DISTRIBUTION = REPLICATE,
CLUSTERED COLUMNSTORE INDEX
);

---- Setting up an Alternate Key----

ALTER TABLE dbo.mytable1
ADD CONSTRAINT Series_reference UNIQUE(Series_reference)  NOT ENFORCED;
```

Given some input values to the above created table, see below:

```sql
--- Inserting Values----

INSERT INTO dbo.mytable
(Series_reference, Period, Data_value, Suppressed, STATUS, UNITS, Magnitude, Subject, [Group], Series_title_1, Series_title_2, Series_title_3, Series_title_4, Series_title_5)
VALUES
('SR001', '2024-Q1', 123.45, 'No', 'Active', 'Units', 1.0, 'Subject1', 'Group1', 'Title1', 'Title2', 'Title3', 'Title4', 'Title5');

-- Inserting a second row with a new Series_reference
INSERT INTO dbo.mytable
(Series_reference, Period, Data_value, Suppressed, STATUS, UNITS, Magnitude, Subject, [Group], Series_title_1, Series_title_2, Series_title_3, Series_title_4, Series_title_5)
VALUES
('SR002', '2024-Q2', 567.89, 'No', 'Active', 'Units', 1.1, 'Subject2', 'Group2', 'Title6', 'Title7', 'Title8', 'Title9', 'Title10');

INSERT INTO  dbo.mytable
(Series_reference, Period, Data_value, Suppressed, STATUS, UNITS, Magnitude, Subject, [Group], Series_title_1, Series_title_2, Series_title_3, Series_title_4, Series_title_5)
VALUES
('SR001', '2024-Q3', 789.01, 'Yes', 'Inactive', 'Units', 2.0, 'Subject3', 'Group3', 'Title11', 'Title12', 'Title13', 'Title14', 'Title15');

SELECT * FROM dbo.mytable
```

Below output, alternate key can be seen…

| Series_reference | Period | Data_value | Suppressed | STATUS | UNITS | Magnitude | Subject | Group | Series_ti |
|---|---|---|---|---|---|---|---|---|---|
| SR002 | 2024-Q2 | 567.89 | No | Active | Units | 1.5 | Subject2 | Group2 | Title6 |
| SR001 | 2024-Q1 | 123.45 | No | Active | Units | 1 | Subject1 | Group1 | Title1 |
| SR001 | 2024-Q3 | 780.01 | Yes | Inactive | Units | 2 | Subject3 | Group3 | Title11 |

Now created another table to demonstrate Slowly Changing Dimension Type 2:

Created a staging table as below and copied the data the from ADLS container dataset to the table.

```sql
--------------------------------------------
--- SCD-- Created a Staging Table----
CREATE TABLE StagingCustomer (
    Series_reference VARCHAR(50),
    Period VARCHAR(50),
    Data_value FLOAT,
    Suppressed VARCHAR(50),
    STATUS VARCHAR(50),
    UNITS VARCHAR(50),
    Magnitude VARCHAR(50),
    Subject VARCHAR(100),
    [Group] VARCHAR(100),
    Series_title_1 VARCHAR(200),
    Series_title_2 VARCHAR(200),
    Series_title_3 VARCHAR(200),
    Series_title_4 VARCHAR(200),
    Series_title_5 VARCHAR(200)
);

---Copied the values into the table from ADLS container dataset---
COPY INTO StagingCustomer
FROM 'https://synapsestorageadls12.dfs.core.windows.net/project/Data/business_employment.csv'
WITH (
    FILE_TYPE = 'CSV',
    FIELDQUOTE = '"',
    FIELDTERMINATOR = ',',
    ROWTERMINATOR = '0x0A',
    FIRSTROW = 2
);
--- Alternate method can also be used to move the values from external table to our table---
INSERT INTO StagingCustomer
SELECT * FROM Staging.externaltable;
```

| Series_reference | Period | Data_value | Suppressed | STATUS | UNITS | Magnitude | Subject | Group | Series_ti |
|---|---|---|---|---|---|---|---|---|---|
| BDCQ.SEA2BT | 2017.06 | (NULL) | Y | C | Value | 8 | Business Data ... | (NULL) | (NULL) |
| BDCQ.SEA2DS | 2017.03 | (NULL) | Y | R | Value | 6 | Business Data ... | (NULL) | (NULL) |
| BDCQ.SED1RCS | 2019.03 | 983377 | (NULL) | R | Number | 0 | Business Data ... | (NULL) | (NULL) |
| BDCQ.SED1RDT | 2018.09 | 128013 | (NULL) | R | Number | 0 | Business Data ... | (NULL) | (NULL) |
| BDCQ.SED2RPA | 2020.09 | 615.031032 | (NULL) | R | Value | 6 | Business Data ... | (NULL) | (NULL) |
| BDCQ.SED3RCA | 2015.09 | 139822 | (NULL) | F | Number | 0 | Business Data ... | (NULL) | (NULL) |
| BDCQ.SED3RGA | 2016.03 | 44214 | (NULL) | F | Number | 0 | Business Data ... | (NULL) | (NULL) |

✓ 00:00:00 Query executed successfully

Now we are creating  a Dimension table as below:

```sql
--- Created another Dimension Table with Surrogate Key, Start Date, End Date to Show type SCD---
CREATE TABLE CustomerDimension (
    SurrogateKey INT IDENTITY(1,1) NOT NULL,  -- Surrogate key
    Series_reference VARCHAR(50), -- Business key
    Period VARCHAR(50),
    Data_value FLOAT,
    Suppressed VARCHAR(50),
    STATUS VARCHAR(50),
    UNITS VARCHAR(50),
    Magnitude VARCHAR(50),
    Subject VARCHAR(100),
    [Group] VARCHAR(100),
    Series_title_1 VARCHAR(200),
    Series_title_2 VARCHAR(200),
    Series_title_3 VARCHAR(200),
    Series_title_4 VARCHAR(200),
    Series_title_5 VARCHAR(200),
    StartDate DATETIME,                     -- Start date for validity
    EndDate DATETIME,                       -- End date (null if current)
    IsCurrent BIT                           -- Flag to mark the current record (1 = current, 0 = historical)
);

ALTER TABLE CustomerDimension
ADD CONSTRAINT PK_Customerdimension_SurrogateKey PRIMARY KEY NONCLUSTERED (SurrogateKey) NOT ENFORCED;
```

In order to perform Type 2 SCD,
We will compare the data in the staging table with the dimension table,

update existing records in the dimension table as historical by updating the IsCurrent flag and setting the EndDate.

```sql
-- Update existing records in the dimension table
UPDATE CustomerDimension
SET
    IsCurrent = 0,        -- Mark as historical
    EndDate = GETDATE()   -- Set the end date to the current date
FROM
    CustomerDimension dd
INNER JOIN Staging.externaltable sd
    ON dd.Series_reference = sd.Series_reference
    AND dd.Period = sd.Period
WHERE
    dd.IsCurrent = 1 -- Only update current records
    AND (
        dd.Data_value <> sd.Data_value OR
        dd.Suppressed <> sd.Suppressed OR
        dd.STATUS <> sd.STATUS OR
        dd.UNITS <> sd.UNITS OR
        dd.Magnitude <> sd.Magnitude OR
        dd.Subject <> sd.Subject OR
        dd.[Group] <> sd.[Group] OR
        dd.Series_title_1 <> sd.Series_title_1 OR
        dd.Series_title_2 <> sd.Series_title_2 OR
        dd.Series_title_3 <> sd.Series_title_3 OR
        dd.Series_title_4 <> sd.Series_title_4 OR
        dd.Series_title_5 <> sd.Series_title_5
    );
```

Once the existing records are marked as historical, we will insert the updated records from the staging table into the dimension table as new current records (IsCurrent = 1).

```sql
    -- Insert new records for changed rows or new rows
INSERT INTO CustomerDimension (
    Series_reference,
    Period,
    Data_value,
    Suppressed,
    STATUS,
    UNITS,
    Magnitude,
    Subject,
    [Group],
    Series_title_1,
    Series_title_2,
    Series_title_3,
    Series_title_4,
    Series_title_5,
    StartDate,
    EndDate,
    IsCurrent
)
SELECT
    sd.Series_reference,
    sd.Period,
    sd.Data_value,
    sd.Suppressed,
    sd.STATUS,
    sd.UNITS,
    sd.Magnitude,
    sd.Subject,
    sd.[Group],
    sd.Series_title_1,
    sd.Series_title_2,
    sd.Series_title_3,
    sd.Series_title_4,
    sd.Series_title_5,
    GETDATE(),  -- StartDate for new records
    NULL,       -- EndDate (NULL for current records)
    1           -- IsCurrent = 1 (this is the current record)
FROM
    Staging.externaltable sd
LEFT JOIN CustomerDimension dd
    ON sd.Series_reference = dd.Series_reference
    AND sd.Period = dd.Period

    AND dd.IsCurrent = 1
WHERE
    dd.Series_reference IS NULL -- Insert new rows
    OR (
        dd.Data_value <> sd.Data_value OR
        dd.Suppressed <> sd.Suppressed OR
        dd.STATUS <> sd.STATUS OR
        dd.UNITS <> sd.UNITS OR
        dd.Magnitude <> sd.Magnitude OR
        dd.Subject <> sd.Subject OR
        dd.[Group] <> sd.[Group] OR
        dd.Series_title_1 <> sd.Series_title_1 OR
        dd.Series_title_2 <> sd.Series_title_2 OR
        dd.Series_title_3 <> sd.Series_title_3 OR
        dd.Series_title_4 <> sd.Series_title_4 OR
        dd.Series_title_5 <> sd.Series_title_5
    );

Select * From CustomerDimension
```

We get the below output with Start Date, End Date and Surrogate Key:

Did some transformation using where clause and selected few columns only

```sql
Select * From CustomerDimension
WHERE Series_title_1 = 'Filled jobs'
```

```sql
Select Series_reference, Period, Data_value, Series_title_1, Series_title_2
FROM CustomerDimension
WHERE Series_title_1 = 'Filled jobs'
```
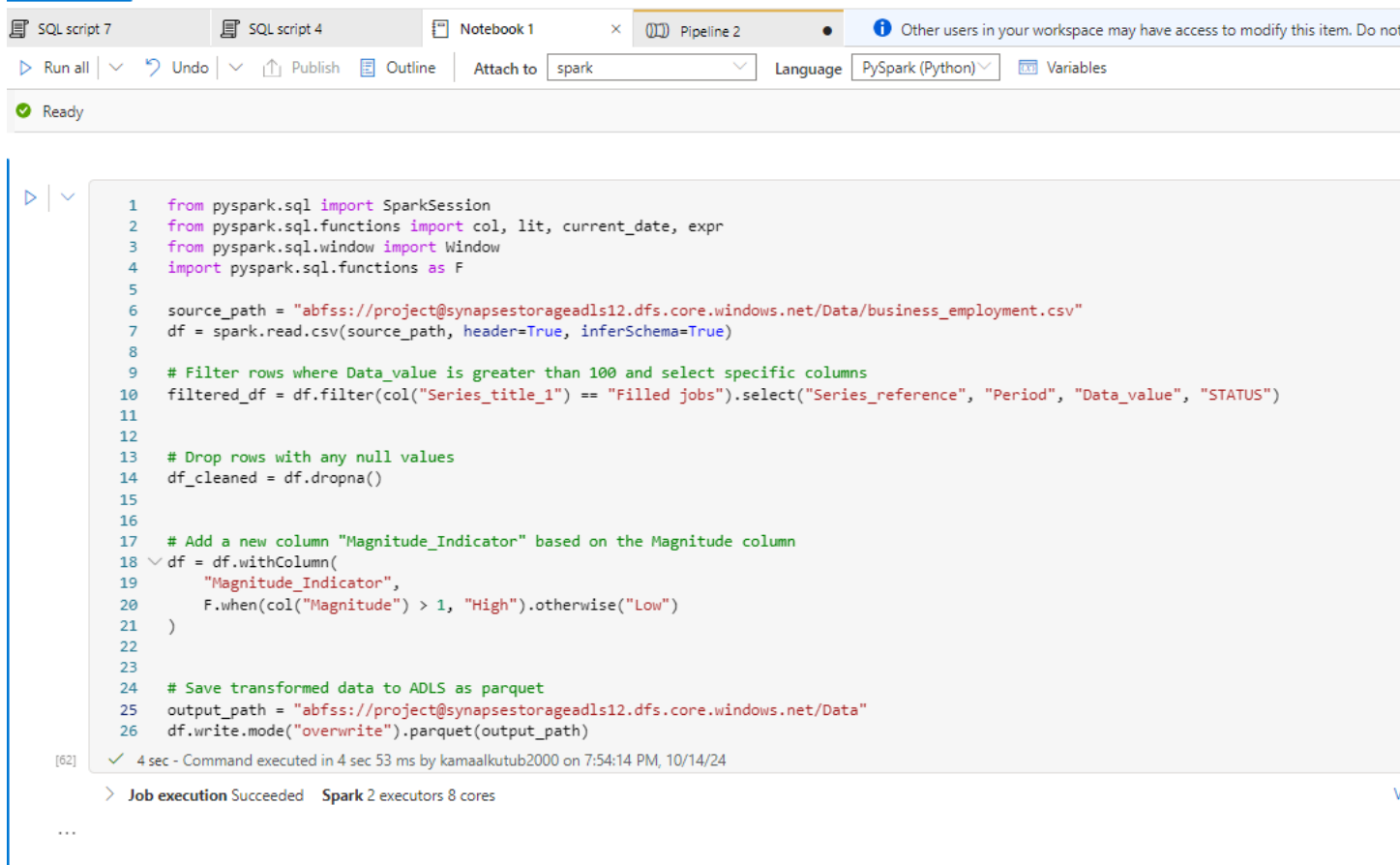
| Series_reference | Period | Data_value | Series_title_1 | Series_title_2 |
|---|---|---|---|---|
| BDCQ.SEA1AA | 2011.06 | 80078 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2016.03 | 99291 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2020.12 | 103593 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2011.09 | 78324 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2016.06 | 88716 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2021.03 | 102002 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2011.12 | 85850 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2016.09 | 85933 | Filled jobs | Agriculture, Forestry and Fishing |
| BDCQ.SEA1AA | 2021.06 | 93431 | Filled jobs | Agriculture, Forestry and Fishing |

**Apache Spark Pool:** Enable large-scale data processing and transformations using Spark.

- Create and configure the Spark Pool.
- Write Python scripts for data transformation.
- Use Spark to integrate and process data in near real-time or batch mode.

Created spark pool and launched the Notebook,

Gave few queries and transformed the data and imported it to a pipeline



```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, lit, current_date, expr
from pyspark.sql.window import Window
import pyspark.sql.functions as F

source_path = "abfss://project@synapsestorageadls12.dfs.core.windows.net/Data/business_employment.csv"
df = spark.read.csv(source_path, header=True, inferSchema=True)

# Filter rows where Data_value is greater than 100 and select specific columns
filtered_df = df.filter(col("Series_title_1") == "Filled jobs").select("Series_reference", "Period", "Data_value", "STATUS")


# Drop rows with any null values
df_cleaned = df.dropna()


# Add a new column "Magnitude_Indicator" based on the Magnitude column
df = df.withColumn(
    "Magnitude_Indicator",
    F.when(col("Magnitude") > 1, "High").otherwise("Low")
)


# Save transformed data to ADLS as parquet
output_path = "abfss://project@synapsestorageadls12.dfs.core.windows.net/Data"
df.write.mode("overwrite").parquet(output_path)
```

[62]  ✓ 4 sec - Command executed in 4 sec 53 ms by kamaalkutub2000 on 7:54:14 PM, 10/14/24

> Job execution Succeeded   Spark 2 executors 8 cores

Imported the Notebook and ran the pipeline along with data exists in Get Metadata.



The Output of the pipeline can be seen in the ADLS Container as below:

**Azure Event Hubs:** To Ingest real-time data from various sources such as IoT devices and applications.

- Set up Event Hub for real-time data ingestion.
- Integrate Event Hub with Stream Analytics for real-time processing.

Created a new event hub workspace and created a new event:



Used the sample event to send the data from vehicle toll booth:

## Send events

These events will be sent to event hub newevent

Select Dataset *

Vehicle toll booth

(i) Following properties from Vehicle toll booth dataset are going to be dynamic:
**carModel.make, carModel.model, licensePlate, state, tag, tollAmount, tollId**

Sample event

```
 1  {
 2      "entryTime": "2023-05-09T04:49:15.0189703Z",
 3      "carModel": {
 4          "make": "Honda",
 5          "model": "Civic",
 6          "vehicleType": 1,
 7          "vehicleWeight": 0
 8      },
 9      "state": "NJ",
10      "tollAmount": 10,
11      "tag": 584666966,
12      "tollId": 4,
13      "licensePlate": "A9T IL7N",
14      "eventProcessedUtcTime": "2023-05-09T04:52:54.3513112Z",
15      "partitionId": 0,
16      "eventEnqueuedUtcTime": "2023-05-09T04:49:16.0750000Z"
17  }
```
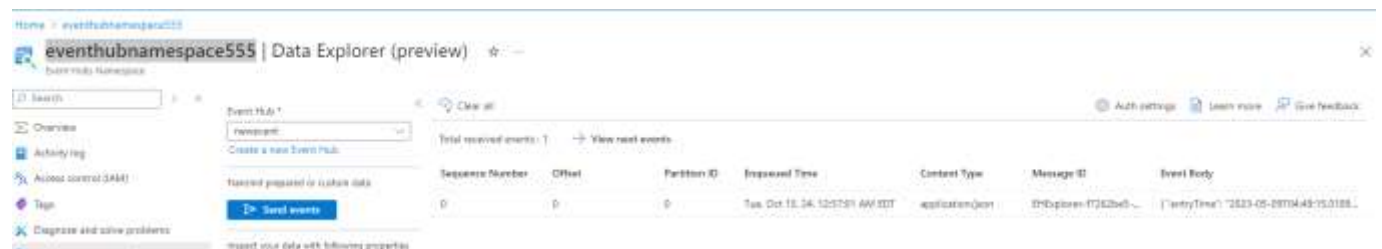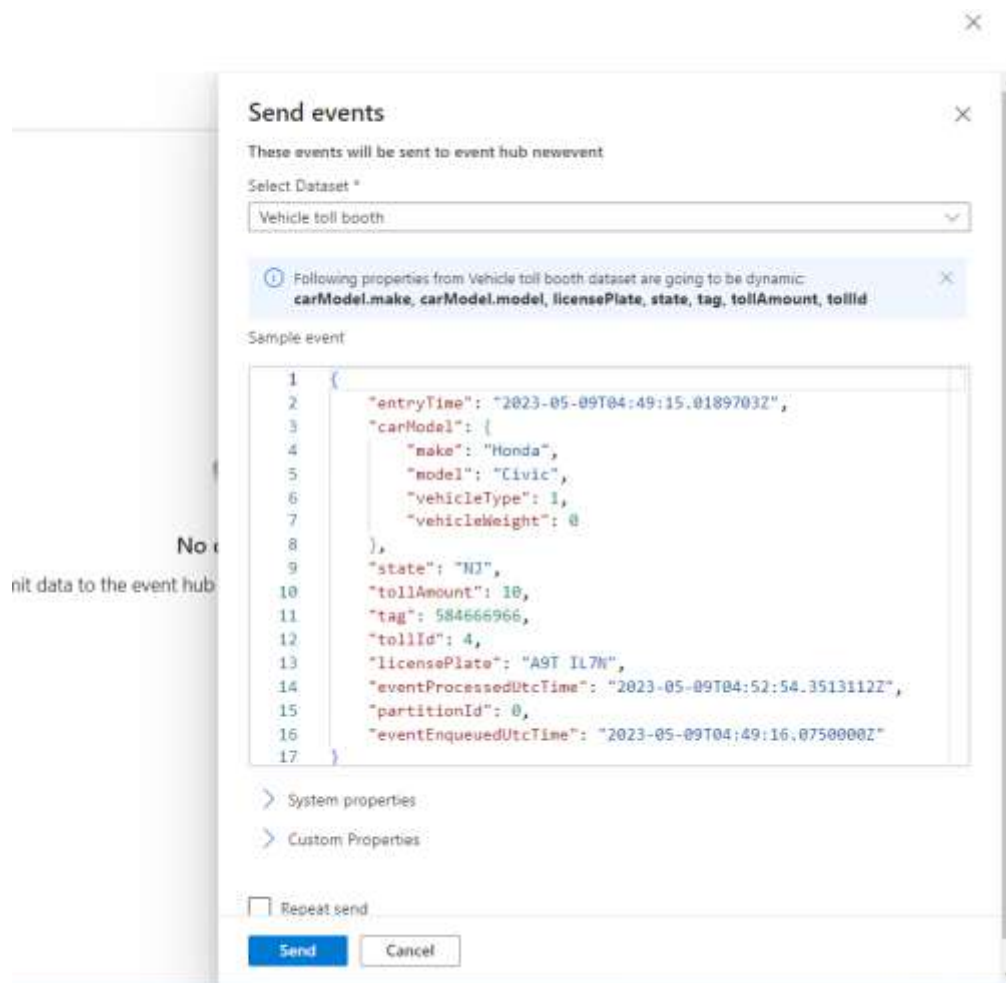
> System properties

> Custom Properties

[ ] Repeat send

**Send**    Cancel

---

Home > eventhubnamespace555

# eventhubnamespace555 | Data Explorer (preview)

Event Hubs Namespace

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Event Hub * | | Clear all | | | | Auth settings  Learn more  Give feedback | |
| newevent | | Total received events: 1  → View next events | | | | | |
| Create a new Event Hub | | | | | | | |
| Received prepared or custom data | Sequence Number | Offset | Partition ID | Enqueued Time | Content Type | Message ID | Event Body |
| ⊳ Send events | 0 | 0 | 0 | Tue, Oct 15, 24, 12:57:51 AM EDT | application/json | EHExplorer-f7262be5-... | {"entryTime": "2023-05-09T04:49:15.0189... |
| Inspect your data with following properties | | | | | | | |

**Azure Stream Analytics:** To Perform real-time stream processing and transformation of data from Event Hub and route it to Dedicated SQL Pool.

- Set up Stream Analytics jobs for data transformation.
- Define queries to process data from Event Hub and route output to Cosmos DB or Dedicated SQL Pool.
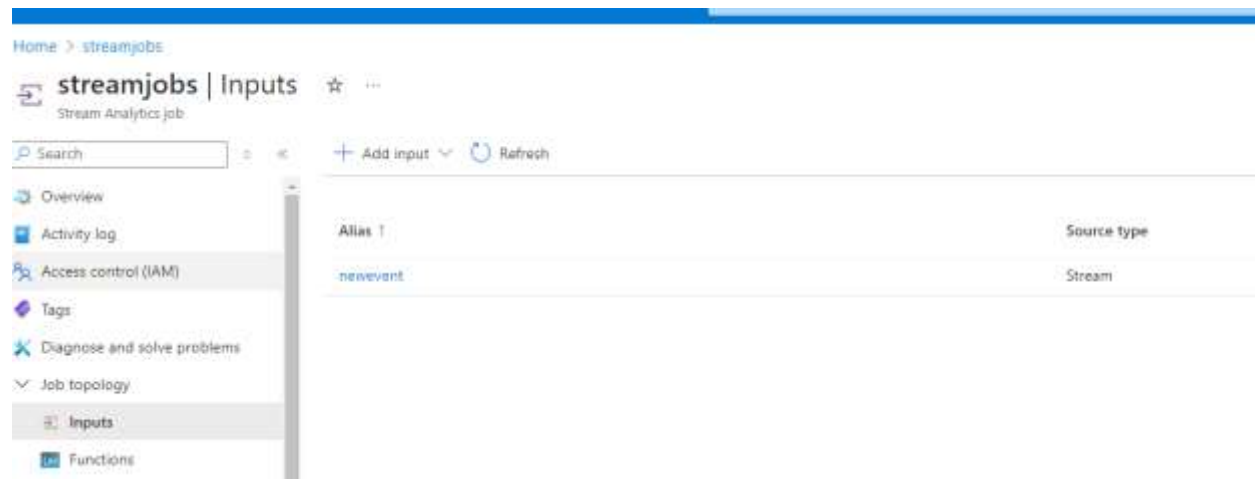
Created a stream job for the stream analytics real time data transformation
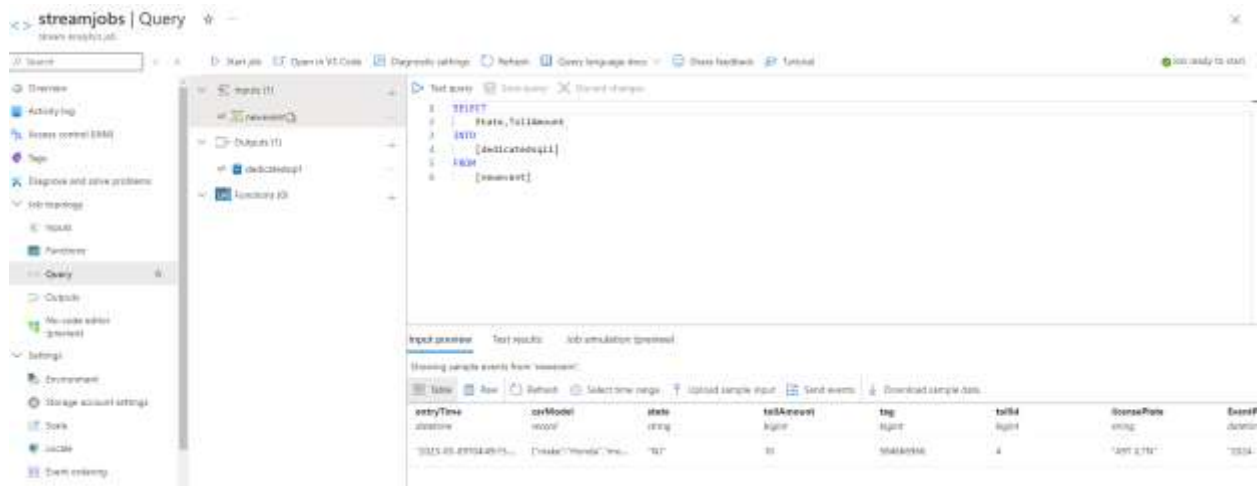
Created the below for stream jobs:

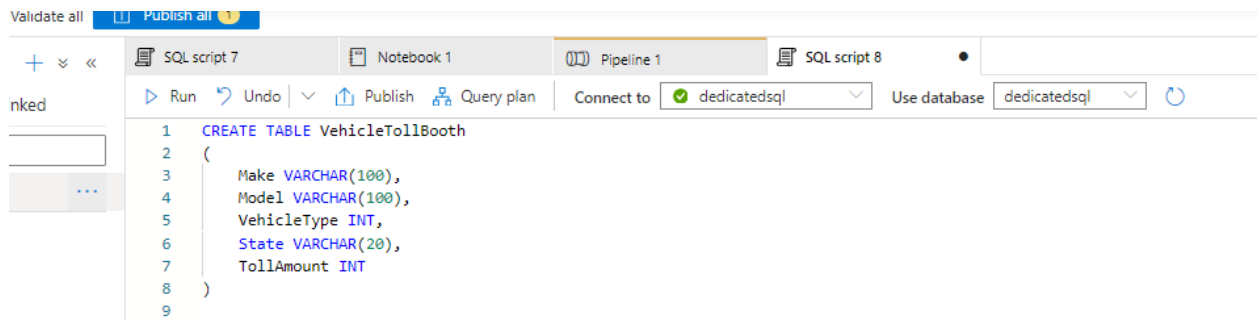Input – to fetch from our Event Hubs created event

Output – To our dedicated sql Pool

Query – Gave the query to run the stream jobs for gathering the data only for two columns.
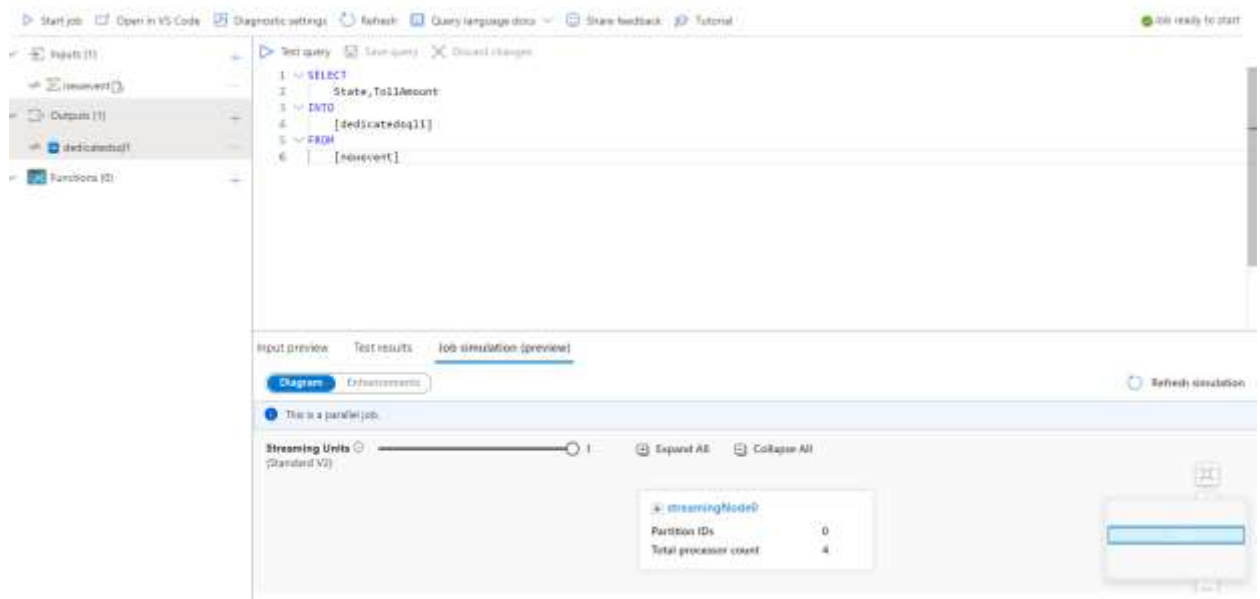
Also created a table 'vehicletollbooth' in dedicated pool:



```
1    CREATE TABLE VehicleTollBooth
2    (
3        Make VARCHAR(100),
4        Model VARCHAR(100),
5        VehicleType INT,
6        State VARCHAR(20),
7        TollAmount INT
8    )
9
```

And final ran the below query, saved the query, hit the start job.

**Error Log:**

| Error Faced | Work around |
|---|---|
| Spark Pool Memory Error | Created with different number of nodes, raised a ticket to Microsoft and later after a few days started working |
| Dedicated SQL Pool ingestion Error | Autoresolved after a few attempts |
| SQL Alternate Key Error | Given Not Enforced to resolve |
| Stream Jobs error for output not found | Changed from Managed Identity and given SQL authentication , then connection was made succesfully |
| Mege Table Query Error | Used Alternative Update along with joins to compare the staging and dimension table |

**Conclusion:** This project successfully demonstrates how **Azure's data services** can be leveraged to solve complex data challenges involving real-time data ingestion, processing, and analysis. By using a combination of **Cosmos DB**, **SQL Pool**, **Event Hub**, **Stream Analytics**, and **Spark Pool**, we built a robust data architecture capable of addressing the needs of modern data-driven applications, including real-time analytics and large-scale batch processing.