



Azure Data Services in a Day

Implement a Modern Data Platform Architecture

Table of Contents

Authors.....	4
Version History.....	4
Overview	5
Document Structure	5
Data Source References.....	5
Lab Prerequisites.....	6
Lab Guide	7
Lab 0: Provision Azure Resources	8
Lab 1: Load Data into Azure SQL Data Warehouse using Azure Data Factory Pipelines	8
Lab 2: Transform Big Data using Azure Data Factory and Azure SQL Data Warehouse	8
Lab 3: Explore Big Data using Azure Databricks.....	9
Lab 4: Add AI to your Big Data Pipeline with Cognitive Services	9
Lab 5: Ingest and Analyse real-time data with Event Hubs and Stream Analytics	9
Lab 0: Provision Azure Resources	11
Download Lab Files from GitHub	11
Prepare your Azure subscription	11
Deploy Azure Services.....	14
Lab 1: Load Data into Azure SQL Data Warehouse using Azure Data Factory Pipelines	17
Lab Architecture.....	17
Connect to MDWDesktop	19
Install required software onto MDWDesktop	20
Restore NYCDatasets database onto MDWSQLServer.....	22

Create Azure SQL Data Warehouse database objects	24
Configure the Azure Data Factory Self Hosted Integration Runtime	27
Connect to MDWDataGateway and register the Self Hosted Integration Runtime with Azure Data Factory	30
Create Staging Container on Azure Blob Storage	33
Create Azure Data Factory Pipeline to Copy Relational Data	36
Create Linked Service connections	36
Create Source and Destination Data Sets	40
Create and Execute Pipeline	42
Visualize Data with Power BI	45
Lab 2: Transform Big Data using Azure Data Factory and Azure SQL Data Warehouse	46
Lab Architecture.....	46
Create Azure SQL Data Warehouse database objects	47
Create NYCTaxiData Container on Azure Blob Storage	52
Create Linked Service connection to MDWResources.....	54
Create Source and Destination Data Sets	56
Create and Execute Pipeline	63
Visualize Data with Power BI	69
Lab 3: Explore Big Data using Azure Databricks.....	71
Lab Architecture.....	71
Create Azure Databricks Cluster	71
Create an Azure Databricks Notebook.....	73
Lab 4: Add AI to your Big Data Pipeline with Cognitive Services	80
Lab Architecture.....	80

Create NYCIImages and NYCIImageMetadata Containers in Azure Blob Storage	81
Create CosmosDB database and collection	83
Import Databricks Notebook to Invoke Computer Vision Cognitive Services API	85
Create Databricks Linked Service in Azure Data Factory	90
Create CosmosDB Linked Service in Azure Data Factory	93
Create Azure Data Factory data sets	95
Create Azure Data Factory pipeline to generate and save image metadata to Cosmos DB.	102
Explore Image Metadata Documents in CosmosDB	107
Visualize Data with Power BI	110
Lab 5: Ingest and Analyse real-time data with Event Hubs and Stream Analytics	112
Lab Architecture.....	112
Create NYCTweets Container in Azure Blob Storage	113
Create and Configure Event Hubs.....	115
Create Azure Logic App to Read #NYC Tweets and post them to Event Hubs	117
Create and Configure Stream Analytics	123
Create Power BI Dashboard to Visualise Real-Time Data	127
Create Power BI Report to Visualise Real-Time Data	132
References	135
Feedback.....	135
Disclaimer.....	136

Authors

Fabio Braga

Cloud Solution Architect – Data Platform

fabio.braga@microsoft.com

Rod Colledge

Cloud Solution Architect – Data Platform

rodcoll@microsoft.com

Version History

Version	Date	Description	Author
0.1	23/04/2019	Lab Complete. Ready for Review.	Fabio Braga
0.2	25/04/2019	Review	Rod Colledge
0.3	06/05/2019	Added Lab Architecture diagram to each Lab. Fixes from Rod's review. Minor formatting adjustments.	Fabio Braga
0.4	07/05/2019	Fixed Lab2 inconsistencies and added new screenshots.	Fabio Braga

Overview

Today you will learn about the main concepts related to advanced analytics and Big Data processing and how Azure Data Services can be used to implement a modern data warehouse architecture. You will understand what Azure services you can leverage to establish a solid data platform to quickly ingest, process and visualise data from a large variety of data sources. The reference architecture you will build as part of this exercise has been proven to give you the flexibility and scalability to grow and handle large volumes of data and keep an optimal level of performance.

In the exercises in this lab you will build data pipelines using data related to New York City. The workshop was designed to progressively implement an extended modern data platform architecture starting from a traditional relational data pipeline. Then we introduce big data scenarios with large files and distributed computing. We add non-structured data and AI into the mix and finish with real-time streaming analytics. You will have done all of that by the end of the day.

Document Structure

This document contains detailed step-by-step instructions on how to implement a Modern Data Platform architecture using Azure Data Services. It's recommended you carefully read the detailed description contained in this document for a successful experience with all Azure services.

You will see the label **IMPORTANT** whenever there is a critical step to the lab. Please pay close attention to the instructions given.

Data Source References

New York City data used in this lab was obtained from the New York City Open Data website: <https://opendata.cityofnewyork.us/>. The following datasets were used:

- NYPD Motor Vehicle Collisions: <https://data.cityofnewyork.us/Public-Safety/NYPD-Motor-Vehicle-Collisions/h9gi-nx95>
- TLC Yellow Taxi Trip Data: <https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

Lab Prerequisites

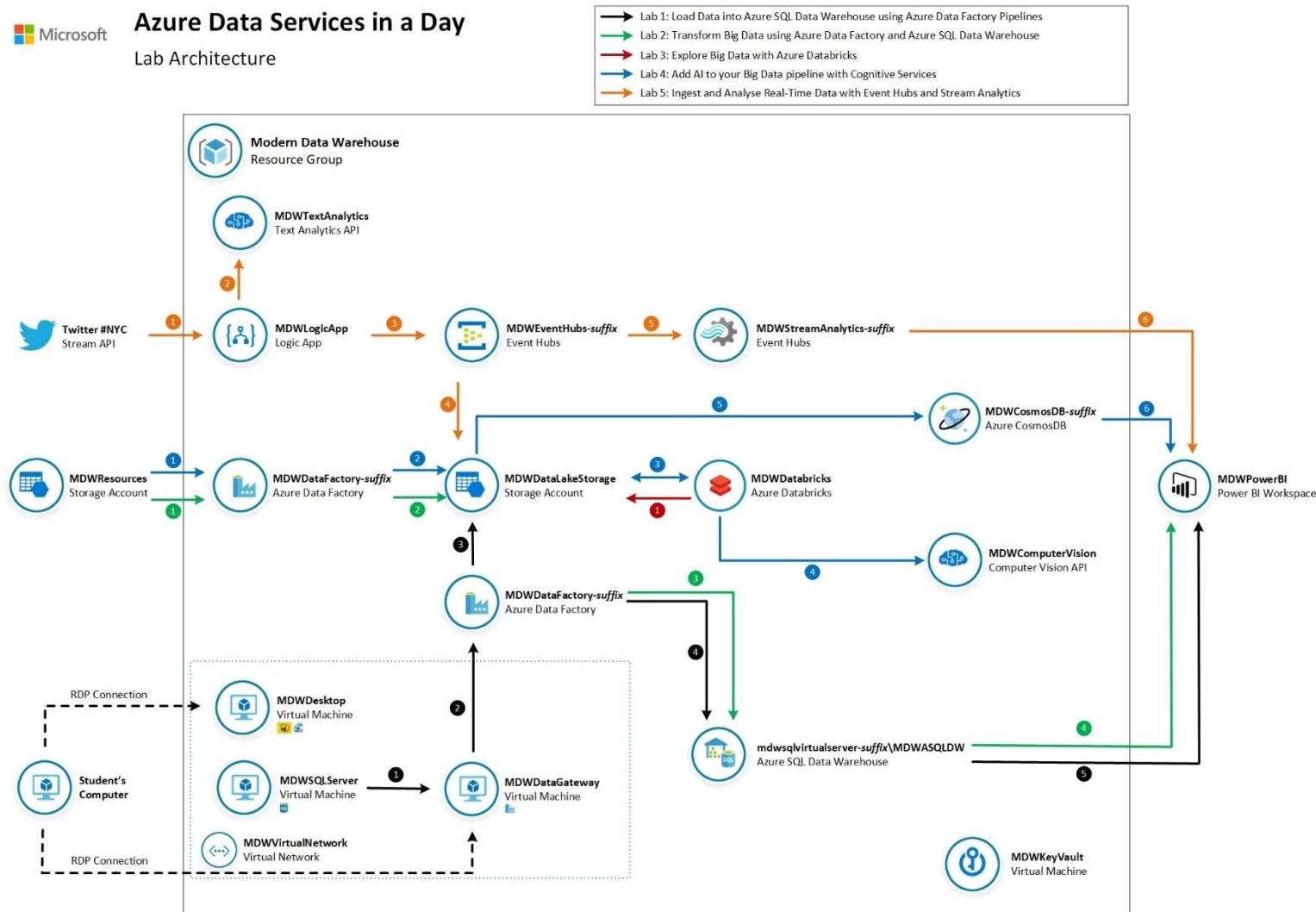
The following prerequisites must be completed before you start these labs:

- You must be connected to the internet;
- You must have an Azure account with administrator-level access to your subscription. If you don't have an account, you can sign up for free following the instructions here: <https://azure.microsoft.com/en-au/free/>
- Download Lab files from <insert GitHub link here> and save them in the local folder C:\ADSIAD\LabFiles;
- Lab 5 requires you to have a Twitter account. If you don't have an account you can sign up for free following the instructions here: <https://twitter.com/signup>.
- Lab 5 requires you to have a Power BI Pro account. If you don't have an account you can sign up for a 60-day trial for free here: <https://powerbi.microsoft.com/en-us/power-bi-pro/>
- You must deploy all Azure services required in each lab by running a PowerShell script as administrator on your local computer. We recommend you follow the instructions defined in **Lab 0 – Provision Azure Resources** for an automated deployment through ARM Templates and PowerShell. You can also deploy the resources manually through the Azure Portal by following the instructions in **Appendix 1 – Manually Provision Azure Resources**.
- The ARM template used to automate the deployment of Azure services attempts to deploy all services required in the resource group default region. See more information on Lab 0.

IMPORTANT: A collection of Azure resources will be provisioned in your subscription and they will incur billing costs. The estimated consumption cost for the execution of the labs included in this workshop is around US \$150.00.

Lab Guide

Throughout a series of 5 labs we will progressively implement the architecture referenced below:



Lab 0: Provision Azure Resources

In this lab you will automatically provision all Azure resources required to complete labs 1 though to 5. You will use a pre-defined ARM template with the definition of all Azure services used to ingest, store, process and visualise data. The estimated time to complete this lab is: 30 minutes.

Lab 1: Load Data into Azure SQL Data Warehouse using Azure Data Factory Pipelines

In this lab you will configure the Azure environment to allow relational data to be transferred from a SQL Server 2017 database to an Azure SQL Data Warehouse database using Azure Data Factory. The dataset you will use contains data about motor vehicle collisions that happened in New York City from 2012 to 2019. You will use Power BI to visualise collision data loaded from Azure SQL Data Warehouse.

- 1 Restore SQL Server backup from Azure Storage and Configure Azure Data Factory Self-Hosted Integration Runtime
- 2 Build an Azure Data Factory Pipeline to copy data from a SQL Server table
- 3 Use Azure Storage as a staging area for Polybase
- 4 Load data to a Azure SQL Data Warehouse table using Polybase
- 5 Visualize data from Azure SQL Data Warehouse using Power BI

Lab 2: Transform Big Data using Azure Data Factory and Azure SQL Data Warehouse

In this lab you will use Azure Data Factory to download large data files to your data lake and use Azure SQL Data Warehouse stored procedure to generate a summary dataset and store it in the final table. The dataset you will use contains detailed New York City Yellow Taxi rides for 2018. You will generate a daily aggregated summary of all rides and save the result in your data warehouse. You will then use Power BI to visualise summarised data. The estimated time to complete this lab is: **xx minutes.**

- 1 Build an Azure Data Factory Pipeline to copy big data files from shared Azure Storage
- 2 Save data files to your data lake
- 3 Use Polybase to load data into staging tables in your Azure SQL Data Warehouse.
Call a Stored Procedure to perform data aggregations and save results in the final table.
- 4 Visualize data from your Azure SQL Data Warehouse using Power BI

Lab 3: Explore Big Data using Azure Databricks

In this lab you will use Azure Databricks to explore the New York Taxi data files you saved in your data lake in Lab 2. Using a Databricks notebook you will connect to the data lake and query taxi ride details. The estimated time to complete this lab is: **xx minutes.**

- 1 Build an Azure Databricks notebook to explore the data files you saved in your data lake in the previous exercise. You will use Python and SQL commands to open a connection to your data lake and query data from data files.

Lab 4: Add AI to your Big Data Pipeline with Cognitive Services

In this lab you will use Azure Data Factory to download New York City images to your data lake. Then, as part of the same pipeline, you are going to use an Azure Databricks notebook to invoke Computer Vision Cognitive Service to generate metadata documents and save them back in your data lake. The Azure Data Factory pipeline then finishes by saving all metadata information in a Cosmos DB collection. You will use Power BI to visualise NYC images and their AI-generated metadata. The estimated time to complete this lab is: **xx minutes.**

- 1 Build an Azure Data Factory Pipeline to copy image files from shared Azure Storage
- 2 Save image files to your data lake
- 3 For each image in your data lake, invoke an Azure Databricks notebook that will take the image URL as parameter
- 4 For each image call the Azure Computer Vision Cognitive service to generate image metadata. Metadata files are saved back in your data lake
- 5 Copy metadata JSON documents into your Cosmos DB database
- 6 Visualize images and associated metadata using Power BI

Lab 5: Ingest and Analyse real-time data with Event Hubs and Stream Analytics

In this lab you will use an Azure Logic App to connect to Twitter and generate a stream of messages using the hashtag #NYC. The logic app will invoke the Azure Text Analytics Cognitive service to score Tweet sentiment and send the messages to Event Hubs. You will use Stream Analytics to generate the average Tweet sentiment in the last 60 seconds and send the results to a real-time dataset in Power BI. The estimated time to complete this lab is: **xx minutes.**

- 1 Build an Azure Logic App to invoke the Twitter API and retrieve Tweets with the hashtag #NYC

- 2 For each Tweet, invoke the Azure Text Analytics Cognitive service to detect its sentiment score
- 3 Format and send the Tweet's JSON message to Event Hubs
- 4 Save Tweet messages into your data lake for future analysis (cold path)
- 5 Send stream of Tweet messages to Stream Analytics for real-time analytics (hot path)
- 6 Visualize real-time data generated by Stream Analytics with Power BI

Lab 0: Provision Azure Resources

In this lab you will automatically provision all Azure resources required to complete labs 1 though to 5. We will use a pre-defined ARM template with the definition of all Azure services used to ingest, store, process and visualise data. The estimated time to complete this lab is: 30 minutes.

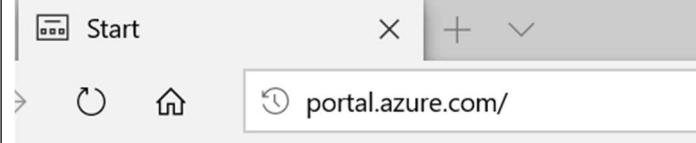
Download Lab Files from GitHub

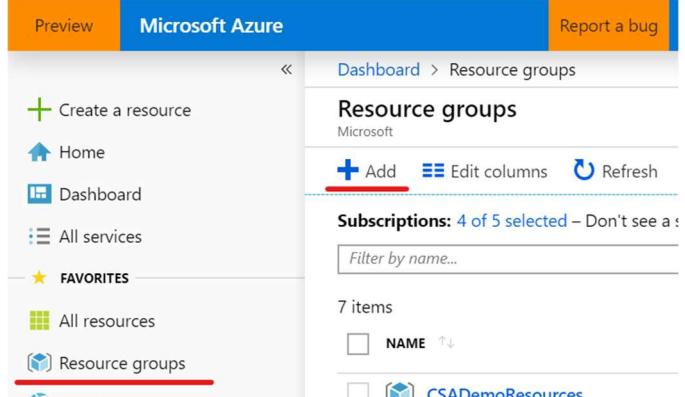
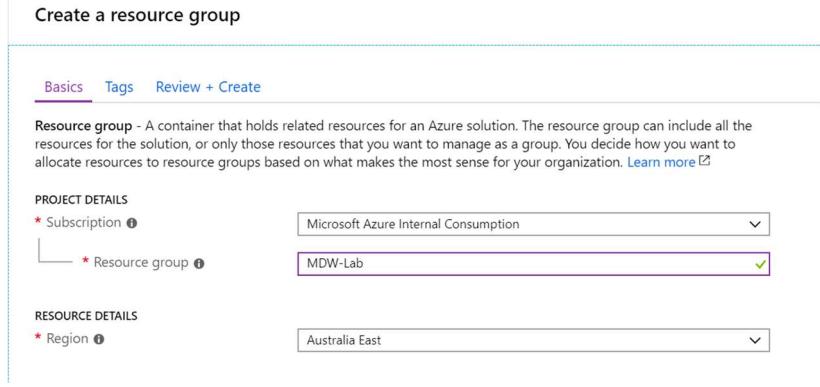
In this section you will download all files used in this lab to deploy Azure Data Services and all other files used in labs 1 through to 5.

Add content here!

Prepare your Azure subscription

In this section you will use the Azure Portal to create a Resource Group that will host the Azure Data Services used in labs 1 through to 5.

IMPORTANT: Execute these steps on your host computer	
1. Open the browser and navigate to https://portal.azure.com	
2. Log on to Azure using your account credentials	

<p>3. Once you have successfully logged on, locate the Favourites menu on the left-hand side panel and click the Resource groups item to open the Resource groups blade.</p> <p>4. On the Resource groups blade, click the “+ Add” button to create a new resource group.</p>	
<p>5. On the Create a resource group blade, select your subscription in Subscription drop down list.</p> <p>6. In the Resource group text box enter “MDW-Lab”</p> <p>7. IMPORTANT: The name of the resource group chosen is *not* relevant to the successful completion of the labs. If you choose to use a different name, then please proceed with the rest of the lab using your unique name for the resource group.</p>	
<p>IMPORTANT</p> <p>The ARM template you will use to deploy the lab components uses the Resource Group region as the default region for all services. To avoid deployment error when services are not available in the region selected, please use one of the recommended regions.</p> <p>8. In the Region drop down list, select one of the regions from the list.</p>	<p>Recommended regions:</p> <ul style="list-style-type: none"> • US East • US East 2 • US Central • US West 2 • Europe West • Asia Southeast • Europe North • US South Central • US West • UK South

	<ul style="list-style-type: none"> ● Australia East ● Canada Central ● Japan East ● Germany Central ● France Central ● India Central ● Brazil South ● Korea Central ● Korea South 						
9. Proceed to create the resource group by clicking Review + Create , and then Create	<p>Create a resource group</p> <p>Basics Tags Review + Create</p> <p>SUMMARY</p> <p>BASICS</p> <table> <tbody> <tr> <td>Subscription</td> <td>Microsoft Azure Internal Consumption</td> </tr> <tr> <td>Resource group</td> <td>MDW-Lab</td> </tr> <tr> <td>Region</td> <td>Australia East</td> </tr> </tbody> </table> <p>Create Previous : Tags</p>	Subscription	Microsoft Azure Internal Consumption	Resource group	MDW-Lab	Region	Australia East
Subscription	Microsoft Azure Internal Consumption						
Resource group	MDW-Lab						
Region	Australia East						

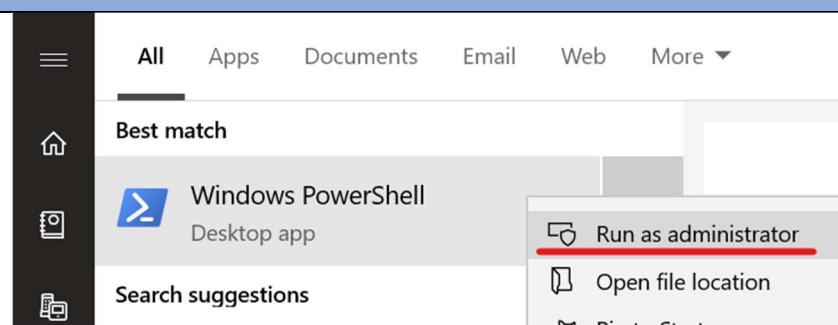
10. Once the MDW-Lab resource group has been created, navigate again to the **Resource groups** menu under **Favourites**.
11. On the **Resource groups** blade, click the MDW-Lab resource group.
12. On the **MDW-Lab Resource Group** blade, take note of your **Subscription ID**. You will need it in the next step to deploy the ARM template.

The screenshot shows the Azure portal interface. On the left, there's a sidebar with options like 'Create a resource', 'Home', 'Dashboard', 'All services', 'FAVORITES' (which contains 'All resources', 'Resource groups', 'App Services', 'Function Apps', and 'SQL databases'), and 'Settings'. The 'Resource groups' option is highlighted with a red underline. The main content area is titled 'MDW-Lab Resource group'. It shows basic information: 'Subscription (change) : Microsoft Azure Internal Consumption', 'Subscription ID : [REDACTED]', and 'Tags (change) : Click here to add tags'. There are buttons for 'Add', 'Edit columns', 'Delete resource group', and 'Refresh'. A search bar at the top says 'Search (Ctrl+J)'. Below the main info, there are sections for 'Overview', 'Activity log', 'Access control (IAM)', 'Tags', 'Events', and 'Quickstart'. A 'Filter by name...' dropdown and a 'NAME' column header are also visible.

Deploy Azure Services

In this section you will use PowerShell scripts and ARM templates to automate the deployment of all Azure Data Services used in labs 1 through to 5.

1. **IMPORTANT:** You must be a member of the local administrators group on the Student's computer.
2. Click the Windows menu logo and start typing **PowerShell until it appears, and then right-click it and select Run as administrator**.
3. Answer **Yes** to the question about allowing this app to make changes to your device.



4. Navigate to the C:\ADSIAD\Lab\Lab0 folder.
 5. Execute the script MDWDeploy.ps1
 6. Answer R to the confirm the execution of the deployment script.

```
Administrator: Windows PowerShell
PS C:\>> cd C:\ADSIAD\Lab\Lab0\
PS C:\ADSIAD\Lab\Lab0>> dir

Directory: C:\ADSIAD\Lab\Lab0

Mode                LastWriteTime         Length Name
----                - - - - - - - - - - - - - - - - - - -
-a---    26/03/2019  8:38 AM           3644 MDWDDeploy.ps1
-a---    22/03/2019  7:19 PM            664 parameters.json
-a---    26/03/2019  8:51 AM          27968 template.json

PS C:\ADSIAD\Lab\Lab0>> .\MDWDDeploy.ps1

Security warning
Run only scripts that you trust. While scripts from the internet can be useful, this script can potentially harm your
computer. If you trust this script, use the Unblock-File cmdlet to allow the script to run without this warning
message. Do you want to run C:\ADSIAD\Lab\Lab0\MDWDDeploy.ps1?
[D] Do not run [R] Run once [S] Suspend [?] Help (default is "D"): R
```

- Paste the **Subscription ID** you retrieved from the previous step in the **subscriptionId** field.
 - Enter “MDW-Lab” in the **resourceGroupName** field.
 - Enter “MDWDDeploy01” in the **deploymentName** field.

```
PS C:\ADSIAD\Lab\Lab0> .\MDWDeploy.ps1

Security warning
Run only scripts that you trust. While scripts from the internet can be useful, this script can potentially harm your
computer. If you trust this script, use the Unblock-File cmdlet to allow the script to run without this warning
message. Do you want to run C:\ADSIAD\Lab\Lab0\MDWDeploy.ps1?
[D] Do not run [R] Run once [S] Suspend [?] Help (default is "D"): R

cmdlet MDWDeploy.ps1 at command pipeline position 1
Supply values for the following parameters:
subscriptionId: 96bd7145-ad7f-445a-9763-862e32480bf1
resourceGroupName: MDW-Lab
deploymentName: MDWDeploy01
```

IMPORTANT: PowerShell AzureRm module is going to be installed as part of the deployment. Please find more info here: [Link](#)

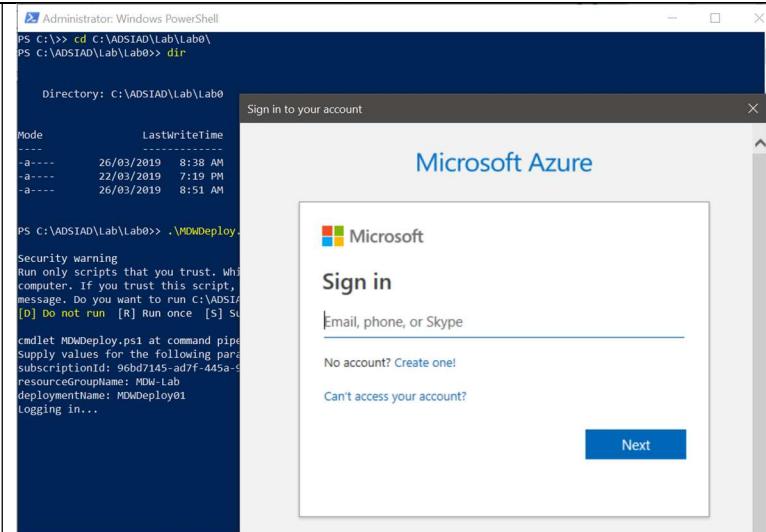
<https://docs.microsoft.com/en-us/powershell/azure/install-az-ps>

10. Answer **Y** to confirm the installation of the NuGet provider.
 11. Answer **A** to proceed with the installation.

```
NuGet provider is required to continue
PowerShellGet requires NuGet provider version '2.8.5.201' or newer to interact with NuGet-based repositories. The NuGet provider must be available in 'C:\Program Files\PackageManagement\ProviderAssemblies' or
'C:\Users\fabraga\AppData\Local\PackageManagement\ProviderAssemblies'. You can also install the NuGet provider by
running 'Install-PackageProvider -Name NuGet -MinimumVersion 2.8.5.201 -Force'. Do you want PowerShellGet to install
and import the NuGet provider now?
[Y] Yes [N] No [S] Suspend [?] Help (default is "Y"): Y

Untrusted repository
You are installing the modules from an untrusted repository. If you trust this repository, change its
InstallationPolicy value by running the Set-PSRepository cmdlet. Are you sure you want to install the modules from
'PSGallery'?
[Y] Yes [A] Yes to All [N] No [L] No to All [S] Suspend [?] Help (default is "N"): A
```

12. Log on to Azure using your credentials.
13. Once credentials are validated the script executions should last for around 15 mins.



14. Once the deployment script finishes successfully, you should be able to see 25 new resources deployed in the MDW-Lab resource group.
15. **IMPORTANT:** Note that some resources had a “-suffix” appended to their names to guarantee uniqueness. Take note of the suffix generated as it will be used very frequently during the lab.

MDW-Lab Resource group

Subscription (change) : Microsoft Azure Internal Consumption Deployments : 1 Succeeded

Subscription ID : 96bd7145-ad7f-445a-9763-862e32480bf1

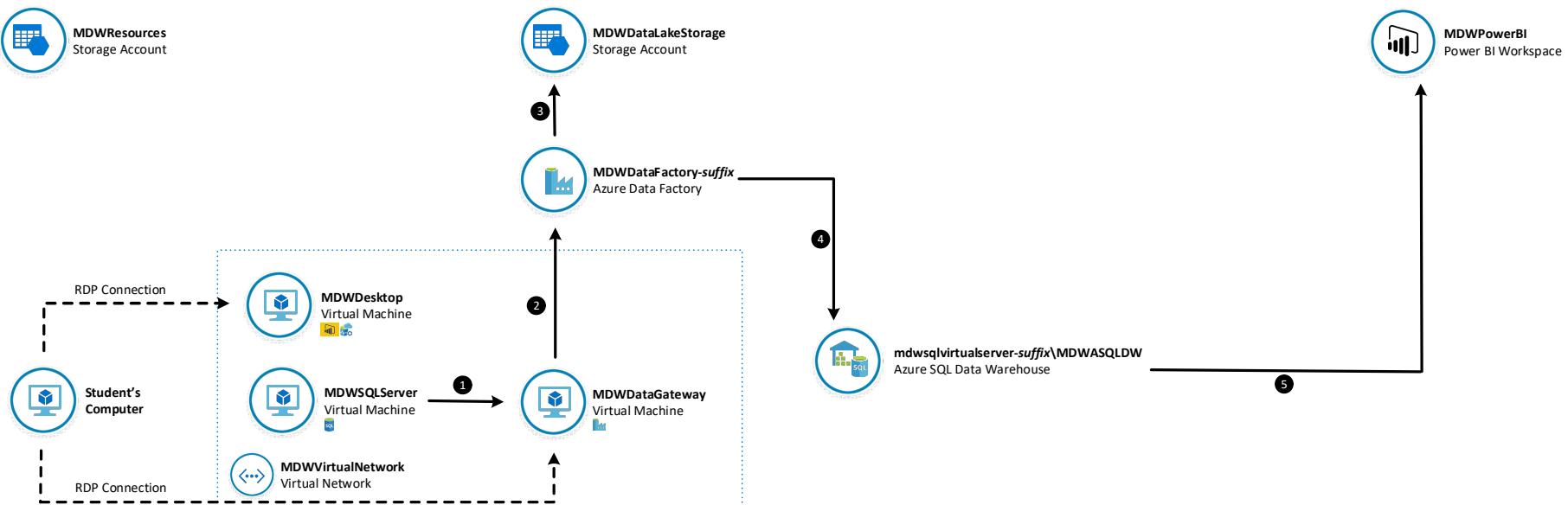
Tags (change) : Click here to add tags

NAME	TYPE	LOCATION	TAGS
MDWComputerVision	Cognitive Services	Australia East	
mdwcosmosdb-zhlnm	Azure Cosmos DB account	Australia East	1
MDWDatabricks	Azure Databricks Service	Australia East	
MDWDataFactory-zhlnm	Data factory (V2)	Southeast Asia	
MDWDataGateway	Virtual machine	Australia East	
MDWDataGateway_OsDisk	Disk	Australia East	
MDWDataGateway-nic	Network interface	Australia East	
mdwdatalakezhlnm	Storage account	Australia East	
MDWDesktop	Virtual machine	Australia East	
MDWDesktop_OsDisk	Disk	Australia East	
MDWDesktop-nic	Network interface	Australia East	
MDWDesktop-publicip	Public IP address	Australia East	

Lab 1: Load Data into Azure SQL Data Warehouse using Azure Data Factory Pipelines

In this lab you will configure the Azure environment to allow relational data to be transferred from a SQL Server 2017 database to an Azure SQL Data Warehouse database using Azure Data Factory. The dataset you will use contains data about motor vehicle collisions that happened in New York City from 2012 to 2019. You will use Power BI to visualise collision data loaded from Azure SQL Data Warehouse.

Lab Architecture



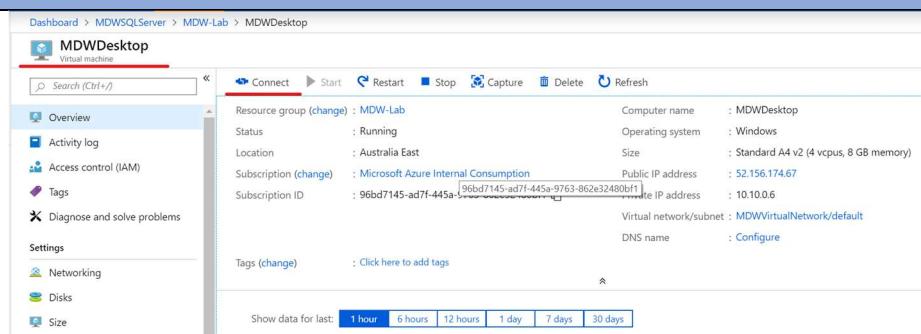
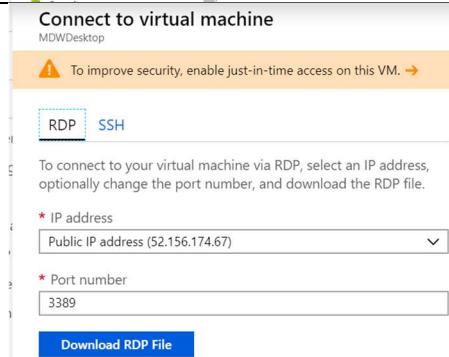
- ① Restore SQL Server backup from Azure Storage and Configure Azure Data Factory Self-Hosted Integration Runtime
- ② Build an Azure Data Factory Pipeline to copy data from a SQL Server table
- ③ Use Azure Storage as a staging area for Polybase
- ④ Load data to an Azure SQL Data Warehouse table using Polybase
- ⑤ Visualize data from Azure SQL Data Warehouse using Power BI

IMPORTANT: Some of the Azure services provisioned by Lab0 require globally unique name and a “-suffix” has been added to their names to ensure this uniqueness. Please take note of the suffix generated as you will need it for the following resources:

Name	Type
mdwcosmosdb-suffix	Cosmos DB account
MDWDataFactory-suffix	Data Factory (V2)
mdwdatalakesuffix	Storage Account
MDWEVENTHUBS-suffix	Event Hubs Namespace
MDWKeyVault-suffix	Key vault
mdwsqvirtualserver-suffix	SQL server
MDWStreamAnalytics-suffix	Stream Analytics job

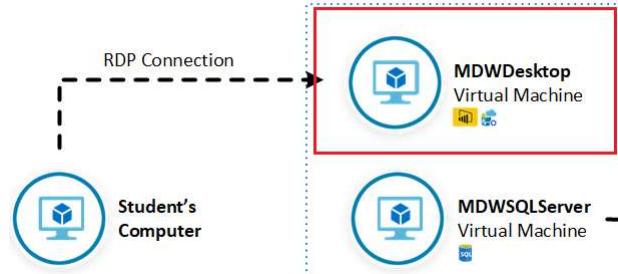
Connect to MDWDesktop

In this section you are going to establish a Remote Desktop Connection to MDWDesktop virtual machine.

IMPORTANT: Execute these steps on your host computer	
<ol style="list-style-type: none">1. In the Azure Portal, navigate to the MDW-Lab resource group and click the MDWDesktop virtual machine.2. On the MDWDesktop blade, from the Overview menu, click the Connect button.	
<ol style="list-style-type: none">3. On the Connect to virtual machine blade, click Download RDP File. This will download a .rdp file that you can use to establish a Remote Desktop Connection with the virtual machine.	

Install required software onto MDWDesktop

In this section you are going to install Power BI Desktop and Azure Data Studio on MDWDesktop.



IMPORTANT: Execute these steps inside the MDWDesktop remote desktop connection

1. Once the file is downloaded, click on file to establish the RDP connection with MDWDesktop
2. Use the following credentials to authenticate:
 - a. User Name: MDWAdmin
 - b. Password: P@ssw0rd123!
3. Once logged in, accept the default privacy settings.
4. Using the browser, download and install the latest version of following software. During the setup, accept all default settings.

Azure Data Studio (User Installer)

<https://docs.microsoft.com/en-us/sql/azure-data-studio/download>

Download and install Azure Data Studio

04/19/2019 • 4 minutes to read • Contributors all

Azure Data Studio runs on Windows, macOS, and Linux.

Download and install the latest release, the *April Release*:

Note

If you're updating from SQL Operations Studio and want to keep your settings, keyboard shortcuts, or code snippets, see [Move user settings](#).

Platform	Download	Release date	Version
Windows	User installer (recommended) System Installer .zip	April 18, 2019	1.6.0

Power BI Desktop (64-bit)

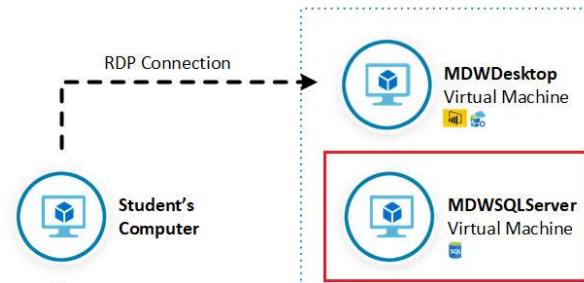
<https://www.microsoft.com/en-us/download/details.aspx?id=45331>

Choose the download you want

File Name	Size
PBIDesktop.msi	186.5 MB
<input checked="" type="checkbox"/> PBIDesktop_x64.msi	205.0 MB

Restore NYCDataSets database onto MDWSQLServer

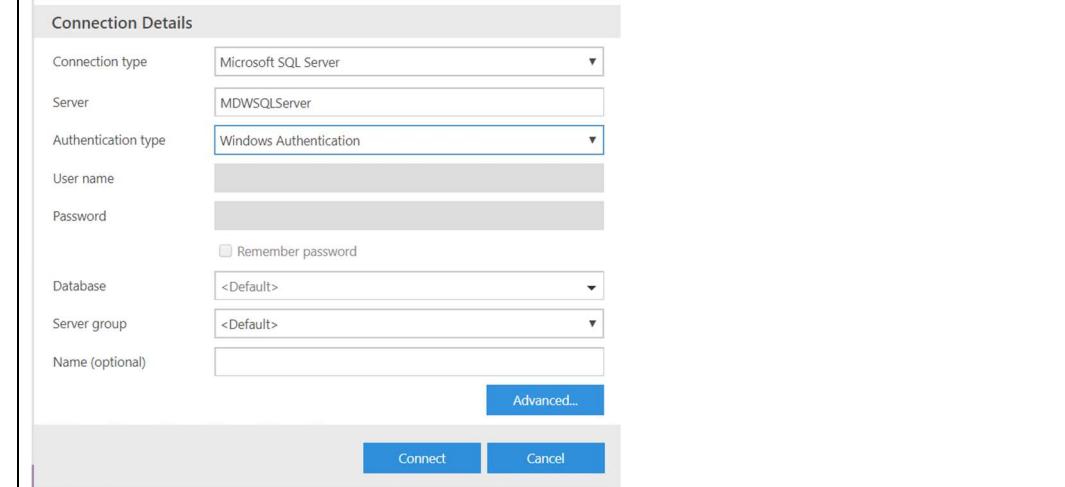
In this section you are going to connect to MDWSQLServer to restore the NYCDataSets database from backup stored in an Azure Storage Account.

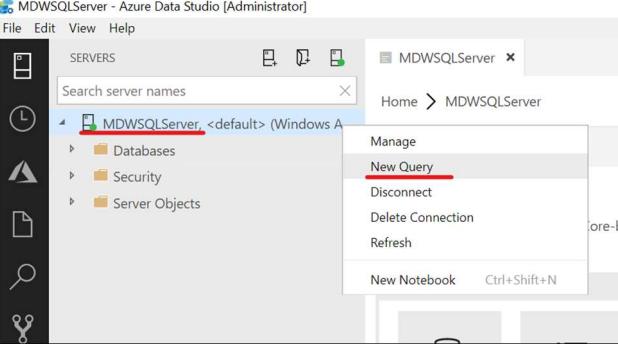


IMPORTANT: The full solution script is saved in the Student's Computer path **C:\ADSIAD\Lab\Lab1\Solution\Restore NYCDataSets.sql**.

IMPORTANT: Execute these steps inside the MDWDesktop remote desktop connection

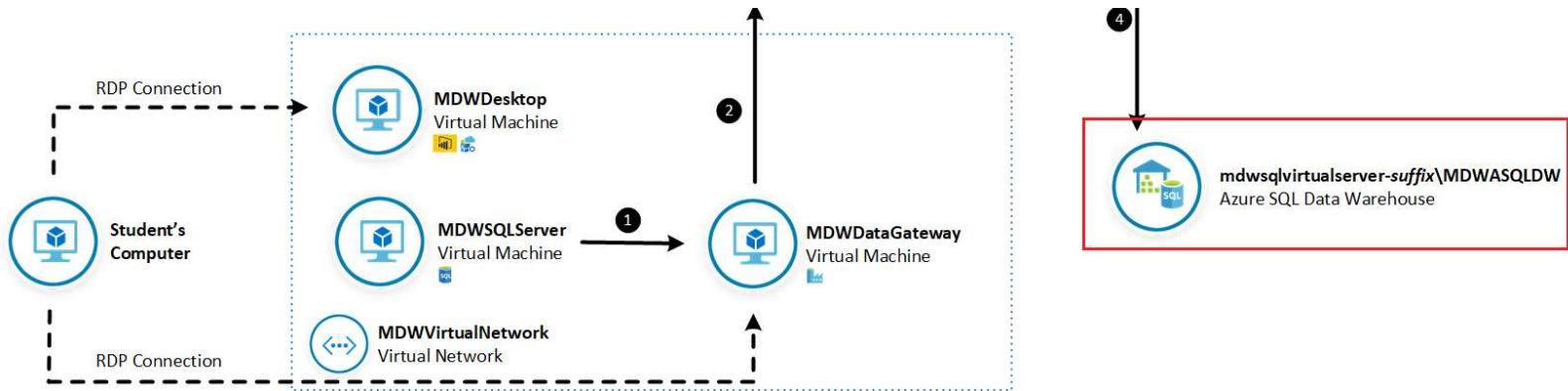
1. Open Azure Data Studio and establish a new connection to **MDWSQLServer** using Windows Authentication



<p>2. Press Ctrl+G to expand the Servers panel 3. Right-click the MDWSQLServer server name on the SERVERS panel and select New Query</p>	
<p>4. On the Query Editor window, create a new credential named [https://mdwresources.blob.core.windows.net/nycdatasets] using a Shared Access Signature (SAS). See SQL command:</p>	<p>SQL command:</p> <pre>create credential [https://mdwresources.blob.core.windows.net/nycdatasets] with identity = 'SHARED ACCESS SIGNATURE', secret = 'sv=2018-03-28&ss=b&srt=sco&sp=rwl&se=2050-12- 30T17:25:52Z&st=2019-04- 05T09:25:52Z&spr=https&sig=4qrD8NmhaSmRFu2gKja67ayohfIDEQH3LdVMa2Utykc%3D ' go</pre>
<p>5. Restore the NYCDataSets database from the backup file stored in the Azure Storage Account. The backup file name is NYCDataSets.Full.bak. The restore command should move the data file to the 'F:\Data' folder and the log file to the 'F:\Log' folder. See SQL command:</p>	<p>SQL command:</p> <pre>restore database NYCDataSets from url = 'https://mdwresources.blob.core.windows.net/nycdatasets/NYCDataSets.Full. bak' with move 'NYCDataSets' to 'F:\Data\NYCDataSets.mdf' , move 'NYCDataSets_log' to 'F:\Log\NYCDataSets_log.ldf' Go</pre>
<p>6. The full solution script is saved in the Student's Computer path C:\ADSIAD\Lab\Lab0\Solution\Restore NYCDataSets.sql.</p>	

Create Azure SQL Data Warehouse database objects

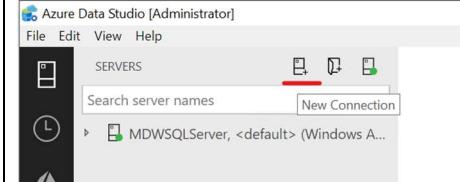
In this section you will connect to Azure SQL Data Warehouse to create the database objects used to host and process data.

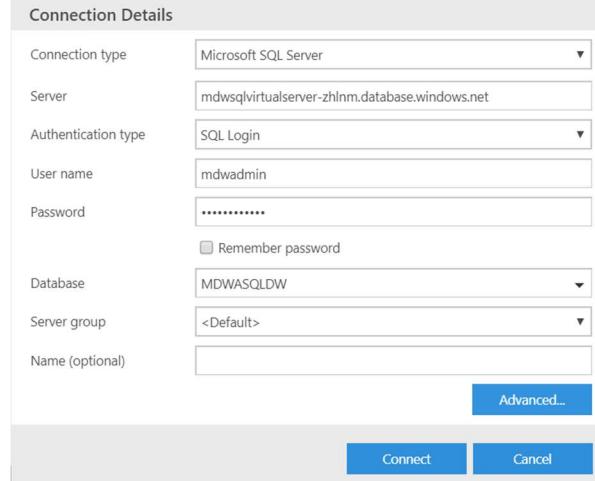
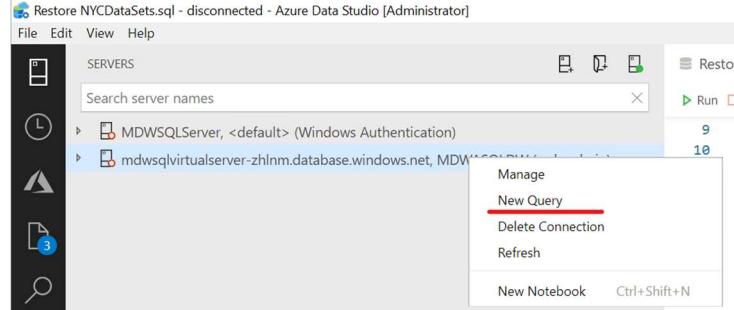


IMPORTANT: The full solution script is saved in the Student's Computer path **C:\ADSIAD\Lab\Lab1\Solution\Create MotorVehicleCollisions Table.sql**.

IMPORTANT: Execute these steps inside the MDWDesktop remote desktop connection

1. Open Azure Data Studio. On the **Servers** panel, click **New Connection**.



<p>2. On the Connection Details panel, enter the following connection details:</p> <ol style="list-style-type: none"> Server: mdwsqlvirtualserver-suffix.database.windows.net Authentication Type: SQL Login User Name: mdwadmin Password: P@ssw0rd123! Database: MDWASQLDW <p>3. Click Connect.</p>	
<p>4. Right-click the server name and click New Query.</p>	
<p>5. On the new query window, create a new database schema named [NYC]. See SQL Command.</p>	<p>SQL Command:</p> <pre>create schema [NYC] go</pre>
<p>6. Create a new round robin distributed table named NYC.NYPD_MotorVehicleCollisions, see column definitions on the SQL Command.</p>	<p>SQL Command:</p> <pre>create table [NYC].[NYPD_MotorVehicleCollisions]([UniqueKey] [int] NULL, [CollisionDate] [date] NULL, [CollisionDayOfWeek] [varchar](9) NULL, [CollisionTime] [time](7) NULL, [CollisionTimeAMPM] [varchar](2) NOT NULL,</pre>

```
[CollisionTimeBin] [varchar](11) NULL,  
[Borough] [varchar](200) NULL,  
[ZipCode] [varchar](20) NULL,  
[Latitude] [float] NULL,  
[Longitude] [float] NULL,  
[Location] [varchar](200) NULL,  
[OnStreetName] [varchar](200) NULL,  
[CrossStreetName] [varchar](200) NULL,  
[OffStreetName] [varchar](200) NULL,  
[NumberPersonsInjured] [int] NULL,  
[NumberPersonsKilled] [int] NULL,  
[IsFatalCollision] [int] NOT NULL,  
[NumberPedestriansInjured] [int] NULL,  
[NumberPedestriansKilled] [int] NULL,  
[NumberCyclistInjured] [int] NULL,  
[NumberCyclistKilled] [int] NULL,  
[NumberMotoristInjured] [int] NULL,  
[NumberMotoristKilled] [int] NULL,  
[ContributingFactorVehicle1] [varchar](200) NULL,  
[ContributingFactorVehicle2] [varchar](200) NULL,  
[ContributingFactorVehicle3] [varchar](200) NULL,  
[ContributingFactorVehicle4] [varchar](200) NULL,  
[ContributingFactorVehicle5] [varchar](200) NULL,  
[VehicleTypeCode1] [varchar](200) NULL,  
[VehicleTypeCode2] [varchar](200) NULL,  
[VehicleTypeCode3] [varchar](200) NULL,  
[VehicleTypeCode4] [varchar](200) NULL,  
[VehicleTypeCode5] [varchar](200) NULL  
)  
with (distribution = round_robin)  
go
```

Configure the Azure Data Factory Self Hosted Integration Runtime

In this section you are going to install and configure required software onto MDWDataGateway.

IMPORTANT: Execute these steps on your host computer

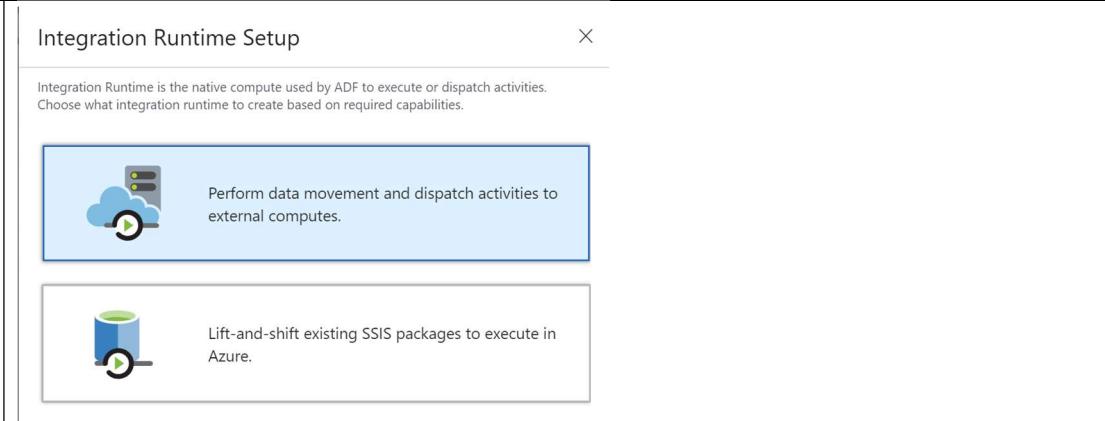
1. In the Azure Portal, navigate to the MDW-Lab resource group and locate the Azure Data Factory **MDWDataFactory-suffix**.
2. On the **MDWDataFactory-suffix** blade, click the **Author & Monitor** button. The Azure Data Factory portal will open on a new browser tab.

The screenshot shows the Azure Data Factory blade for the 'MDWDataFactory-zhlmn' resource. The 'Author & Monitor' button is highlighted with a red box. Other buttons like 'Documentation' and 'Delete' are also visible.

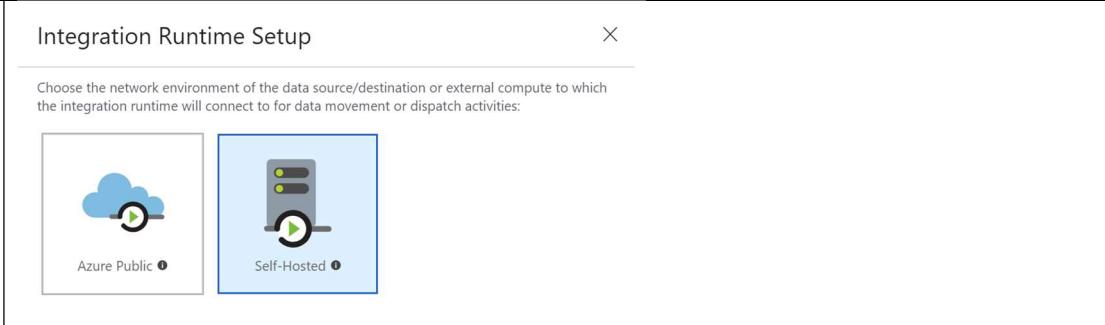
3. On the **Azure Data Factory** portal, click the **Author** (pencil icon) button on the left-hand side menu. On the **Connections** tab, click **Integration Runtimes**.
4. Click the **+ New** button to create a new Integration Runtime.

The screenshot shows the Azure Data Factory portal for the 'MDWDataFactory-zhlmn' resource. The 'Author' button on the left is highlighted with a red box. The 'Connections' tab is selected, and the 'Integration Runtimes' sub-tab is highlighted with a red box. The '+ New' button is also highlighted with a red box. The 'AutoResolveIntegrationRu...' and 'MDWDataGateway' entries are listed under the 'Azure' type.

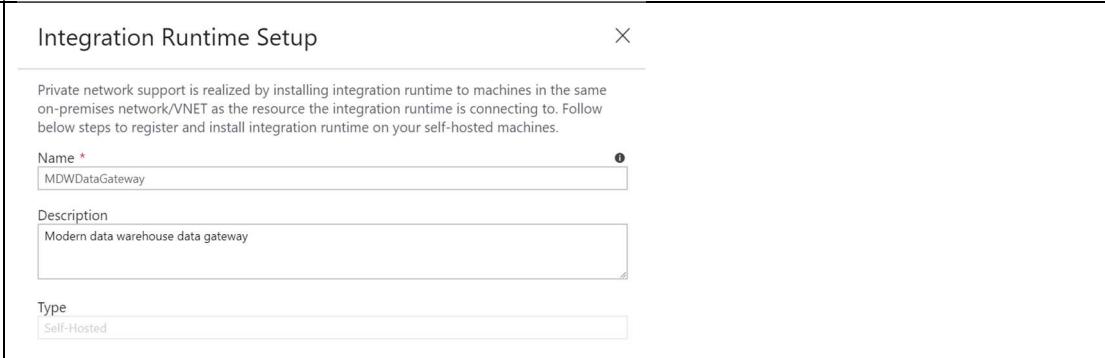
5. On the **Integration Runtime Setup** blade, select **Perform data movement and dispatch activities to external computers** and click **Next**.



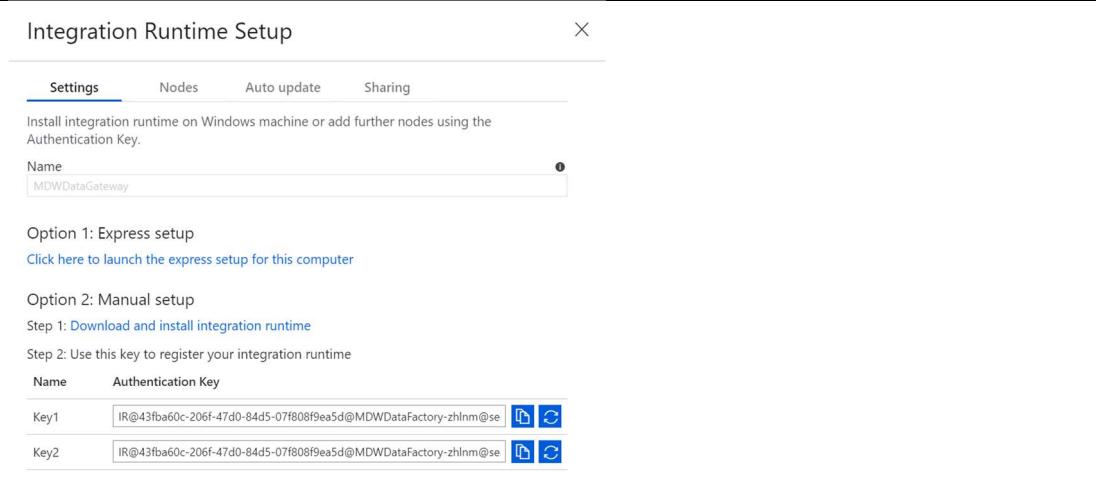
6. When prompted to choose what network environment the integration runtime will connect to, select **Self-Hosted** and click **Next**.



7. Type **MDWDataGateway** in the **Name** text box and give it a meaningful description such as the example here. Click **Next**.

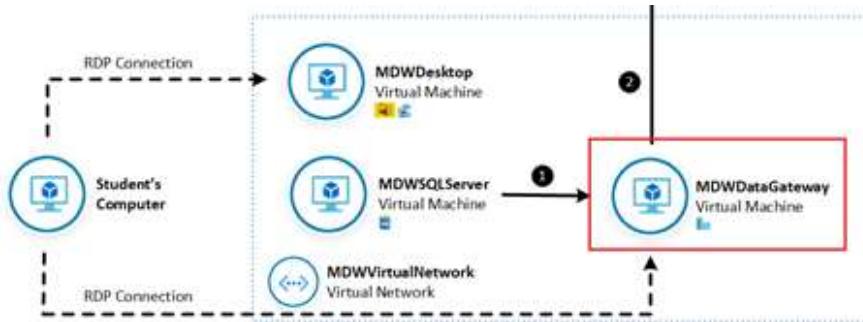


8. Copy any of the generated **Authentication Key** keys (Key 1 or Key 2) to Notepad. You are going to need it in the next step.
9. Click **Finish**.



Connect to MDWDataGateway and register the Self Hosted Integration Runtime with Azure Data Factory

In this section you are going to establish a Remote Desktop Connection to MDWDataGateway virtual machine.



IMPORTANT: Execute these steps on your host computer

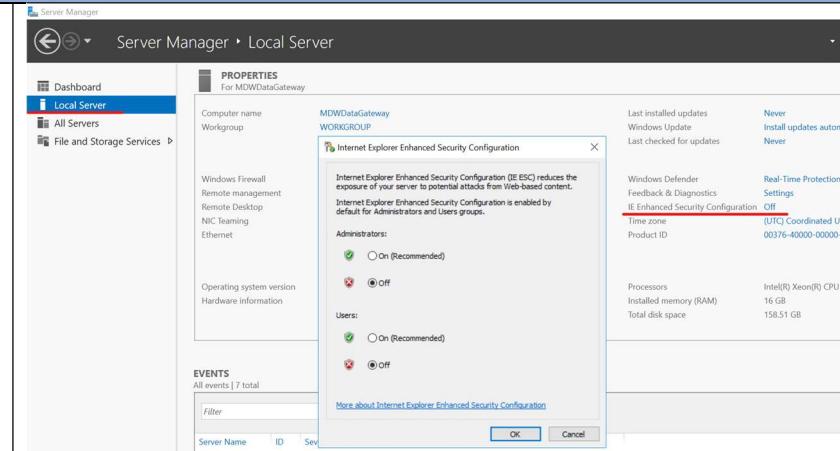
1. On the Azure Portal, navigate to the **MDW-Lab** resource group and locate the **MDWDataGateway** virtual machine.
2. On the **MDWDataGateway** blade, from the **Overview** menu, click the **Connect** button.
3. On the **Connect to virtual machine** blade, click **Download RDP File**. This will download a .rdp file that you can use to establish a Remote Desktop Connection with the virtual machine.

	<p>Dashboard > MDW-Lab > MDWDataGateway</p> <p>MDWDataGateway Virtual machine</p> <p>Search (Ctrl+F)</p> <p>Overview</p> <p>Activity log</p> <p>Access control (IAM)</p> <p>Tags</p> <p>Diagnose and solve problems</p> <p>Settings</p> <p>Resource group (change) : MDW-Lab</p> <p>Status : Running</p> <p>Location : Australia East</p> <p>Subscription (change) : Microsoft Azure Internal Consumption</p> <p>Subscription ID : 96bd7145-ad7f-445a-9763-862e32480bf1</p> <p>Computer name : MDWDataGateway</p> <p>Operating system : Windows</p> <p>Size : Standard D4s v3 (4 vcpus, 16 GB memory)</p> <p>Public IP address : 13.72.224.119</p> <p>Private IP address : 10.10.0.5</p> <p>Virtual network/subnet : MDWVirtualNetwork/default</p> <p>DNS name : Configure</p>
	<p>Connect to virtual machine</p> <p>MDWDataGateway</p> <p>To improve security, enable just-in-time access on this VM.</p> <p>RDP SSH</p> <p>* IP address Public IP address (13.72.224.119)</p> <p>* Port number 3389</p> <p>Download RDP File</p>

4. Once the file is downloaded, click on file to establish the RDP connection with **MDWDataGateway**
5. User the following credentials to authenticate:
 - a. **User Name:** MDWAdmin
 - b. **Password:** P@ssw0rd123!

IMPORTANT: Execute these steps inside the MDWDataGateway remote desktop connection

1. Once logged in, on the **Server Manager**, select **Local Server** on the left-hand side menu. On the right-hand side panel, locate the **IE Enhanced Security Configuration** and click the **On** link.
2. Turn the setting **Off** for both Administrators and Users.
3. Close **Server Manager**.

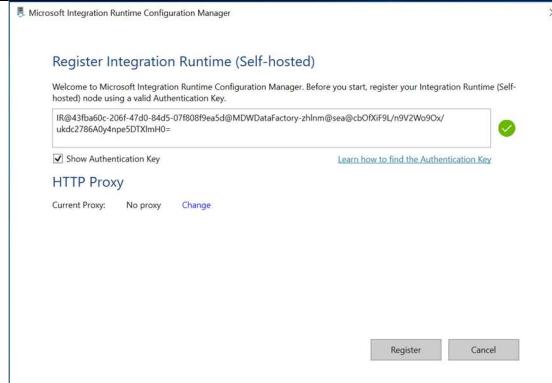


4. Open the browser and download and execute the latest version of the Azure Data Factory Integration Runtime.

Azure Data Factory Integration Runtime

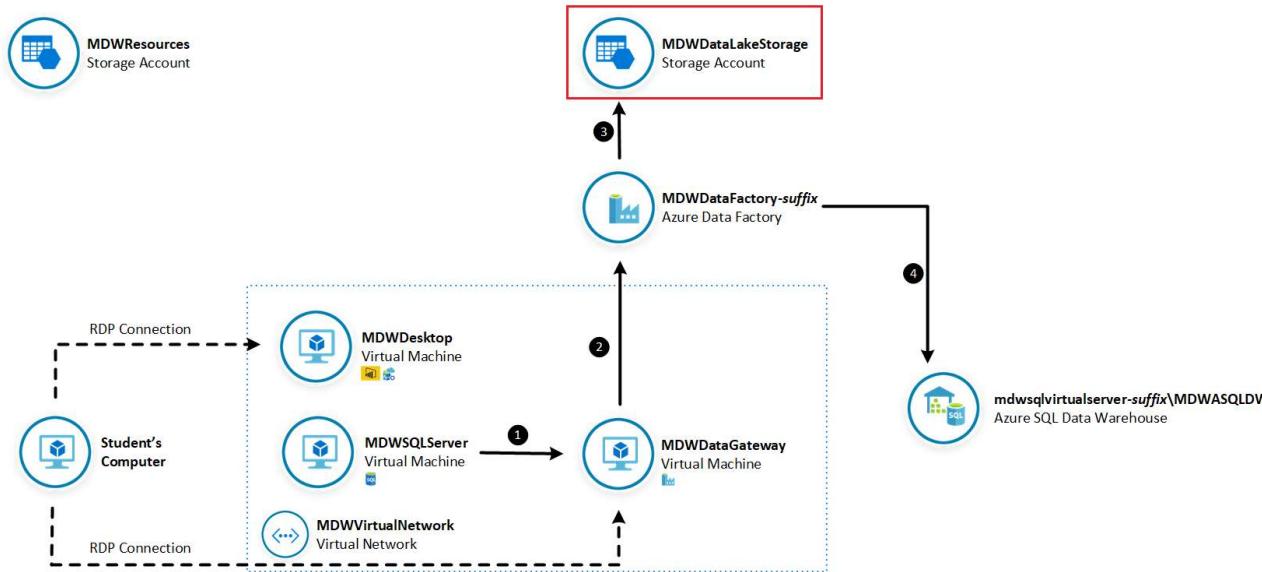
<https://www.microsoft.com/en-ie/download/details.aspx?id=39717>

5. Accept all default options during the setup wizard.
Once the setup if finished, the **Microsoft Integration Runtime Configuration Manager** will pop up asking you to enter a valid authentication key.
6. Enter the authentication key generated in the previous exercise and click **Register**.
7. Once registration is confirmed, click **Finish**.



Create Staging Container on Azure Blob Storage

In this section you create a staging container in your MDWDataLake that will be used as a staging environment for Polybase before data can be copied to Azure SQL Data Warehouse.



IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, go to the lab resource group and locate the Azure Storage account **mdwdatalakesuffix**.
2. On the **Overview** panel, click **Blobs**.

The screenshot shows the Azure Storage account 'mdwdatalakezhlnm' in the MDW-Lab resource group. The 'Overview' tab is selected. The 'Services' section highlights the 'Blobs' service, which is described as 'REST-based object storage for unstructured data'. Other services listed include 'Activity log', 'Access control (IAM)', 'Tags', 'Diagnose and solve problems', 'Events', 'Storage Explorer (preview)', 'Access keys', 'Geo-replication', and 'CORS'.

3. On the **mdwdatalakesuffix - Blobs** blade, click **+ Container**.

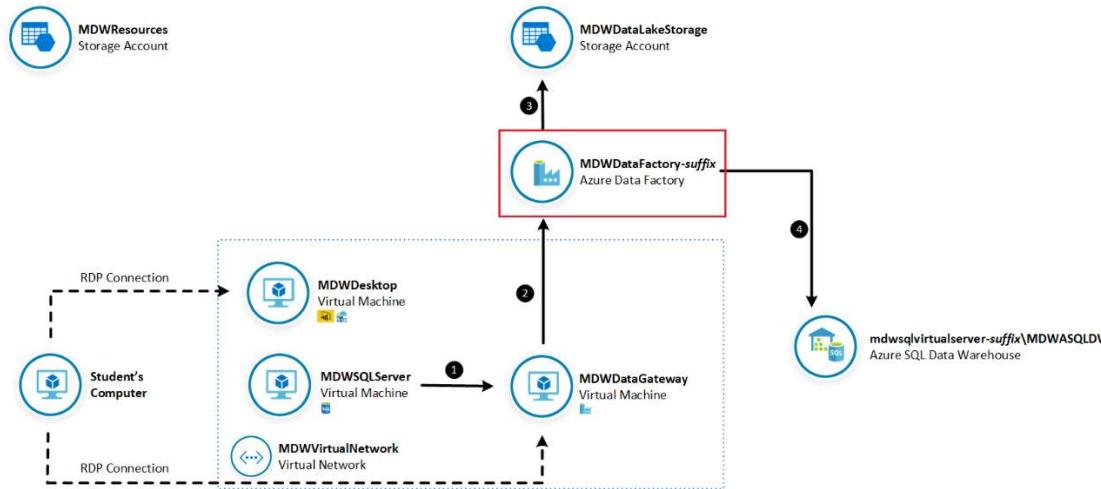
The screenshot shows the 'Blobs' blade for the 'mdwdatalakezhlnm' storage account. A new container named 'polybase' has been created and is listed in the 'NAME' column. The 'Container' button is highlighted in red at the top of the blade.

4. On the **New container** blade, enter the following details:
 - a. **Name:** polybase
 - b. **Public access level:** Private (no anonymous access)
5. Click **OK** to create the new container.

The screenshot shows the 'New container' blade. At the top, there are buttons for Container, Refresh, Delete, and Change access level. Below that, the title 'New container' is displayed. There are two main input fields: one for 'Name' containing 'polybase' and another for 'Public access level' set to 'Private (no anonymous access)'. At the bottom of the blade are two buttons: 'OK' and 'Cancel'.

Create Azure Data Factory Pipeline to Copy Relational Data

In this section you will build an Azure Data Factory pipeline to copy a table from MDWSQLServer to Azure SQL Data Warehouse.



Create Linked Service connections

IMPORTANT: Execute these steps on your host computer

1. Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel (pencil icon). Under **Connections** tab, click **Linked Services** and then click **+ New** to create a new linked service connection.

Microsoft Azure | Data Factory > MDWDataFactory-zhlnm

Search resource

Factory Resources

Pipelines

Datasets

Data Flows (Preview)

Connections

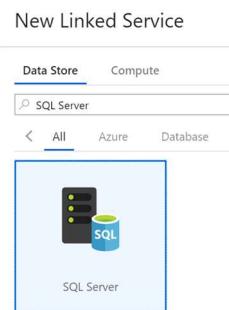
Linked Services

New

Name

Action

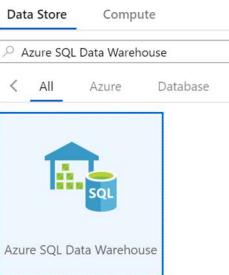
2. On the **New Linked Service** blade, type “SQL Server” in the search box to find the **SQL Server** linked service. Click **Continue**.



3. On the **New Linked Service (SQL Server)** blade, enter the following details:
- Name:** MDWSQLServer_NYCDatasets
 - Connect via integration runtime:** MDWDataGateway
 - Server Name:** MDWSQLServer
 - DatabaseName:** NYCDatasets
 - Authentication Type:** Windows Authentication
 - User Name:** MDWAdmin
 - Password:** P@ssw0rd123!
4. Click **Test connection** to make sure you entered the correct connection details and then click **Finish**.

5. Repeat the process to create an Azure SQL Data Warehouse linked service connection.

New Linked Service



6. On the **New Linked Service (Azure SQL Data Warehouse)** blade, enter the following details:

- Name:** MDWVirtualSQLServer_MDWASQLDW
- Connect via integration runtime:** AutoResolveIntegrationRuntime
- Account selection method:** From Azure subscription
- Azure subscription:** <your subscription>
- Server Name:** mdwsqlvirtualserver-suffix
- DatabaseName:** MDWASQLDW
- Authentication Type:** SQL Authentication
- User Name:** MDWAdmin
- Password:** P@ssw0rd123!

7. Click **Test connection** to make sure you entered the correct connection details and then click **Finish**.

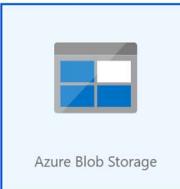
8. Repeat the process once again to create an Azure Storage linked service connection.

New Linked Service

Data Store Compute

Azure Blob Storage

All Azure Database



Azure Blob Storage

9. On the **New Linked Service (Azure Blob Storage)** blade, enter the following details:
- j. **Name:** MDWDataLake
 - k. **Connect via integration runtime:** AutoResolveIntegrationRuntime
 - l. **Authentication method:** Account key
 - m. **Account selection method:** From Azure subscription
 - n. **Azure subscription:** <your subscription>
 - o. **Storage account name:** mdwdatalakesuffix
10. Click **Test connection** to make sure you entered the correct connection details and then click **Finish**.

New Linked Service (Azure Blob Storage)

Name *
MDWDataLake

Description

Connect via integration runtime *
AutoResolveIntegrationRuntime

Authentication method
Account key

Connection String Azure Key Vault

Account selection method
 From Azure subscription Enter manually

Azure subscription
Microsoft Azure Internal Consumption (96bd7145-ad7f-445a-9763-862e32480bf1)

Storage account name *
mdwdatalakezhlnm

Additional connection properties
+ New

Annotations
+ New

Advanced

Cancel Test connection Finish

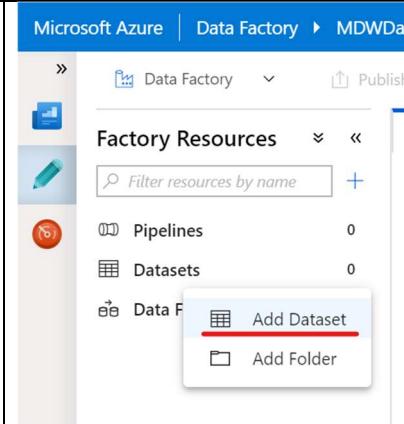
11. You should now see 3 linked services connections that will be used as source, destination and staging.

Connections X		
Linked Services		Integration Runtimes
+ New	Name	Actions
	MDWDataLake	
	MDWSQLServer_NYCDatasets	
	MDWSQLVirtualServer_MDWASQLDW	

Create Source and Destination Data Sets

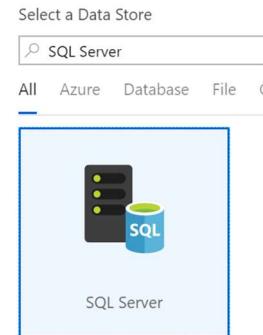
IMPORTANT: Execute these steps on your host computer

1. Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel (pencil icon). Under **Factory Resources** tab, click the ellipsis (...) next to **Datasets** and then click **Add Dataset** to create a new dataset.



2. Type SQL Server in the search box and select **SQL Server**. Click **Finish**.

New Dataset



3. On the **New Data Set** tab, enter the following details:
- Name:** NYCDataSets_MotorVehicleCollisions
 - Linked Service:** MDWSQLServer_NYCDatasets
 - Table:** [NYC].[NYPD_MotorVehicleCollisions]
4. Leave remaining fields with default values and click Continue

← Set Properties

X

Name
NYCDataSets_MotorVehicleCollisions

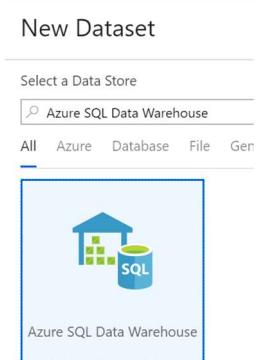
Linked service *
MDWSQLServer_NYCDatasets

Edit Connection

Table
[NYC].[NYPD_MotorVehicleCollisions]

Import schema
 From connection/store None

5. Repeat the process to create a new Azure SQL Data Warehouse data set.



6. On the **New Data Set** tab, enter the following details:
- Name:** MDWASQLDW_MotorVehicleCollisions
 - Linked Service:** MDWSQLVirtualServer_MDWASQLDW
 - Table:** [NYC].[NYPD_MotorVehicleCollisions]
7. Leave remaining fields with default values and click Continue.

← Set Properties

Name: MDWASQLDW_MotorVehicleCollisions

Linked service *: MDWVirtualSQLServer_MDWASQLDW

[Edit Connection](#)

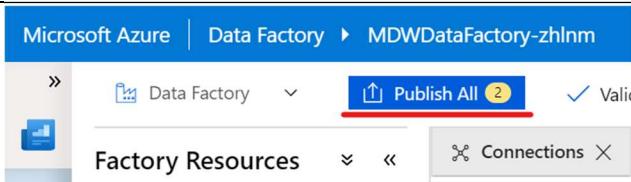
Table: [NYC].[NYPD_MotorVehicleCollisions]

Edit

Import schema

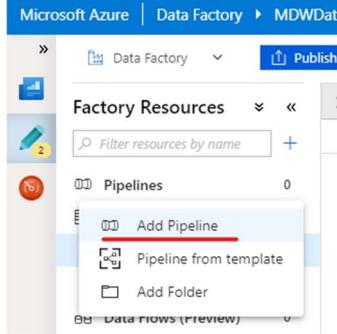
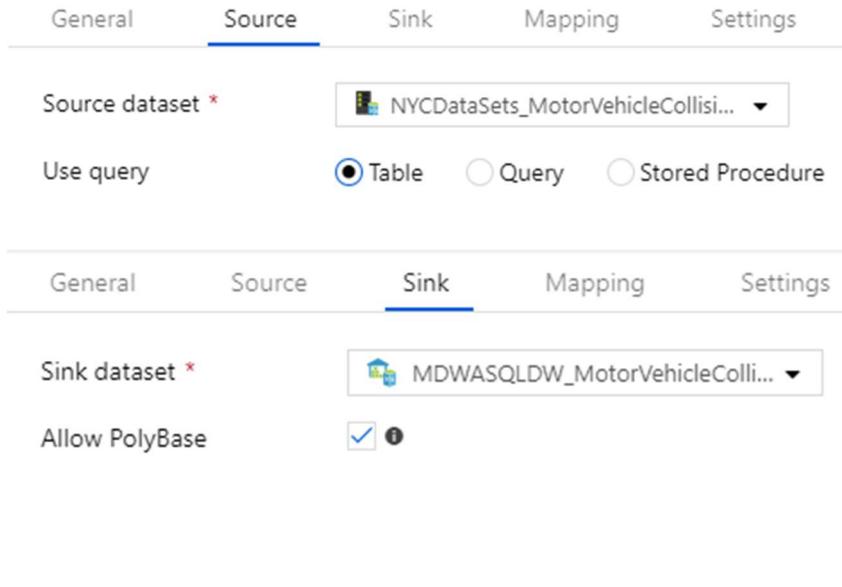
From connection/store None

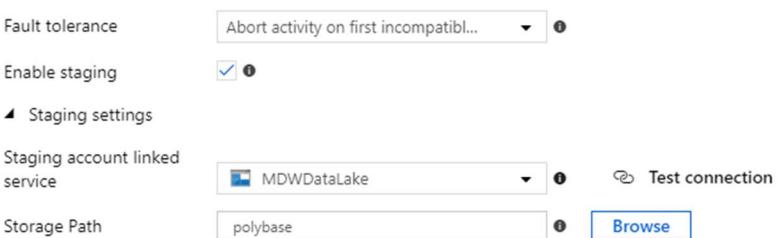
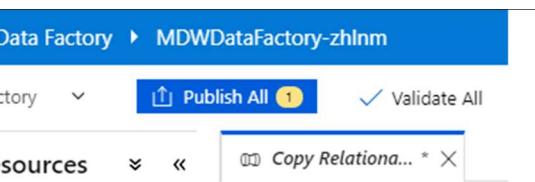
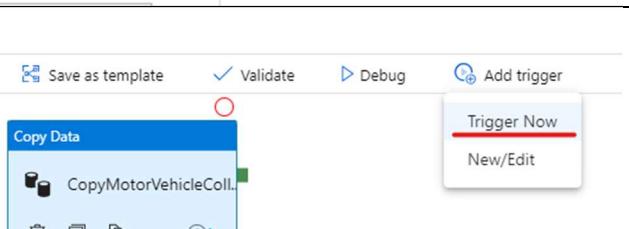
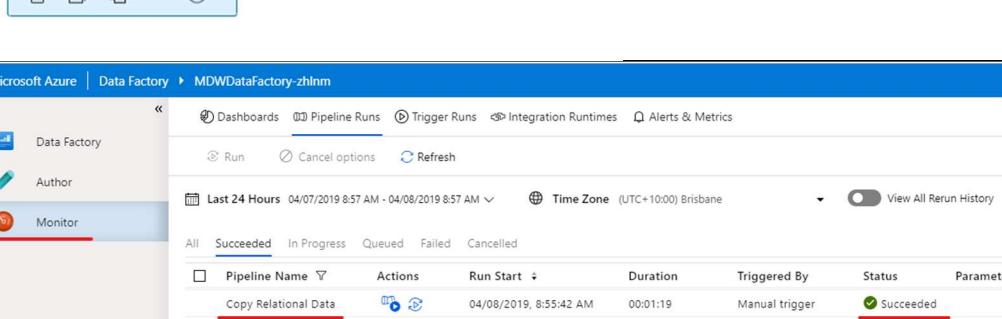
8. Publish your dataset changes by clicking the **Publish All** button on the top of the screen.



Create and Execute Pipeline

IMPORTANT: Execute these steps on your host computer

<ol style="list-style-type: none"> 1. Open the Azure Data Factory portal and click the Author option on the left-hand side panel (pencil icon). Under Factory Resources tab, click the ellipsis (...) next to Pipelines and then click Add Pipeline to create a new dataset. 2. On the New Pipeline tab, enter the following details: <ol style="list-style-type: none"> a. General > Name: Copy Relational Data 3. Leave remaining fields with default values. 	
<ol style="list-style-type: none"> 4. From the Activities panel, type “Copy Data” in the search box. Drag the Copy Data activity on to the design surface. 5. Select the Copy Data activity and enter the following details: <ol style="list-style-type: none"> a. General > Name: CopyMotorVehicleCollisions b. Source > Source dataset: NYCDataSets_MotorVehicleCollisions c. Sink > Sink dataset: MDWASQLDW_MotorVehicleCollisions d. Sink > Allow PolyBase: Checked e. Settings > Enable staging: Checked f. Settings > Staging account linked service: MDWDataLake g. Settings > Storage Path: polybase 6. Leave remaining fields with default values. 	

	
7. Publish your pipeline changes by clicking the Publish all button.	
8. To execute the pipeline, click on Add trigger menu and then Trigger Now . 9. On the Pipeline Run blade, click Finish .	
10. To monitor the execution of your pipeline, click on the Monitor menu on the left-hand side panel. 11. You should be able to see the Status of your pipeline execution on the right-hand side panel.	

Visualize Data with Power BI

In this section you are going to use Power BI to visualize data from Azure SQL Data Warehouse. The Power BI report will use an Import connection to query Azure SQL Data Warehouse and visualise Motor Vehicle Collision data from the table you loaded in the previous exercise.

IMPORTANT: Execute these steps inside the MDWDesktop remote desktop connection

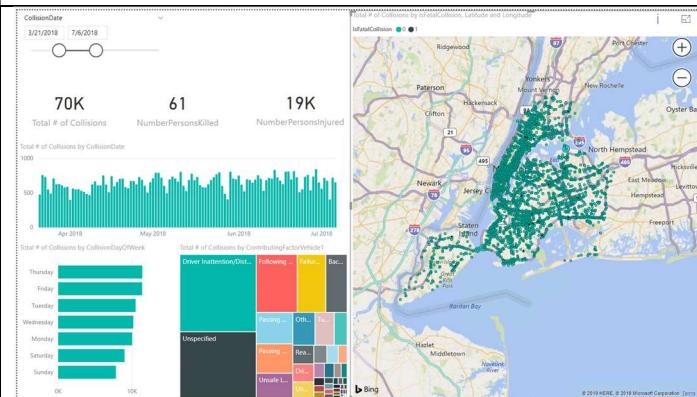
1. On MDWDesktop, download the Power BI report from the link <https://aka.ms/MDWLab1> and save it on the Desktop.
2. Open the file MDWLab1.pbix with Power BI Desktop. Optionally sign up for the Power BI tips and tricks email, or to dismiss this, click to sign in with an existing account, and then hit the escape key.
3. When prompted to enter the value of the **MDWSQLVirtualServer** parameter, type the full server name: **mdwsqvirtualserver-suffix.database.windows.net**
4. Click **Load**, and then **Run** to acknowledge the **Native Database Query message**
5. When prompted, enter the credentials to connect to the database – Username is **MDWAdmin**, Password is **P@ssw0rd123!**
6. Once the data is finished loading, interact with the report by changing the CollisionDate slicer and by clicking on the other visualisations.
7. Save your work and close Power BI Desktop.

MDWLab1

Modern Data Warehouse Power BI Report Template.

MDWSQLVirtualServer

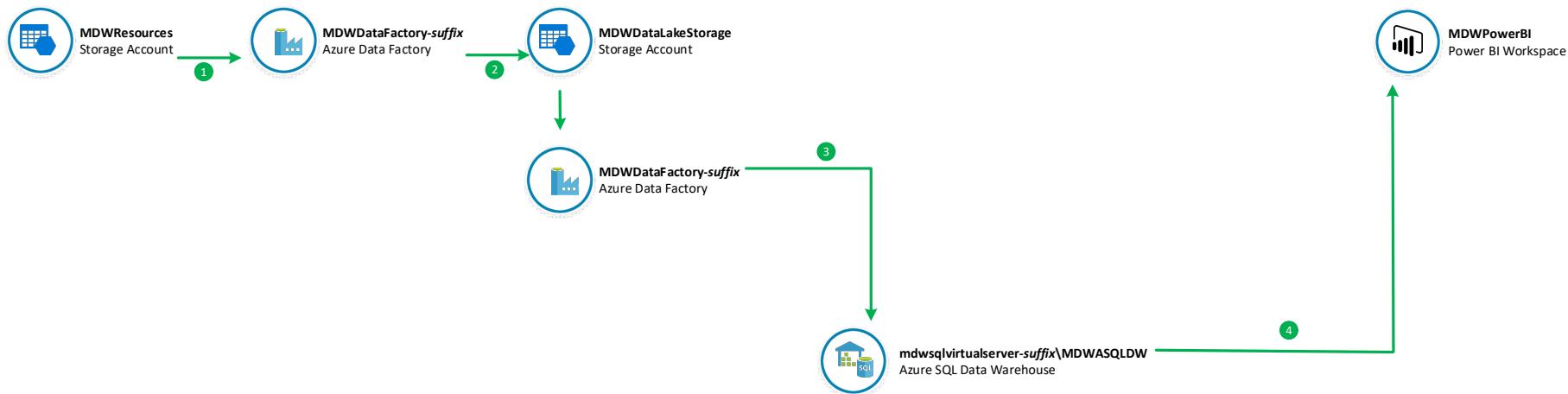
mdwsqvirtualserver-zhlm.database.windows.net



Lab 2: Transform Big Data using Azure Data Factory and Azure SQL Data Warehouse

In this lab you will use Azure Data Factory to download large data files to your data lake and use Azure SQL Data Warehouse to generate a summary dataset and store it. The dataset you will use contains detailed New York City Yellow Taxi rides for 2018. You will generate a daily aggregated summary of all rides and save the result in your data warehouse. You will use Power BI to visualise summarised taxi ride data.

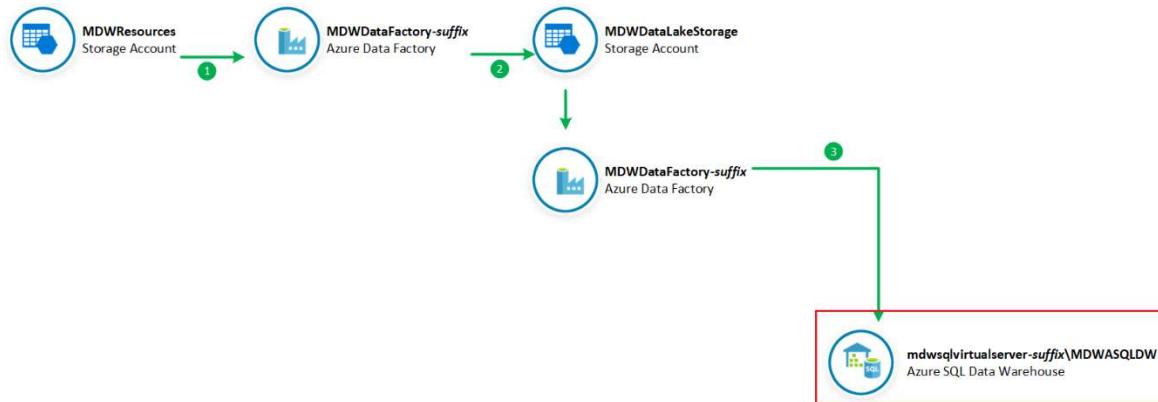
Lab Architecture



- 1 Build an Azure Data Factory Pipeline to copy big data files from shared Azure Storage
- 2 Save data files to your data lake
- 3 Use Polybase to load data into staging tables in your Azure SQL Data Warehouse.
Call a Stored Procedure to perform data aggregations and save results in the final table.
- 4 Visualize data from your Azure SQL Data Warehouse using Power BI

Create Azure SQL Data Warehouse database objects

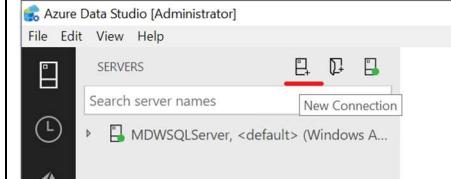
In this section you will connect to Azure SQL Data Warehouse to create the database objects used to host and process data.

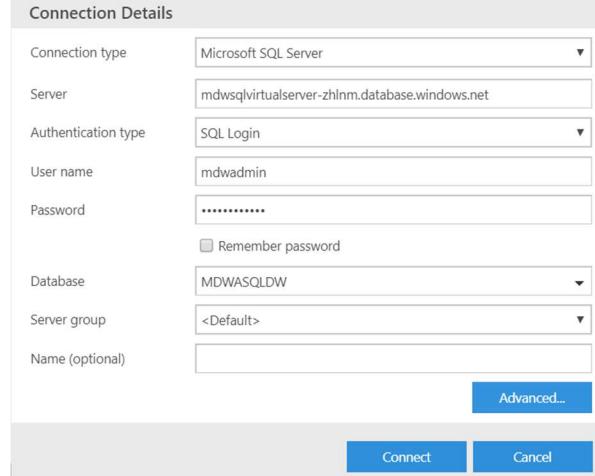


IMPORTANT: The full solution script is saved in the Student's Computer path **C:\ADSIAD\Lab\Lab2\Solution\Create Staging NYCTaxiData.sql**.

IMPORTANT: Execute these steps inside the MDWDesktop remote desktop connection

1. Open Azure Data Studio.
2. If you already have a connection to MDWSQLVirtualServer, then go to step 6.
3. On the **Servers** panel, click **New Connection**.



<p>4. On the Connection Details panel, enter the following connection details:</p> <ol style="list-style-type: none"> Server: mdwsqlvirtualserver-suffix.database.windows.net Authentication Type: SQL Login User Name: mdwadmin Password: P@ssw0rd123! Database: MDWASQLDW <p>5. Click Connect.</p>	
<p>6. Right click the MDWSQLVirtualServer name and then click New Query.</p>	
<p>7. Create a new round robin distributed table named NYC.TaxiDataSummary – execute the script shown to the right</p>	<p>SQL Command</p> <p>IMPORTANT: The full solution script is saved in the Student's Computer path C:\ADSIAD\Lab\Lab2\Solution\Create Staging NYCTaxiData.sql.</p> <pre>create schema [Staging] go create table [Staging].[NYCTaxiData] ([VendorID] [int] NULL, [tpep_pickup_datetime] [datetime] NULL, [tpep_dropoff_datetime] [datetime] NULL, [passenger_count] [smallint] NULL, [trip_distance] [decimal](8, 2) NULL, [RatecodeID] [smallint] NULL, [store_and_fwd_flag] [char](1) NULL, [PULocationID] [int] NULL,</pre>

	<pre>[DOLocationID] [int] NULL, [payment_type] [smallint] NULL, [fare_amount] [decimal](10, 2) NULL, [extra] [decimal](10, 2) NULL, [mta_tax] [decimal](10, 2) NULL, [tip_amount] [decimal](10, 2) NULL, [tolls_amount] [decimal](10, 2) NULL, [improvement_surcharge] [decimal](10, 2) NULL, [total_amount] [decimal](10, 2) NULL) with (distribution = round_robin, heap) go</pre>
8. Create a new Staging table called NYCTaxiLocationLookup to host NYC taxi location lookup data – execute the script shown to the right	<p>SQL Command:</p> <pre>create table [Staging].[NYCTaxiLocationLookup] ([LocationID] [int] NULL, [Borough] [varchar](200) NULL, [Zone] [varchar](200) NULL, [service_zone] [varchar](200) NULL) with (distribution = round_robin, clustered columnstore index)</pre>
9. Create a stored procedure that will transform and aggregate data and save the result in a new table	<p>SQL Command:</p> <pre>create procedure Staging.spNYCLoadTaxiDataSummary</pre>

called NYC.TaxiDataSummary – execute the script shown to the right

```
as
    --Drop and re-create Staging.idx_NYCTaxiData index
    if (exists(select top 1 1 from sys.indexes where name =
'idx_NYCTaxiData' and object_id = object_id('Staging.NYCTaxiData')))
        drop index idx_NYCTaxiData on Staging.NYCTaxiData

    create index idx_NYCTaxiData on
        Staging.NYCTaxiData(tpipe_pickup_datetime, PULocationID, payment_type,
        passenger_count, trip_distance, tip_amount, fare_amount,
        total_amount)

    --Drop and re-create NYC.TaxiDataSummary table
    if (exists(select top 1 1 from sys.objects where name =
'TaxiDataSummary' and schema_id = schema_id('NYC') and type = 'U'))
        drop table NYC.TaxiDataSummary

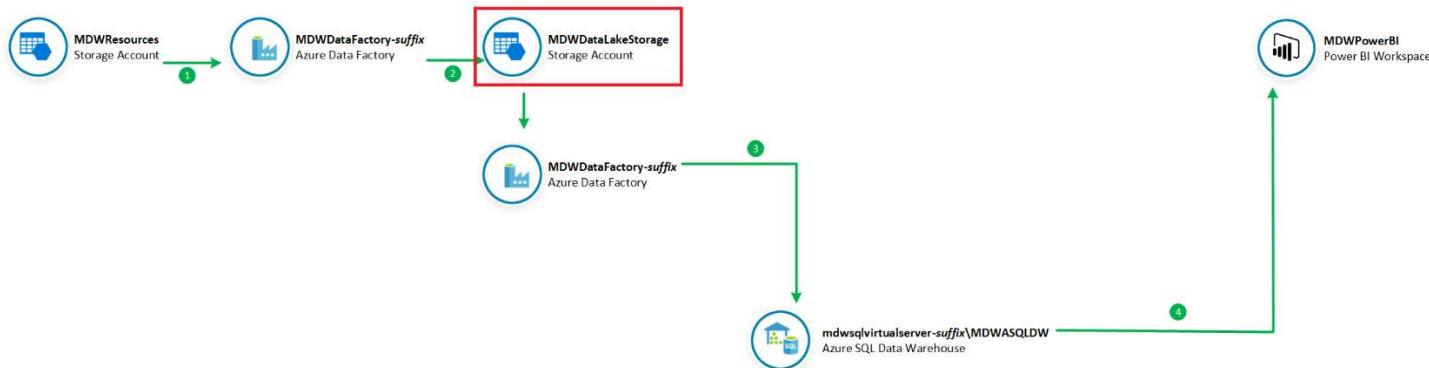
    create table NYC.TaxiDataSummary
        with (
            distribution = round_robin
        )
        as
        select
            cast(tpipe_pickup_datetime as date) as PickUpDate
            , PickUp.Borough as PickUpBorough
            , PickUp.Zone as PickUpZone
            , case payment_type
                when 1 then 'Credit card'
                when 2 then 'Cash'
                when 3 then 'No charge'
                when 4 then 'Dispute'
                when 5 then 'Unknown'
                when 6 then 'Voided trip'
```

```
        end as PaymentType
        , count(*) as TotalTripCount
        , sum(passenger_count) as TotalPassengerCount
        , sum(trip_distance) as TotalDistanceTravelled
        , sum(tip_amount) as TotalTipAmount
        , sum(fare_amount) as TotalFareAmount
        , sum(total_amount) as TotalTripAmount
from Staging.NYCTaxiData
    inner join Staging.NYCTaxiLocationLookup as PickUp
        on NYCTaxiData.PULocationID = PickUp.LocationID
group by cast(tpep_pickup_datetime as date)
        , PickUp.Borough
        , PickUp.Zone
        , payment_type

--drop index idx_NYCTaxiData so it does not impact future loads
drop index idx_NYCTaxiData on Staging.NYCTaxiData
go
```

Create NYCTaxiData Container on Azure Blob Storage

In this section you create a container in your MDWDataLake that will be used as a repository for the NYC Taxi Data files. You will copy 12 files from the MDWResources Storage Account into your NYCTaxiData container. These files contain data for all Yellow Taxi rides in 2018, one file for each month of the year.



IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, go to the lab resource group and locate the Azure Storage account **mdwdatalakesuffix**.
2. On the **Overview** panel, click **Blobs**.

The screenshot shows the Azure Storage Account Overview page for **mdwdatalakezhlnm**. The page includes the following sections:

- Overview** (highlighted in red): Includes links for **Activity log**, **Access control (IAM)**, **Tags**, **Diagnose and solve problems**, **Events**, and **Storage Explorer (preview)**.
- Resource group (change) : MDW-Lab**
- Status**: Primary; Available, Secondary: Avail
- Location**: Australia East, Australia Southeast
- Subscription (change)**: Microsoft Azure Internal Consumption
- Subscription ID**: 96bd7145-ad7f-445a-9763-862e32
- Tags (change)**: Click here to add tags
- Services** section:
 - Blobs** (highlighted in red): REST-based object storage for unstructured data. A link to **Learn more** is provided.

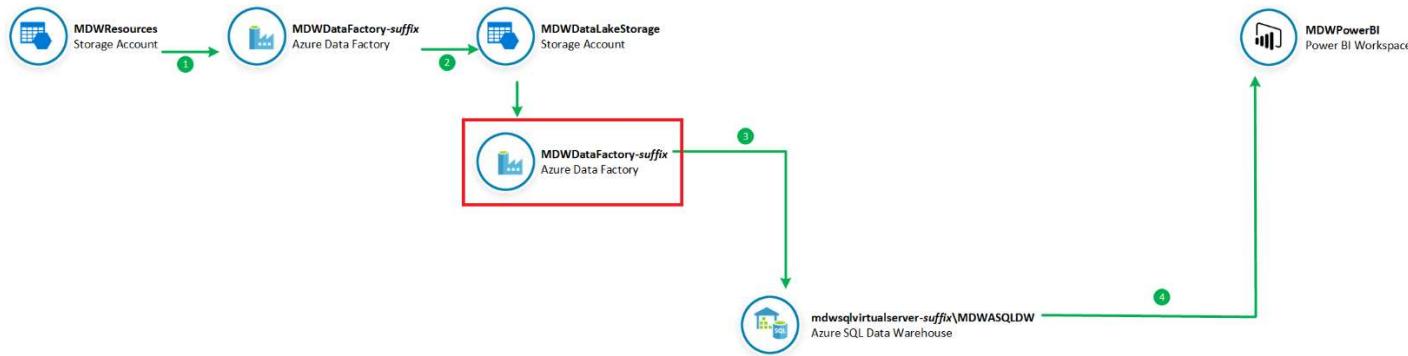
3. On the **mdwdatalakesuffix – Blobs** blade, click **+ Container**.

The screenshot shows the 'Blobs' blade for the storage account 'mdwdatalakezhlnm'. At the top right, there is a red box highlighting the '+ Container' button. Below it, the storage account name 'mdwdatalakezhlnm' is displayed. To the right, there is a search bar labeled 'Search containers by prefix' and a 'NAME' field containing 'polybase'. The left sidebar includes links for 'Overview', 'Activity log', 'Access control (IAM)', 'Tags', and 'Diagnose and solve problems'.

4. On the **New container** blade, enter the following details:
- Name:** nyctaxidata
 - Public access level:** Private (no anonymous access)
5. Click **OK** to create the new container.

The screenshot shows the 'New container' blade. It has a 'Name' field containing 'nyctaxidata' with a green checkmark indicating validation. Below it is a 'Public access level' dropdown set to 'Private (no anonymous access)'. At the bottom are 'OK' and 'Cancel' buttons.

Create Linked Service connection to MDWResources



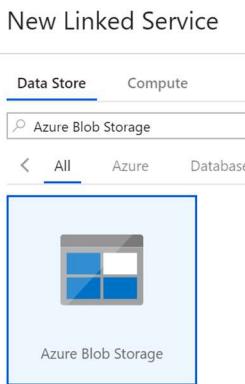
IMPORTANT: Execute these steps on your host computer

1. Open the **Azure Data Factory** portal and click the **Author** option (pencil icon) on the left-hand side panel. Under **Connections** tab, click **Linked Services** and then click **+ New** to create a new linked service connection.

Screenshot of the Microsoft Azure Data Factory portal, specifically the 'Connections' tab. The interface shows:

- Left sidebar: Factory Resources (Pipelines: 0, Datasets: 0, Data Flows (Preview): 0).
- Top navigation: Microsoft Azure | Data Factory > MDWDataFactory-zhlnm.
- Top right: Search resource, Publish All, Validate All, Refresh, Discard A.
- Connections tab: Sub-tabs include 'Connections' (selected), 'Integration Runtimes', and 'New'. A red box highlights the '+ New' button under the 'Connections' tab.

- On the **New Linked Service** blade, type “Azure Blob Storage” in the search box to find the **Azure Blob Storage** linked service. Click **Continue**.



- On the **New Linked Service (Azure Blob Storage)** blade, enter the following details:
 - Name:** MDWResources
 - Connect via integration runtime:** AutoResolveIntegrationRuntime
 - Authentication method:** SAS URI
 - SAS URL:**
<https://mdwresources.blob.core.windows.net/?sv=2018-03-28&ss=b&srt=sco&sp=rwl&se=2050-12-30T17:25:52Z&st=2019-04-05T09:25:52Z&spr=https&sig=4qrD8NmhaSmRFu2gKja67ayohfIDEQH3LdVMa2Utykc%3D>
- Alternatively copy SAS URL from **C:\ADSIAD\Lab\Lab2\MDWResources SAS URL.txt**
- Click **Test connection** to make sure you entered the correct connection details and then click **Finish**.

← New Linked Service (Azure Blob Storage) X

Name *	MDWResources
Description	
Connect via integration runtime *	AutoResolveIntegrationRuntime
Authentication method	SAS URI
SAS URI	Azure Key Vault
SAS URL *	https://mdwresources.blob.core.windows.net/?sv=2018-03-28&ss=b&srt=sco&sp=rwl&se=2050-12-30T17:25:52Z&st=2019-04-05T09:25:52Z&spr=https&sig=4qrD8NmhaSmRFu2gKja67ayohfIDEQH3LdVMa2Utykc%3D
SAS Token	Azure Key Vault
SAS Token	<small>sample: ?sv=<storage services version>&st=<start time>&se=<expire time>&sr=<resource>&sp=<per</small>
Annotations	<ul style="list-style-type: none"> + New
Advanced	•

Cancel Connection successful Test connection Finish

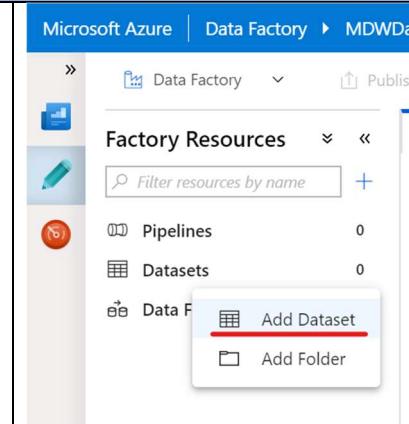
Create Source and Destination Data Sets

In this section you are going to create 5 datasets that will be used by your data pipeline:

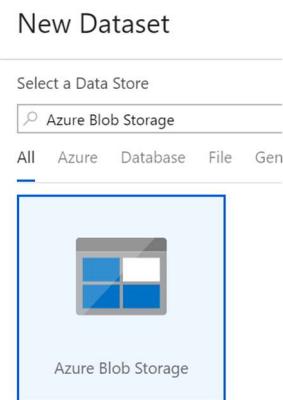
- **MDWResources_NYCTaxiData**: references MDWResources storage account container that contains source data files.
- **MDWResources_NYCTaxiLookup**: references MDWResources storage account that contains a .csv file with all taxi location codes and names.
- **MDWASQLDW_StagingNYCTaxiData**: references the table Staging.NYCTaxiData in the Azure SQL Data Warehouse database MDWASQLDW.
- **MDWASQLDW_StagingNYCLocationLookup**: references the table [Staging].[NYCTaxiLocationLookup] in the Azure SQL Data Warehouse database MDWASQLDW and acts as destination of lookup data copied from **MDWResources_NYCTaxiLookup**.
- **MDWDataLake_NYCTaxiData**: references your MDWDataLake-suffix storage account. It acts as the destination of files copied from **MDWResources_NYCTaxiData** and also as a data source when copying data to **MDWASQLDW_StagingNYCTaxiData**.

IMPORTANT: Execute these steps on your host computer

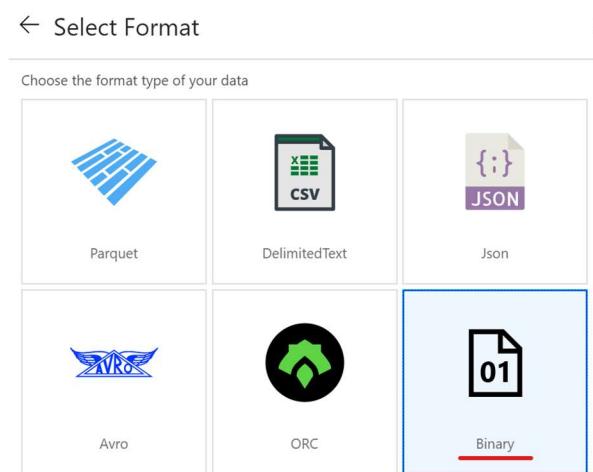
1. Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel. Under **Factory Resources** tab, click the ellipsis (...) next to **Datasets** and then click **Add Dataset** to create a new dataset.



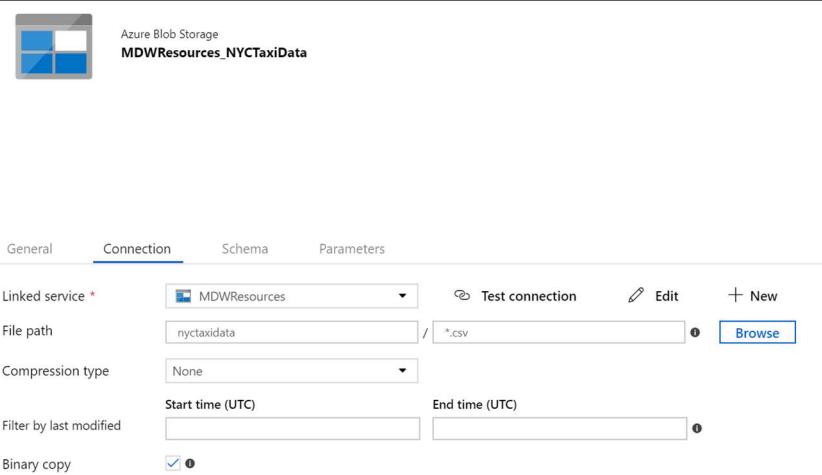
2. Type “Azure Blob Storage” in the search box and select **Azure Blob Storage**. Click **Continue**.



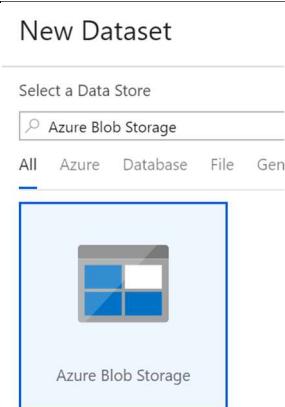
3. On the **Select Format** blade, select **Binary** and click **Continue**.



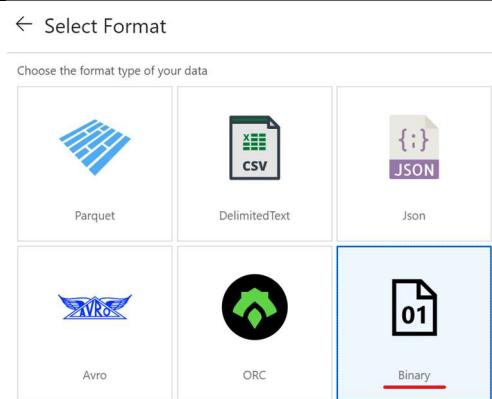
- On the **New Data Set** tab, enter the following details:
 - General > Name:**
MDWResources_NYCTaxiData
 - Connection > Linked Service:**
MDWResources
 - Connection > File Path:** nyctaxidata / *.csv
 - Connection > Binary Copy:** Checked
- Leave remaining fields with default values.



- Repeat the process to create another dataset, this time referencing the NYCTaxiData container in your MDWDataLake storage account.
- Type “Azure Blob Storage” in the search box and select **Azure Blob Storage**. Click **Continue**.



8. On the Select Format blade, select **Binary** and click **Continue**.

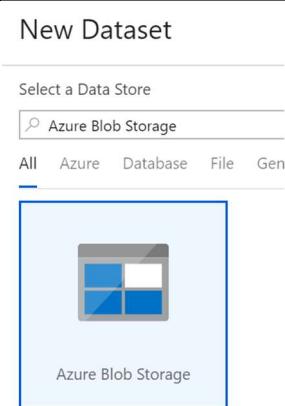


9. On the **New Data Set** tab, enter the following details:
- General > Name:** MDWDataLake_NYCTaxiData
 - Connection > Linked Service:** MDWDataLake
 - Connection > File Path:** nyctaxidata
 - Connection > Binary Copy:** Unchecked
 - Connection > Column names in the first row:** Checked

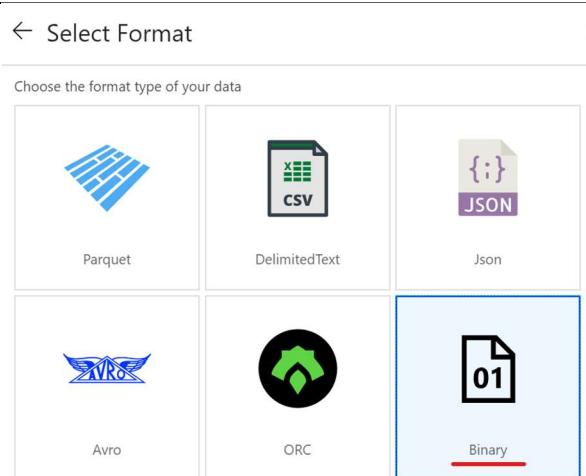
10. Leave remaining fields with default values.

General	Connection	Schema	Parameters
	<p>Connection successful</p> <p>Linked service * <input type="button" value="MDWDataLake"/></p> <p>File path <input type="text" value="nyctaxidata"/> / <input type="button" value="File"/></p> <p>Compression type <input type="button" value="None"/></p> <p>Start time (UTC) <input type="text"/></p> <p>End time (UTC) <input type="text"/></p> <p>Filter by last modified <input type="checkbox"/></p> <p>Binary copy <input type="checkbox"/></p> <p>File format settings</p> <p>File format <input type="button" value="Text format"/> <input type="button" value="Detect Text Format"/> <input type="button" value="Preview data"/></p> <p>Column delimiter <input type="button" value="Comma (,)"/> <input type="checkbox" value="Use custom delimiter"/></p> <p>Row delimiter <input type="button" value="Carriage Return + Line feed (\r\n)"/> <input type="checkbox" value="Use custom delimiter"/></p> <p>Skip line count <input type="text" value="0"/></p> <p>Column names in the first row <input checked="" type="checkbox"/></p>		

11. Repeat the process to create another dataset, this time referencing the NYCTaxiLookup container in your MDWResources storage account.
12. Type “Azure Blob Storage” in the search box and select **Azure Blob Storage**. Click **Continue**.

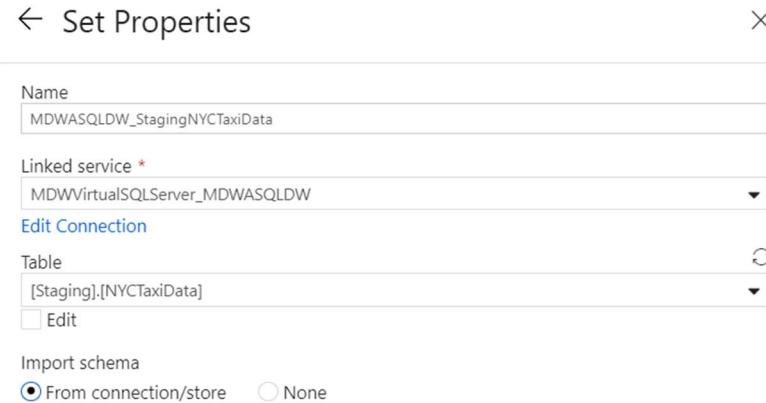


13. On the Select Format blade, select **Binary** and click **Continue**.



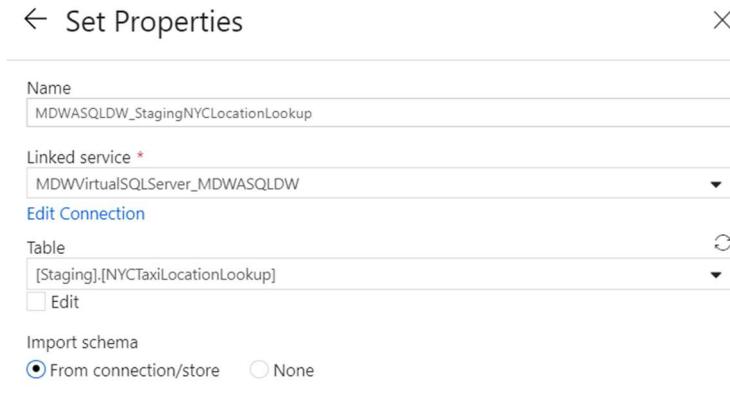
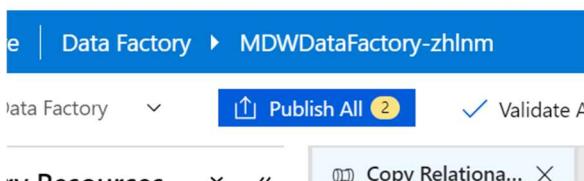
<p>14. On the New Data Set tab, enter the following details:</p> <ul style="list-style-type: none"> a. General > Name: MDWResources_NYCTaxiLookup b. Connection > Linked Service: MDWResources c. Connection > File Path: nyctaxilookup / taxi_zone_lookup.csv d. Connection > Binary Copy: Unchecked e. Connection > Column names in the first row: Checked. f. Connection > Quote character: " (double-quote) <p>15. Leave remaining fields with default values.</p>	
<p>16. Repeat the process to create another dataset, this time referencing the Staging.NYCTaxiData in your Azure SQL Data Warehouse database.</p> <p>17. Type “Azure SQL Data Warehouse” in the search box and select Azure SQL Data Warehouse. Click Continue.</p>	

18. On the **Set Properties** blade, enter the following details:
- Name:** MDWASQLDW_StagingNYCTaxiData
 - Linked Service:** MDWSQLVirtualServer_MDWASQLDW
 - Table:** [Staging].[NYCTaxiData]
19. Leave remaining fields with default values.



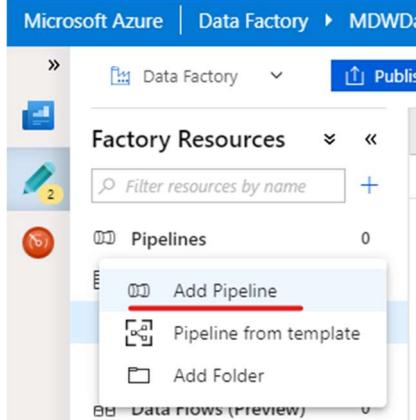
20. Repeat the process to create another dataset, this time referencing the Staging.NYCLocationLookup in your Azure SQL Data Warehouse database.
21. Type “Azure SQL Data Warehouse” in the search box and select **Azure SQL Data Warehouse**. Click **Finish**.



<p>22. On the Set Properties blade, enter the following details:</p> <ul style="list-style-type: none"> a. Name: MDWASQLDW_StagingNYCTaxiLocationLookup b. Linked Service: MDWSQLVirtualServer_MDWASQLDW c. Table: [Staging].[NYCTaxiLocationLookup] <p>23. Leave remaining fields with default values.</p>	
<p>24. Publish your dataset changes by clicking the Publish all button.</p>	

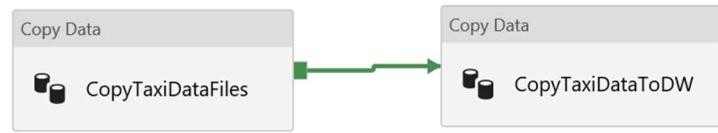
Create and Execute Pipeline

IMPORTANT: Execute these steps on your host computer

<ol style="list-style-type: none"> 1. Open the Azure Data Factory portal and click the Author option on the left-hand side panel. Under Factory Resources tab, click the ellipsis (...) next to Pipelines and then click Add Pipeline to create a new dataset. 2. On the New Pipeline tab, enter the following details: <ul style="list-style-type: none"> a. General > Name: Copy Taxi Data 3. Leave remaining fields with default values. 	
---	---

<p>4. From the Activities panel, type “Copy Data” in the search box. Drag the Copy Data activity on to the design surface. This copy activity will copy data files from MDWResources to MDWDatalake.</p> <p>5. Select the Copy Data activity and enter the following details:</p> <ol style="list-style-type: none"> General > Name: CopyTaxiDataFiles Source > Source dataset: MDWResources_NYCTaxiData Sink > Sink dataset: MDWDatalake_NYCTaxiData Sink > Copy Behavior: Preserve Hierarchy <p>6. Leave remaining fields with default values.</p>	
<p>7. Repeat the process to create another Copy Data Activity, this time to copy data from the files in your data lake to your SQL Data Warehouse.</p> <p>8. From the Activities panel, type “Copy Data” in the search box. Drag the Copy Data activity on to the design surface.</p> <p>9. Select the Copy Data activity and enter the following details:</p> <ol style="list-style-type: none"> General > Name: CopyTaxiDataToDW Source > Source dataset: MDWDatalake_NYCTaxiData Sink > Sink dataset: MDWASQLDW_StagingNYCTaxiData Sink > Pre Copy Script: truncate table Staging.NYCTaxiData <p>10. Leave remaining fields with default values.</p>	

11. Create a **Success** (green) precedence constraint between CopyTaxiDataFiles and CopyTaxiDataToDW. You can do it by dragging the green connector from CopyTaxiDataFiles and landing the arrow onto CopyTaxiDataToDW.



12. Repeat the process to create another Copy Data Activity, this time to copy taxi location lookup data from MDWResources to your SQL Data Warehouse.

13. From the Activities panel, type “Copy Data” in the search box. Drag the Copy Data activity on to the design surface.

14. Select the Copy Data activity and enter the following details:

a. **General > Name:**

CopyTaxiLookupDataToDW

b. **Source > Source dataset:**

MDWResources_NYCTaxiLookup

c. **Sink > Sink dataset:**

MDWASQLDW_StagingNYCLocationLookup

d. **Sink > Pre Copy Script:** truncate table

Staging.NYCTaxiLocationLookup

e. **Settings > Enable staging:** Checked

f. **Settings > Staging account linked service:**

MDWDataLake

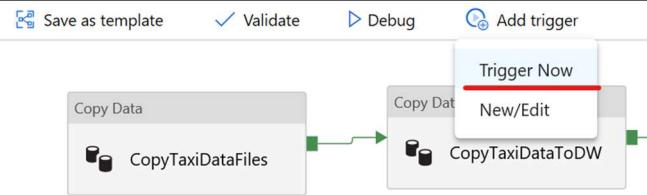
g. **Settings > Storage Path:** polybase

15. Leave remaining fields with default values.

General	Source	Sink 1	Mapping	Settings
<p>Source dataset * MDWResources_NYCTaxiLookup</p> <p><input checked="" type="checkbox"/> Copy file recursively</p> <p>Max concurrent connections</p>				
<p>Sink</p> <p>Sink dataset * MDWASQLDW_StagingNYCLocationLookup</p> <p><input checked="" type="checkbox"/> Allow PolyBase</p> <p>Reject type: Value</p> <p>Reject value: 0</p> <p><input checked="" type="checkbox"/> Use Type default</p> <p>Pre-copy script: truncate table Staging.NYCTaxiLocationLookup</p>				
<p>Settings</p> <p>Degree of copy parallelism: Auto</p> <p>Fault tolerance: Abort activity on first incompatibl...</p> <p><input checked="" type="checkbox"/> Enable staging</p> <p>Staging settings</p> <p>Staging account linked service: MDWDataLake</p> <p>Storage Path: polybase</p>				

<p>16. From the Activities panel, type “Stored Procedure” in the search box. Drag the Stored Procedure activity onto the design surface. This activity will execute the Staging.spNYCLoadTaxiDataSummary to generate aggregated taxi ride information.</p> <p>17. Select the Stored Procedure activity and enter the following details:</p> <ul style="list-style-type: none"> a. General > Name: LoadTaxiDataSummary b. SQL account > Linked service: MDWSQLVirtualServer_MDWASQLDW c. Stored procedure > Stored procedure name: [Staging].[spNYCLoadTaxiDataSummary] <p>18. Leave remaining fields with default values.</p>	
<p>19. Create Success (green) from CopyDataToDW to LoadTaxiDataSummary and from CopyTaxiLookupData to LoadTaxiDataSummary.</p>	<pre> graph LR A[Copy Data
CopyTaxiDataFiles] --> B[Copy Data
CopyTaxiDataToDW] C[Copy Data
CopyTaxiLookupData...] --> B B --> D[Stored Procedure
LoadTaxiDataSummary] </pre>
<p>20. Publish your pipeline changes by clicking the Publish all button.</p>	<p>Data Factory ▶ MDWDataFactory-zhlm</p> <p>Factory ▾ ↑ Publish All 1 ✓ Validate All</p> <p>Resources ▾ << ⌂ Copy Relationa... *</p>

21. To execute the pipeline, click on **Add trigger** menu and then **Trigger Now**.
 22. On the **Pipeline Run** blade, click **Finish**.

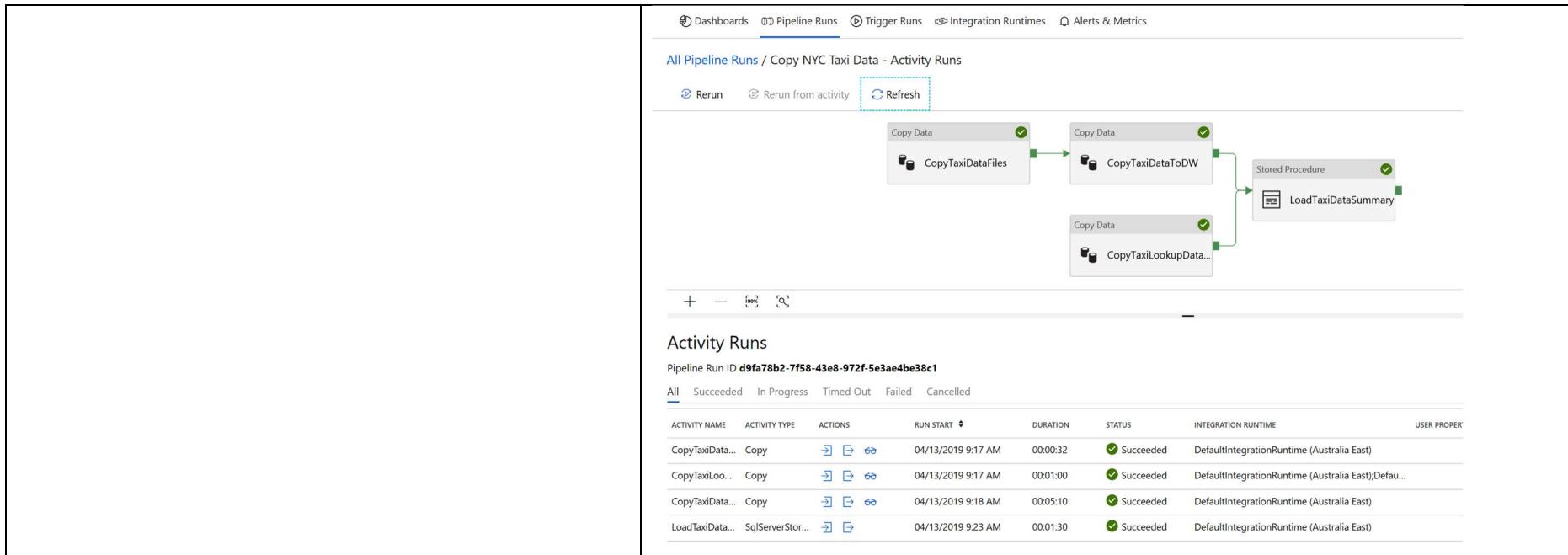


23. To monitor the execution of your pipeline, click on the **Monitor** menu on the left-hand side panel.
 24. You should be able to see the **Status** of your pipeline execution on the right-hand side panel.

Pipeline Name	Actions	Run Start	Duration	Triggered By	Status	Parameters
Copy Relational Data		04/08/2019, 8:55:42 AM	00:01:19	Manual trigger	Succeeded	

25. Click the **View Activity Runs** button for detailed information about each activity execution in the pipeline. The whole execution should last between 7-8 minutes.

Pipeline Name	View Activity Runs	R
Copy NYC Taxi Data		0



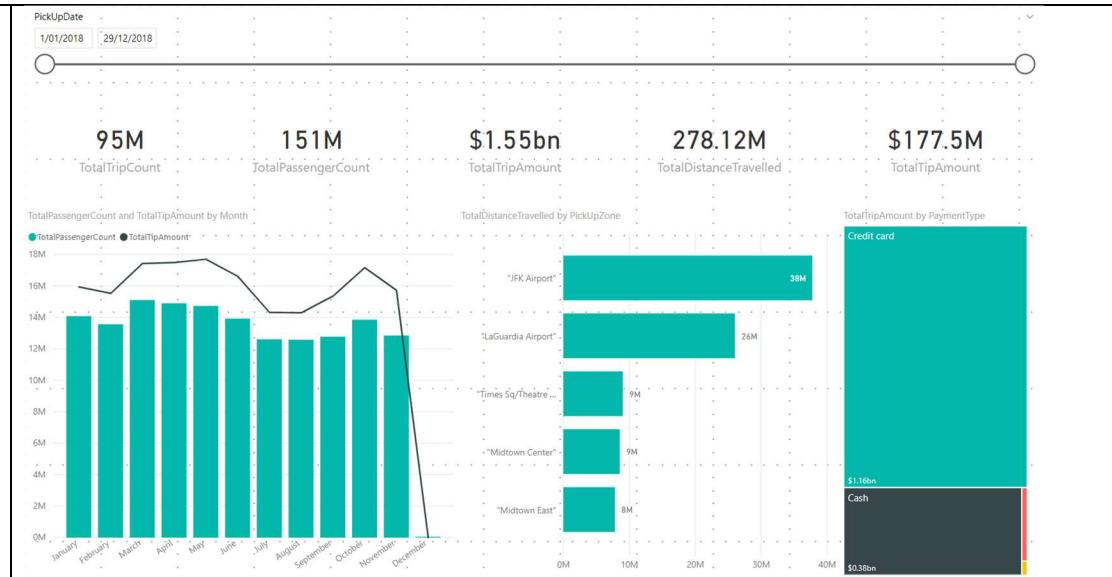
Visualize Data with Power BI

In this section you are going to use Power BI to visualize data from Azure SQL Data Warehouse. The Power BI report will use an Import connection to query Azure SQL Data Warehouse and visualise Motor Vehicle Collision data from the table you loaded in the previous exercise.

IMPORTANT: Steps to be executed on MDWDesktop	
<ol style="list-style-type: none">8. On MDWDesktop, download the Power BI report from the link https://aka.ms/MDWLab2 and save it in the Desktop.9. Open the file MDWLab2.pbit with Power BI Desktop.10. When prompted to enter the value of the MDWSQLVirtualServer parameter, type the full server name: mdwsqvirtualserver-suffix.database.windows.net11. Click Load.	<p>MDWLab2</p> <p>Modern Data Warehouse Power BI Template - Lab 2</p> <p>MDWSQLVirtualServer</p> <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">mdwsqvirtualserver-zhlnm.database.windows.net</div>
<ol style="list-style-type: none">12. When prompted to enter credentials, select Database from the left-hand side panel and enter the following details:<ol style="list-style-type: none">a. User name: mdwadminb. Password: P@ssw0rd123!13. Leave remaining fields with their default values.14. Click Connect.	<p>The screenshot shows the Power BI credential dialog. On the left, there's a sidebar with 'Windows', 'Database' (which is selected and highlighted in grey), and 'Microsoft account'. The main area has a title 'SQL Server database' with a close button 'X'. Below it, there's a connection icon and the text 'mdqsqlvirtuslserver-cf3xh.database.windows.net;M...'. A 'User name' field contains 'mdwadmin' and a 'Password' field contains 'P@ssw0rd123!'. At the bottom, there's a dropdown menu labeled 'Select which level to apply these settings to' with 'mdqsqlvirtuslserver-cf3xh.database.windows.net' selected.</p>

15. Once data finish loading interact with the report by changing the PickUpDate slicer and by clicking on the other visualisations.

16. Save your work and close Power BI Desktop.



Lab 3: Explore Big Data using Azure Databricks

In this lab you will use Azure Databricks to explore the New York Taxi data files you saved in your data lake in the previous exercise. Using a Databricks you will connect to the data lake and query taxi ride details.

Lab Architecture



Build an Azure Databricks notebook to explore the data files you saved in your data lake in the previous exercise. You will use Python and SQL commands to open a connection to your data lake and query data from data files.

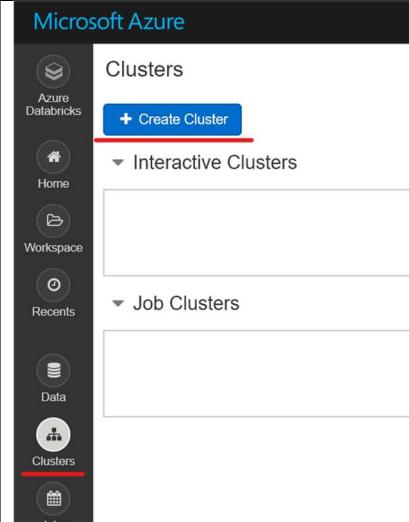
Create Azure Databricks Cluster

In this section you are going to create an Azure Databricks cluster that will be used to execute notebooks.

IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, navigate to the MDW-Lab resource group and locate the Azure Databricks MDWDatabricks.
2. On the **MDWDatabricks** blade, click the **Launch Workspace** button. The Azure Databricks portal will open on a new browser tab.

3. On the **Azure Databricks** portal, click the **Clusters** button on the left-hand side menu.
4. On the **Clusters** blade, click **+ Create Cluster**.



5. On the Create Cluster blade, type "MDWDatabricksCluster" in the **Cluster Name** field. Leave all other fields with their default values.
6. Click **Create Cluster**. It should take around 5 minutes for the cluster to be fully operational.

Create Cluster

New Cluster

2-8 Workers: 28.0-112.0 GB Memory, 8-32 Cores, 1.5-6 DBU
1 Driver: 14.0 GB Memory, 4 Cores, 0.75 DBU

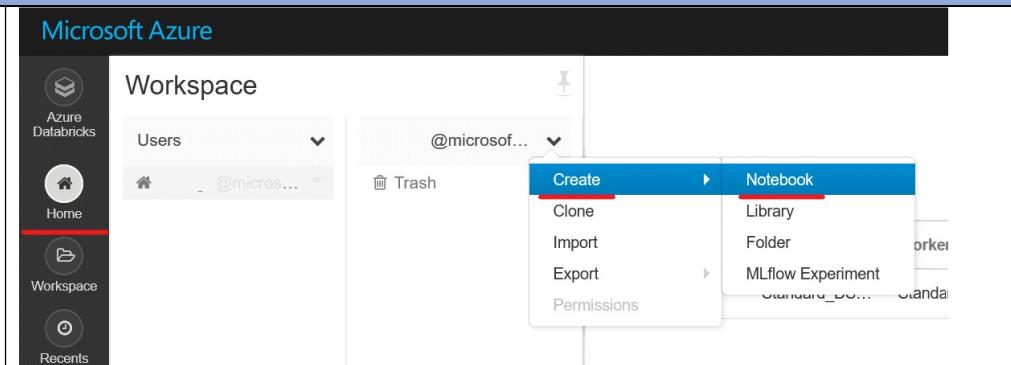
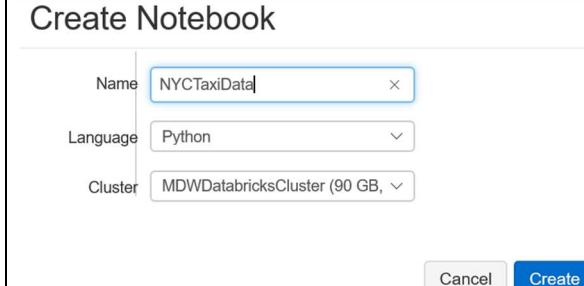
Cluster Name	MDWDatabricksCluster
Cluster Mode	Standard
Databricks Runtime Version	Runtime: 5.2 (Scala 2.11, Spark 2.4.0)
Python Version	3
Autopilot Options	
<input checked="" type="checkbox"/> Enable autoscaling	
<input checked="" type="checkbox"/> Terminate after 120 minutes of inactivity	
Worker Type	Standard_DS3_v2 14.0 GB Memory, 4 Cores, 0.75 DBU
Min Workers	2
Max Workers	8
Driver Type	Same as worker 14.0 GB Memory, 4 Cores, 0.75 DBU

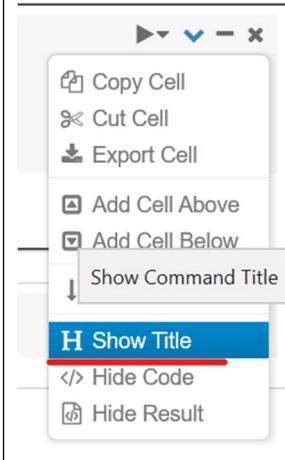
New The default Python version for clusters was changed from major version 2 to 3.

Create an Azure Databricks Notebook

In this section you are going to create an Azure Databricks notebook that will be used to explore the taxi data files you copied to your data lake in Lab 2.

IMPORTANT: The full solution scripts for each notebook command are saved in the Student's Computer path **C:\ADSIAD\Lab\Lab3\Solution**

IMPORTANT: Execute these steps on your host computer	
<ol style="list-style-type: none">1. On the Azure Databricks portal, click the Home button on the left-hand side menu.2. On the Workspace blade, click the down arrow next to your user name and then click Create > Notebook.	
<ol style="list-style-type: none">3. On the Create Notebook pop-up window type "NYCTaxiData" in the Name field.4. Ensure you have the Language field set to Python and the Cluster field is set to MDWDatabricksCluster.5. Click Create.	

<p>6. On the Cmd 1 cell, click the Edit button on the top right-hand corner of the cell and then click Show Title.</p> <p>7. Type “Setup connection to MDWDataLake storage account” in the cell title.</p>	 <p>The screenshot shows a Jupyter Notebook cell titled "Cmd 1". A context menu is open over the cell, with the "Show Title" option highlighted in blue and underlined with a red line. Other options in the menu include "Copy Cell", "Cut Cell", "Export Cell", "Add Cell Above", "Add Cell Below", "Show Command Title", "Hide Code", and "Hide Result".</p> <p>Cmd 1</p> <p>Setup connection to MDWDataLake storage account</p> <p>1 </p>
<p>8. On the Cmd 1 cell, you will invoke the Spark API to establish a connection to your MDWDalaLake storage account. For this you will need to retrieve the name and key of your MDWDataLake storage account from the Azure Portal.</p> <p>9. Remember to replace <your MDWDataLake storage account name> with mdwdatalakesuffix and to replate <your MDWDataLake storage account key> with the storage account key.</p>	<p>Python command:</p> <pre>spark.conf.set("fs.azure.account.key.<your MDWDataLake storage account name>.blob.core.windows.net", "<your MDWDataLake storage account key>")</pre>
<p>10. Press Shift + Enter to execute and create a new notebook cell. Set the title of the Cmd 2 cell to “Define NYCTaxiData schema and load data into a Data Frame”</p>	

<p>11. In the Cmd 2 cell, define a new StructType object that will contain the definition of the data frame schema.</p> <p>12. Using the schema defined above, initialise a new data frame by invoking the Spark API to read the contents of the nyctaxidata container in the MDWDataLake storage account.</p> <p>13. Remember to replace <your MDWDataLake storage account name> with mdwdatalakesuffix and to replate <your MDWDataLake storage account key> with the storage account key.</p>	<p>Python command:</p> <pre>from pyspark.sql.types import * nycTaxiDataSchema = StructType([StructField("VendorID", IntegerType(), True), StructField("tpep_pickup_datetime", DateType(), True), StructField("tpep_dropoff_datetime", DateType(), True), StructField("passenger_count", IntegerType(), True), StructField("trip_distance", DoubleType(), True), StructField("RatecodeID", IntegerType(), True), StructField("store_and_fwd_flag", StringType(), True), StructField("PULocationID", IntegerType(), True), StructField("DOLocationID", IntegerType(), True), StructField("payment_type", IntegerType(), True), StructField("fare_amount", DoubleType(), True), StructField("extra", DoubleType(), True), StructField("mta_tax", DoubleType(), True), StructField("tip_amount", DoubleType(), True), StructField("tolls_amount", DoubleType(), True), StructField("improvement_surcharge", DoubleType(), True), StructField("total_amount", DoubleType(), True)])</pre> <p>dfNYCTaxiData = spark.read.format('csv').options(header='true', schema=nycTaxiDataSchema).load('wasbs://nyctaxidata@<your MDWDataLake storage account name>.blob.core.windows.net/')</p>
<p>14. Hit Shift + Enter to execute the command and create a new cell.</p>	<p>Define NYCTaxiData schema and load data into a Data Frame</p> <pre>1 from pyspark.sql.types import * 2 3 nycTaxiDataSchema = StructType([4 StructField("VendorID", IntegerType(), True), 5 StructField("tpep_pickup_datetime", DateType(), True), 6 StructField("tpep_dropoff_datetime", DateType(), True), 7 StructField("passenger_count", IntegerType(), True), 8 StructField("trip_distance", DoubleType(), True), 9 StructField("RatecodeID", IntegerType(), True), 10 StructField("store_and_fwd_flag", StringType(), True), 11 StructField("PULocationID", IntegerType(), True), 12 StructField("DOLocationID", IntegerType(), True), 13 StructField("payment_type", IntegerType(), True), 14 StructField("fare_amount", DoubleType(), True), 15 StructField("extra", DoubleType(), True), 16 StructField("mta_tax", DoubleType(), True), 17 StructField("tip_amount", DoubleType(), True), 18 StructField("tolls_amount", DoubleType(), True), 19 StructField("improvement_surcharge", DoubleType(), True), 20 StructField("total_amount", DoubleType(), True)]) 21 22 dfNYCTaxiData = spark.read.format('csv').options(header='true', schema=nycTaxiDataSchema).load('wasbs://nyctaxidata@mdwdatalakezhlnm.blob.core.windows.net/')</pre>

15. Set the title of the Cmd 3 cell to “Display Data Frame Content”.																															
16. In the Cmd 3 cell, call the display function to show the contents of the data frame <code>dfNYCTaxiData</code> .	<p>Python command:</p> <pre>display(dfNYCTaxiData)</pre>																														
17. Hit Shift + Enter to execute the command and create a new cell.	<p>Cmd 3</p> <h3>Display Data Frame Content</h3> <pre>1 display(dfNYCTaxiData) ▶ (1) Spark Jobs</pre> <table border="1"> <thead> <tr> <th>VendorID</th> <th>tpep_pickup_datetime</th> <th>tpep_dropoff_datetime</th> <th>passenger_count</th> <th>trip_dis</th> </tr> </thead> <tbody> <tr><td>1</td><td>2018-01-01 00:21:05</td><td>2018-01-01 00:24:23</td><td>1</td><td>.50</td></tr> <tr><td>1</td><td>2018-01-01 00:44:55</td><td>2018-01-01 01:03:05</td><td>1</td><td>2.70</td></tr> <tr><td>1</td><td>2018-01-01 00:08:26</td><td>2018-01-01 00:14:21</td><td>2</td><td>.80</td></tr> <tr><td>1</td><td>2018-01-01 00:20:22</td><td>2018-01-01 00:52:51</td><td>1</td><td>10.20</td></tr> <tr><td>1</td><td>2018-01-01 00:00:10</td><td>2018-01-01 00:27:00</td><td>2</td><td>2.50</td></tr> </tbody> </table>	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_dis	1	2018-01-01 00:21:05	2018-01-01 00:24:23	1	.50	1	2018-01-01 00:44:55	2018-01-01 01:03:05	1	2.70	1	2018-01-01 00:08:26	2018-01-01 00:14:21	2	.80	1	2018-01-01 00:20:22	2018-01-01 00:52:51	1	10.20	1	2018-01-01 00:00:10	2018-01-01 00:27:00	2	2.50
VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_dis																											
1	2018-01-01 00:21:05	2018-01-01 00:24:23	1	.50																											
1	2018-01-01 00:44:55	2018-01-01 01:03:05	1	2.70																											
1	2018-01-01 00:08:26	2018-01-01 00:14:21	2	.80																											
1	2018-01-01 00:20:22	2018-01-01 00:52:51	1	10.20																											
1	2018-01-01 00:00:10	2018-01-01 00:27:00	2	2.50																											
18. Set the title of the Cmd 4 cell to “Create Temp View”																															
19. In the Cmd 4 cell, call the <code>createOrReplaceTempView</code> method of the data frame object to create a temporary view of the data in memory.	<p>Python command:</p> <pre>dfNYCTaxiData.createOrReplaceTempView('NYCTaxiDataTable')</pre>																														
20. Hit Shift + Enter to execute the command and create a new cell.	<p>Cmd 4</p> <h3>Create Temporary View</h3> <pre>1 dfNYCTaxiData.createOrReplaceTempView('NYCTaxiDataTable')</pre> <p>Command took 0.05 seconds -- by fabraga@microsoft.com at 4/13/2019, 4:50:29 P</p>																														
21. Set the title of the Cmd 5 cell to “Use SQL to count NYC Taxi Data records”																															
22. In the Cmd 5 cell, change the default language to SQL using the <code>%sql</code> command.	<p>Command:</p> <pre>%sql</pre>																														
23. Write a SQL query to retrieve the total number of records in the <code>NYCTaxiDataTable</code> view.	<pre>select count(*) from NYCTaxiDataTable</pre>																														

24. Hit Shift + Enter to execute the command and create a new cell.	<p>Cmd. 5</p> <h3>Use SQL to count NYC Taxi Data records</h3> <pre>1 %sql 2 3 select count(*) from NYCTaxiDataTable</pre> <p>▶ (1) Spark Jobs</p> <table border="1"> <tr><td>count(1)</td></tr> <tr><td>102804250</td></tr> </table> <p> </p> <p>Command took 52.64 seconds -- by fabraga@microsoft.com at 5/7/2019, 10:04:38 PM on MDWDatabricksCluster</p>	count(1)	102804250
count(1)			
102804250			
25. Set the title of the Cmd 6 cell to “Use SQL to filter NYC Taxi Data records”			
26. In the Cmd 6 cell, write a SQL query to filter taxi rides that happened on the Apr, 7 th 2018 that had more than 5 passengers.	<p>Command</p> <pre>%sql select cast(tpep_pickup_datetime as date) as pickup_date , tpep_dropoff_datetime , passenger_count , total_amount from NYCTaxiDataTable where cast(tpep_pickup_datetime as date) = '2018-04-07' and passenger_count > 5</pre>		

27. Hit Shift + Enter to execute the command.

Cmd 6

Use SQL to filter NYC Taxi Data records

```
1 %sql
2
3 select cast(tpep_pickup_datetime as date) as pickup_date
4     , tpep_dropoff_datetime
5     , passenger_count
6     , total_amount
7 from NYCTaxiDataTable
8 where cast(tpep_pickup_datetime as date) = '2018-04-07'
9     and passenger_count > 5
10
```

▶ (3) Spark Jobs

pickup_date	tpep_dropoff_datetime
2018-04-07	2018-04-07 00:27:08
2018-04-07	2018-04-07 00:11:57
2018-04-07	2018-04-07 00:12:12
2018-04-07	2018-04-07 00:17:20

28. Set the title of the **Cmd 7** cell to “Use SQL to aggregate NYC Taxi Data records and visualize data”

29. In the **Cmd 7** cell, write a SQL query to aggregate records and return total number of rides by payment type. Use the following command.
30. Hit Ctrl + Enter to execute the command. Results will be displayed in a grid in the cell.

Command

%sql

```
select case payment_type
        when 1 then 'Credit card'
        when 2 then 'Cash'
        when 3 then 'No charge'
        when 4 then 'Dispute'
        when 5 then 'Unknown'
        when 6 then 'Voided trip'
    end as PaymentType
    , count(*) as TotalRideCount
from NYCTaxiDataTable
group by payment_type
order by TotalRideCount desc
```

31. Click the **Bar chart** button to see results as a bar chart.



Cmd 7

Use SQL to aggregate NYC Taxi Data records and visualize data

```
1 %sql
2
3 select case payment_type
4     when 1 then 'Credit card'
5     when 2 then 'Cash'
6     when 3 then 'No charge'
7     when 4 then 'Dispute'
8     when 5 then 'Unknown'
9     when 6 then 'Voided trip'
10    end as PaymentType
11    , count(*) as TotalRideCount
12 from NYCTaxiDataTable
13 group by payment_type
14 order by TotalRideCount desc
```

▶ (1) Spark Jobs

A bar chart titled 'Use SQL to aggregate NYC Taxi Data records and visualize data'. The Y-axis is labeled 'TotalRideCount' and ranges from 0.00 to 70M. The X-axis is labeled 'PaymentType' and includes categories: Credit card, Cash, No charge, Dispute, and Unknown. The chart shows two bars: 'Credit card' at approximately 65M and 'Cash' at approximately 28M. The other categories have zero or negligible values.

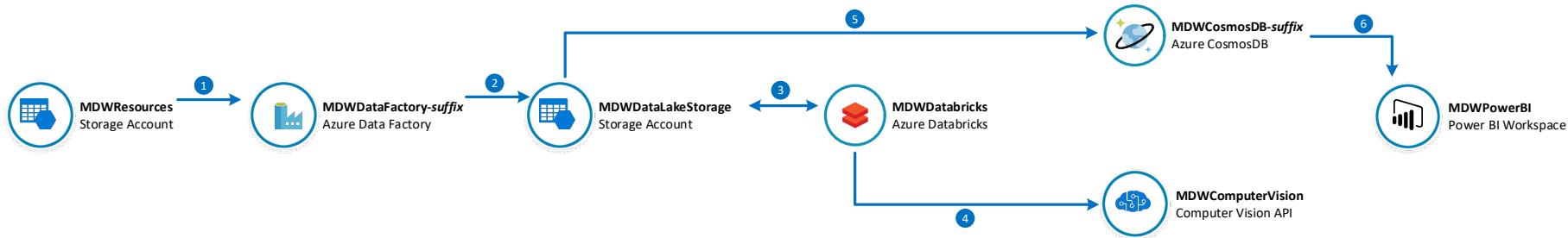
PaymentType	TotalRideCount
Credit card	~65M
Cash	~28M
No charge	0
Dispute	0
Unknown	0

Command took 56.37 seconds -- by fabrige@microsoft.com at 4/14/2019, 2:11:36 PM on MDWDataBricksCluster

Lab 4: Add AI to your Big Data Pipeline with Cognitive Services

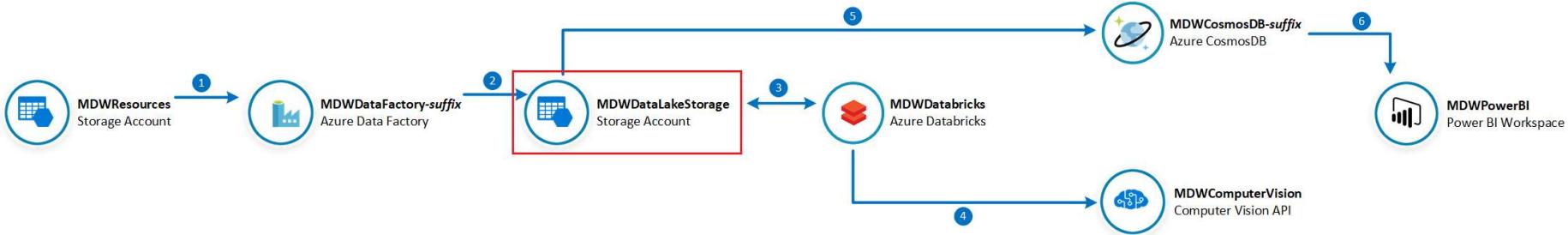
In this lab you will use Azure Data Factory to download New York City images to your data lake. Then, as part of the same pipeline, you are going to use an Azure Databricks notebook to invoke Computer Vision Cognitive Service to generate metadata documents and save them in a CosmosDB database. You will use Power BI to visualise the images and their AI-generated metadata.

Lab Architecture



- 1 Build an Azure Data Factory Pipeline to copy image files from shared Azure Storage
- 2 Save image files to your data lake
- 3 For each image in your data lake, invoke an Azure Databricks notebook that will take the image URL as parameter
- 4 For each image call the Azure Computer Vision Cognitive service to generate image metadata. Metadata files are saved back in your data lake
- 5 Copy metadata JSON documents into your Cosmos DB database
- 6 Visualize images and associated metadata using Power BI

Create NYCImages and NYCImageMetadata Containers in Azure Blob Storage

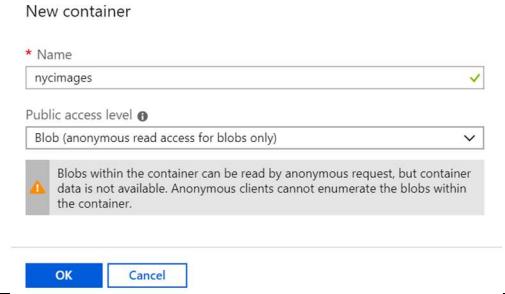
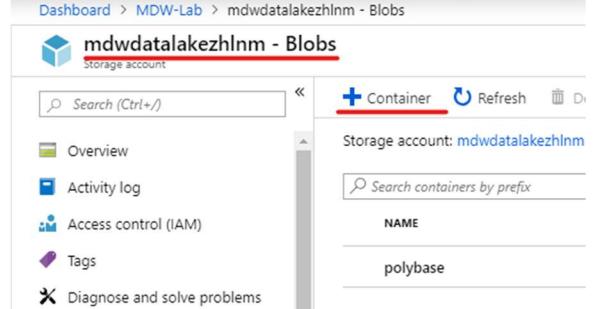
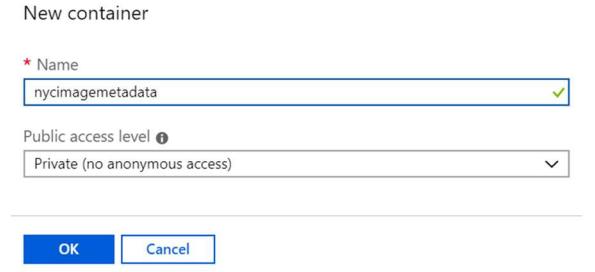


In this section you will create a container in your MDWDataLake that will be used as a repository for the NYC image files. You will copy 30 files from the MDWResources Storage Account into your NYCTaxiData container.

IMPORTANT: Execute these steps on your host computer

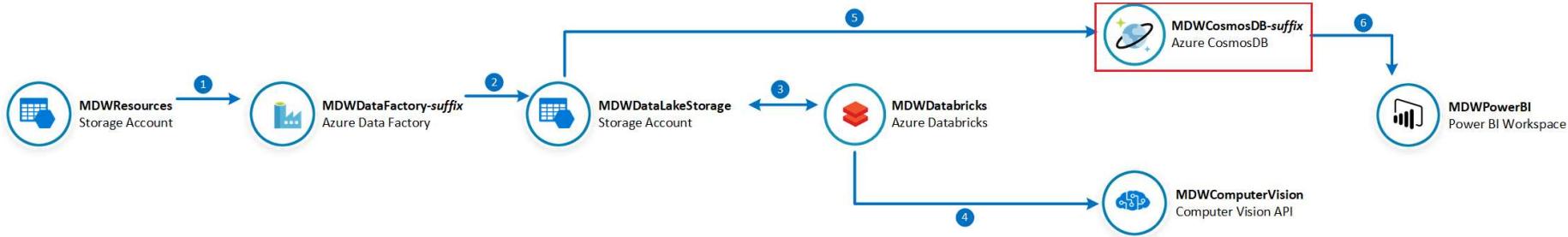
1. In the Azure Portal, go to the lab resource group and locate the Azure Storage account **mdwdatalakesuffix**.
2. On the **Overview** panel, click **Blobs**.

3. On the **mdwdatalakesuffix - Blobs** blade, click **+ Container**.

<p>4. On the New container blade, enter the following details:</p> <ul style="list-style-type: none"> a. Name: nycimages b. Public access level: Blob (anonymous read access for blobs only) <p>5. Click OK to create the new container.</p>	
<p>6. Repeat the process to create the NYCImageMetadata container. This container will be used to host the metadata files generated by Cognitive Services before they can be saved in CosmosDB.</p>	
<p>7. On the New container blade, enter the following details:</p> <ul style="list-style-type: none"> a. Name: nycimagedata b. Public access level: Private (no anonymous access) <p>8. Click OK to create the new container.</p>	

Create CosmosDB database and collection

In this section you will create a CosmosDB database called NYC and a collection called ImageMetadata that will host New York image metadata information.



IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, go to the lab resource group and locate the CosmosDB account **MDWCosmosDB-suffix**.
2. On the **Overview** panel, click **+ Add Container**.

Dashboard > Resource groups > MDW-Lab-S01 > mdwcosmosdb-cf3xh

mdwcosmosdb-cf3xh
Azure Cosmos DB account

Search (Ctrl+ /) <<

Add Container Refresh M

Status : Online

Resource group (change) : MDW-Lab-S01

Subscription (change) : Microsoft Azure

Subscription ID : 96bd7145-ad1...

3. On the **Add Container** blade, enter the following details:
 - a. **Database id > Create new:** NYC
 - b. **Container id:** ImageMetadata
 - c. **Partition key:** /requestId
 - d. **Throughput:** 400
 - e. **Unique keys:** /requestId
4. Click **OK** to create the container.

Add Container X

Start at \$24/mo per database, multiple containers included
[More details](#)

* Database id ?
 Create new Use existing
NYC

Provision database throughput ?

* Container id ?
ImageMetadata

Where did 'fixed' containers go? ?

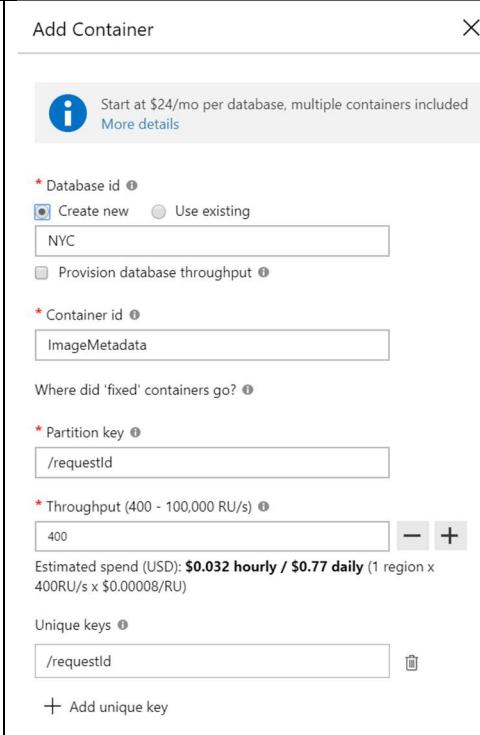
* Partition key ?
/requestId

* Throughput (400 - 100,000 RU/s) ?
400 - +

Estimated spend (USD): **\$0.032 hourly / \$0.77 daily** (1 region x 400RU/s x \$0.00008/RU)

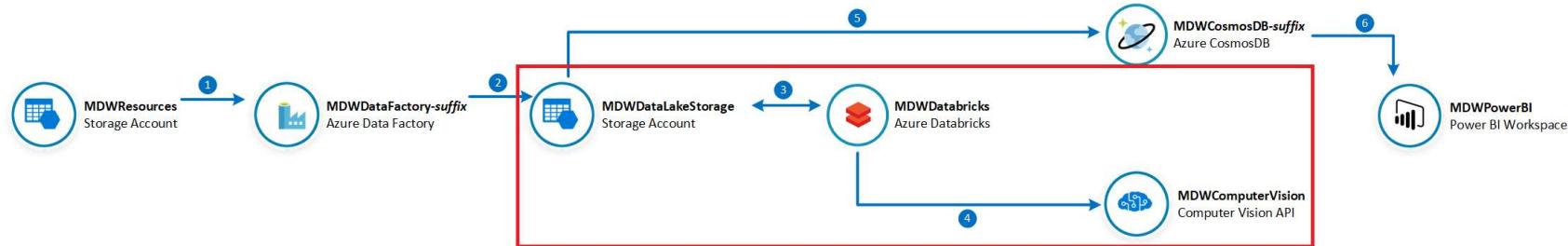
Unique keys ?
/requestId [trash]

+ Add unique key



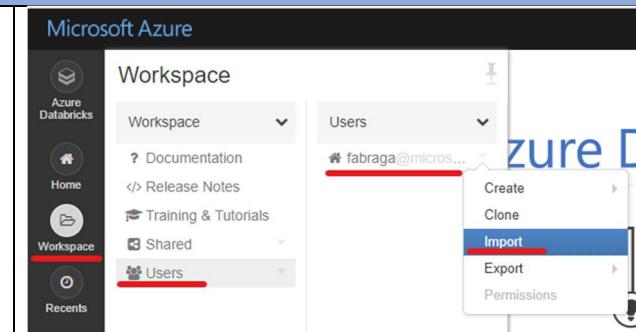
Import Databricks Notebook to Invoke Computer Vision Cognitive Services API

In this section you will import a Databricks notebook to your workspace and fill out the missing details about your Computer Vision API and your Data Lake account. This notebook will be executed from an Azure Data Factory pipeline and it will invoke the Computer Vision API to generate metadata about the images and save the result back to your data lake.

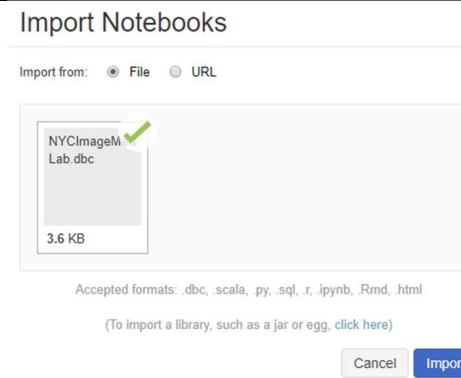


IMPORTANT: Execute these steps on your host computer

1. On the **Azure Databricks** portal, click the **Workspace** button on the left-hand side menu.
2. On the **Workspace** blade, click **Users**.
3. On the **Users** blade, click the arrow next to your user name and then **Import**.



- On the Import Notebooks pop up window, click the browse link to import the notebook file **C:\ADSIAD\Lab\Lab4\NYCImageMetadata-Labdbc**
- Click **Import**.



- On the **NYCImageMetadata-Lab** notebook, go to the cell **Cmd 2 – Define function to invoke Computer Vision API**. You will need to change the function code to include the Computer Vision API details.

Cmd 2

Define function to invoke Computer Vision API

```

1 import requests
2 import json
3
4 def GetNYCImageMetadata(imageUrl):
5     # Replace <MDWComputerVision Subscription Key> with your valid subscription key.
6     subscription_key = "<MDWComputerVision Subscription Key>" REDACTED
7
8     #Replace <MDWComputerVision Base URL> with your valid Computer Vision API base url
9     vision_base_url = "<MDWComputerVision Base URL>" REDACTED #It should look like this one: https://australiaeast.api.cognitive.microsoft.com/
10    analyze_url = vision_base_url + "vision/v2.0/analyze"
11

```

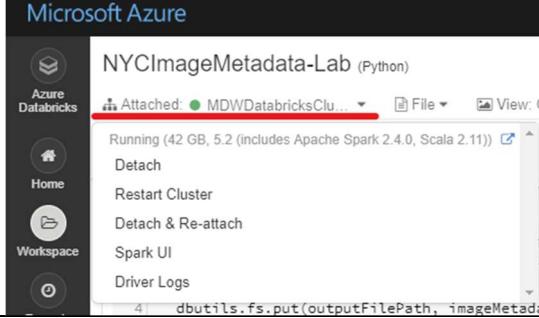
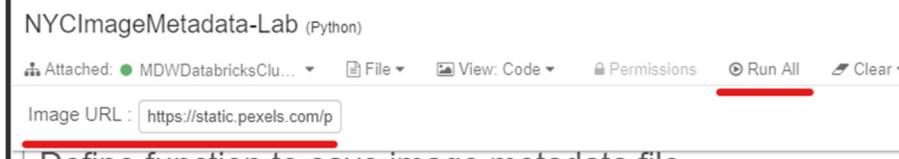
- From the **Azure Portal**, retrieve the **MDWComputerVision** subscription key and base endpoint URL.

Dashboard > MDW-Lab > MDWComputerVision - Keys

MDWComputerVision - Keys REDACTED Cognitive Services

MDWComputerVision - Keys	
<input type="text" value="Search (Ctrl+ /)"/>	↻ Regenerate Key1 ↻ Regenerate Key2
🕒 Overview	NAME
📅 Activity log	MDWComputerVision
👤 Access control (IAM)	KEY 1
🏷️ Tags	REDACTED
✖️ Diagnose and solve problems	KEY 2
💡 Keys	REDACTED

	<p>The screenshot shows the Azure portal's Cognitive Services blade for the 'MDWComputerVision' resource. It displays basic information like Resource group (MDW-Lab), Status (Active), Location (Australia East), Subscription (Microsoft Azure Internal Consumption), and API type (Computer Vision). The Endpoint is listed as https://australiaeast.api.cognitive.microsoft.com/. A 'Manage keys' button is visible.</p>
<p>8. Copy and paste the Key and Endpoint values back in the Databricks notebook.</p>	<p>Cmd 2 Define function to invoke Computer Vision API</p> <pre> 1 import requests 2 import json 3 4 def GetNYCImageMetadata(imageUrl): 5 # Replace <MDWComputerVision Subscription Key> with your valid subscription key. 6 subscription_key = "b094915fea994196bd46042e4253abcf" 7 8 #Replace <MDWComputerVision Base URL> with your valid Computer Vision API base url 9 vision_base_url = "https://australiaeast.api.cognitive.microsoft.com/" #It should look like this one: https://australiaeast.api.cognitive.microsoft.com/ 10 analyze_url = vision_base_url + vision/v2.0/analyze 11 </pre>
<p>9. On the NYCImageMetadata-Lab notebook, go to the cell Cmd 3 – Define function to mount NYC Image Metadata Container. You will need to change the function code to include your data lake storage account details.</p> <p>10. In the dataLakeaccountName variable assignment replace <MDWDataLake storage account name> with mdwdata lakesuffix.</p>	<p>Cmd 3 Define function to mount NYC Image Metadata Container</p> <pre> 1 def MountImageMetadataContainer(): 2 # Replace <MDWDataLake storage account name> with your valid storage account name 3 datalakeAccountName = '<MDWDataLake storage account name>' 4 5 # Replace <MDWDataLake storage account key> with your valid storage account name 6 datalakeAccountKey = '<MDWDataLake storage account key>' 7 </pre>
<p>11. From the Azure Portal, retrieve the MDWDataLakesuffix access key.</p>	<p>The screenshot shows the 'Access keys' section for the 'mdwdatakezhlm' storage account. It displays two access keys: 'key1' and 'key2'. The 'key1' key is shown with its value and a 'Copy' button. Below it is a 'Connection string' field containing the connection string for the storage account.</p>

<p>12. Copy and paste the Access Key Databricks notebook. Replace <MDWDataLake storage account key> with the Access Key you got from the previous step.</p>	 <pre>1 def MountImageMetadataContainer(): 2 # Replace <MDWDataLake storage account name> with your valid storage account name 3 datalakeAccountName = 'mdwdatalakezhlnm' 4 5 # Replace <MDWDataLake storage account key> with your valid storage account name 6 datalakeAccountKey = 'abc2PNNncmrvm7avfeXUc/nSfxb+hGgXG34mz9KyWFB2o8xTla91FoGSbchlP1132iaE1sRVDkaYBBVuMuDkQ==' 7</pre>
<p>13. Attach the notebook to your previously created MDWDatabricksCluster cluster.</p>	
<p>14. Review the notebook code. 15. If you want to test it, you can copy any publicly available image URL and paste it in the Image URL notebook parameter. You can use any of the following image URLs in the list as examples.</p>	<p>Test Image URLs:</p> <p>https://petlifetoday.com/wp-content/uploads/2018/06/wireless-dog-fence.jpg</p> <p>https://static.pexels.com/photos/4204/nature-lawn-blur-flower.jpg</p> <p>https://image.redbull.com/rbcom/052/2017-05-22/89eef344-d24f-4520-8680-8b8f7508b264/0012/0/0/0/2428/3642/800/1/best-beginner-motocross-bikes-ktm-250-sxf.jpg</p> <p>http://www.kiplinger.com/slideshow/investing/T024-S001-the-best-emerging-markets-stocks-for-2019/images/intro.jpg</p>
<p>16. Click Run All to execute the notebook.</p>	

17. After a successful execution you will notice that a new JSON file has been saved in the **NYCImageMetadata** container in your Data Lake.

Cmd 5

Invoke Computer Vision and Save Metadata File

```
1 from urllib.parse import urlparse
2 #invoke Computer Vision API to retrieve image metadata
3 jsonImageMetadata = GetNYCImageMetadata(nycImageUrl)
4
5 #Generate filename for image metadata file.
6 filePath = urlparse(nycImageUrl)[2]
7 filePathParts = filePath.split('/')
8 fileName = filePathParts[len(filePathParts)-1] + '.json'
9
10 #Save image metadata file in the NYCImageMetadata container
11 SaveImageMetadataFile(jsonImageMetadata, fileName)

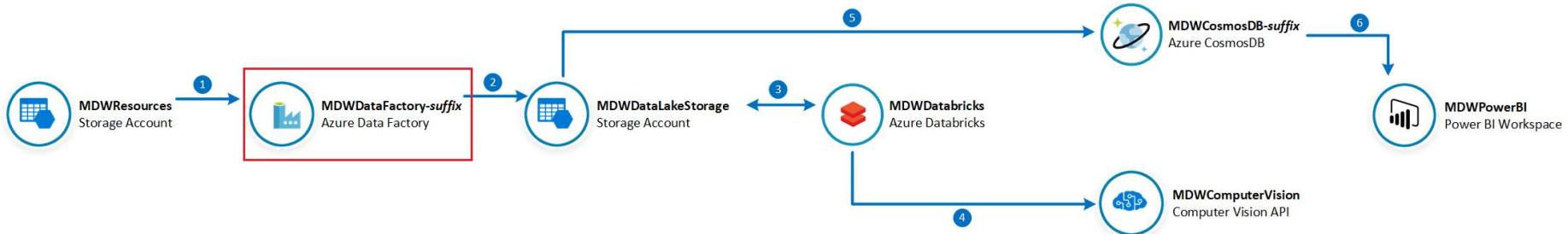
Wrote 939 bytes.
/mnt/NYCImageMetadata/nature-lawn-blur-flower.jpg.json file saved.
Command took 2.44 seconds -- by fabraga@microsoft.com at 4/17/2019, 1:58:49 PM on MDWDatabricksCluster
```

18. Navigate to Azure Portal and check the contents of the **nycimagedata** container in your **MDWDatalakesuffix** storage account.
19. Download the file generated to inspect its contents.
20. **IMPORTANT:** Delete this test file before moving to next steps of this exercise.

The screenshot shows the Azure Storage Explorer interface. The left sidebar lists navigation options: Dashboard, Storage accounts, mdwdatalakezhlm - Blobs, and nycimagedata (which is highlighted with a red underline). The main pane displays the contents of the 'nycimagedata' blob container. It includes a search bar, upload, refresh, and change access level buttons. Below these are sections for Authentication method (Access key) and Location (nycimagedata). A search bar for blobs by prefix is also present. The list of blobs shows one item: 'nature-lawn-blur-flower.jpg.json'. The file name is highlighted with a red underline.

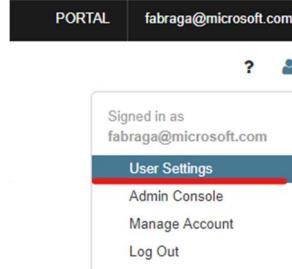
Create Databricks Linked Service in Azure Data Factory

In this section you will create a Databricks linked service in Azure Data Factory. Through this linked service you will be able to create a data pipelines to copy NYC images to your data lake and integrate Databricks notebooks to its execution.

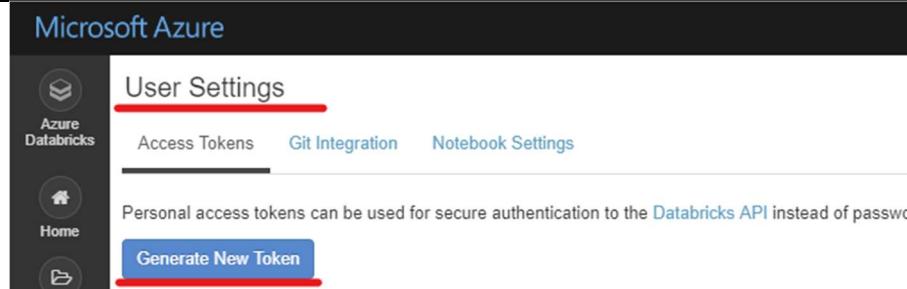


IMPORTANT: Execute these steps on your host computer

1. On the **Azure Databricks** portal, click the **User** icon on the top right-hand corner of the screen.
2. Click on the **User Settings** menu item.



3. On the **User Settings** blade, under the **Access Tokens** tab, click **Generate New Token**.



- On the **Generate New Token** pop-up window, enter “Azure Data Factory Access” in the **Comment** field. Leave **Lifetime (days)** with the default value of 90 days.
- Click **Generate**.

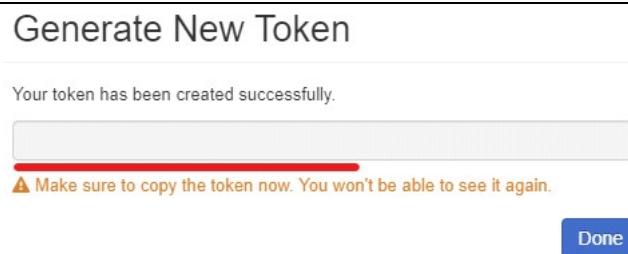
Generate New Token

Comment
Azure Data Factory Access

Lifetime (days) ?
90

Cancel **Generate**

- IMPORTANT:** Copy the generated access token to Notepad and save it. You won't be able to retrieve it once you close this window.



- Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel. Under **Connections** tab, click **Linked Services** and then click **+ New** to create a new linked service connection.

Microsoft Azure | Data Factory > MDWDataFactory-zhlnm

Search resource

Data Factory Publish All Validate All Refresh Discard A

Factory Resources <>

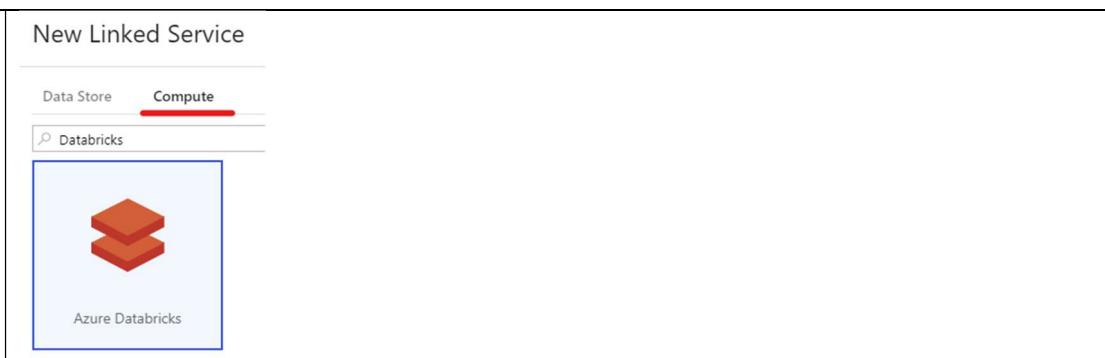
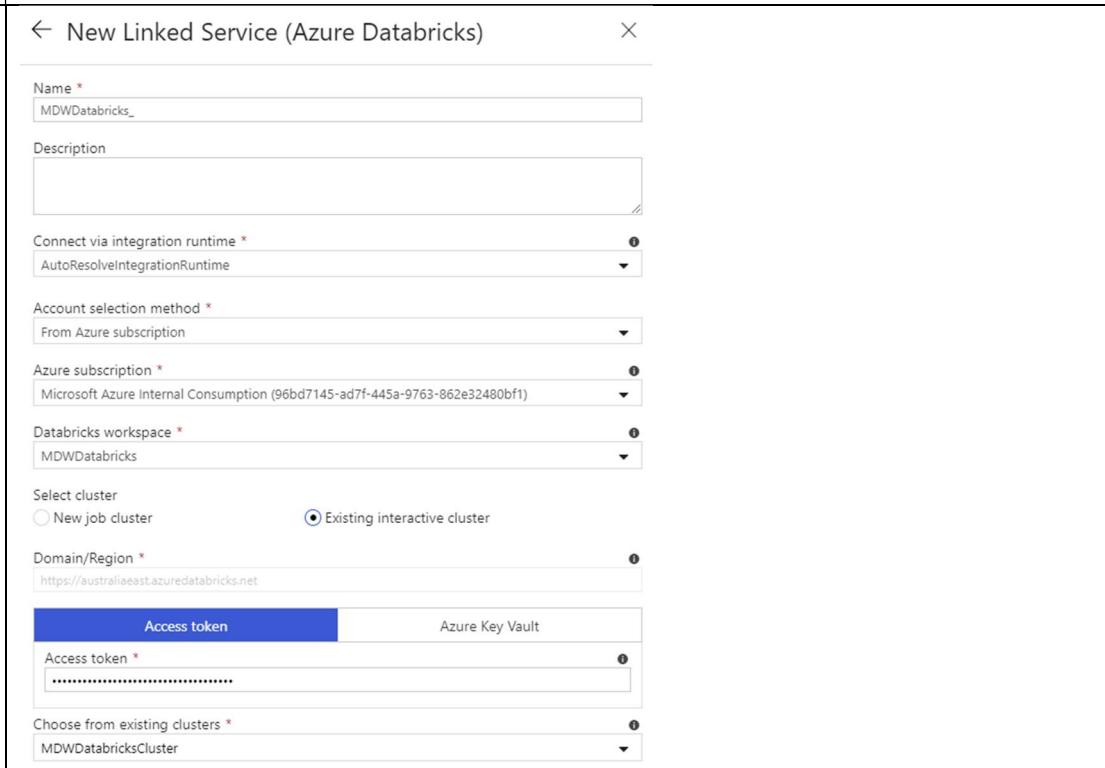
- Pipelines 0
- Datasets 0
- Data Flows (Preview) 0

Connections

Linked Services Integration Runtimes

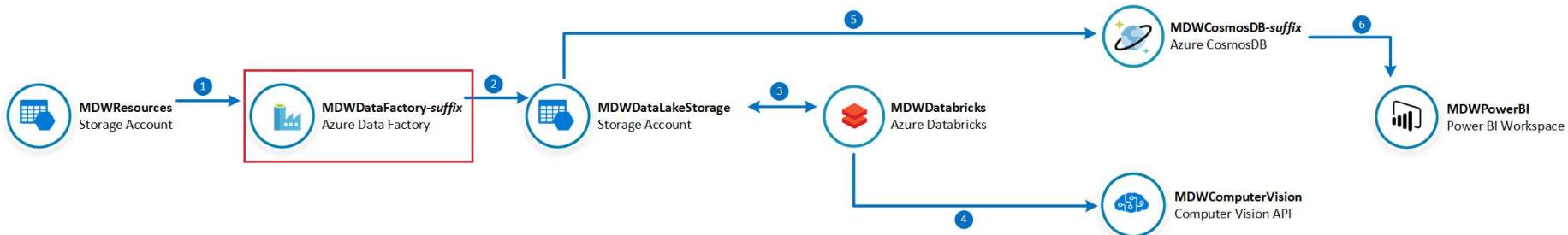
+ New

Name

<p>8. On the New Linked Service blade, click the Compute tab.</p> <p>9. Type “Azure Databricks” in the search box to find the Azure Databricks linked service.</p> <p>10. Click Continue.</p>	
<p>11. On the New Linked Service (Azure Databricks) blade, enter the following details:</p> <ul style="list-style-type: none"> a. Name: MDWDatabricks b. Connect via integration runtime: AutoResolveIntegrationRuntime c. Account selection method: From Azure subscription d. Azure subscription: <select your subscription> e. Databricks workspace: MDWDatabricks f. Select cluster: Existing interactive cluster g. Access token: <copy and paste access token here> h. Choose from existing clusters: MDWDatabricksCluster <p>12. Click Test connection to make sure you entered the correct connection details. You should see a “Connection successful” message above the button.</p> <p>13. If the connection was successful, then click Finish. If you got an error message, please review the connection details above and try again.</p>	

Create CosmosDB Linked Service in Azure Data Factory

In this section you will create a CosmosDB linked service in Azure Data Factory. CosmosDB will be used as the final repository of image metadata information generated by the Computer Vision API. Power BI will then be used to visualise the CosmosDB data.



IMPORTANT: Execute these steps on your host computer

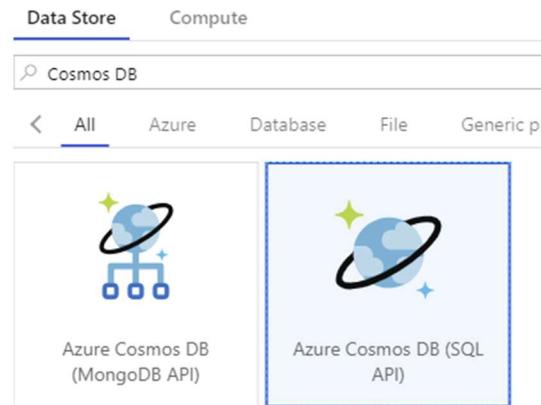
1. Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel. Under **Connections** tab, click **Linked Services** and then click **+ New** to create a new linked service connection.

Screenshot of the Microsoft Azure Data Factory portal:

- The top navigation bar shows "Microsoft Azure | Data Factory > MDWDataFactory-zhlnm".
- The left sidebar displays "Factory Resources" with icons for Data Factory, Pipelines, Datasets, and Data Flows (Preview). Below these are counts: Pipelines (0), Datasets (0), and Data Flows (Preview) (0).
- The main content area is titled "Connections X". It has tabs for "Linked Services" (selected) and "Integration Runtimes".
- Under the "Linked Services" tab, there is a button labeled "+ New".
- Below the button, there is a search bar labeled "Name" with a dropdown arrow.
- At the bottom right of the content area, there is a "Action" button.

2. On the **New Linked Service** blade, click the **Data Store** tab.
3. Type “Cosmos DB” in the search box to find the **Azure Cosmos DB (SQL API)** linked service.
4. Click **Continue**.

New Linked Service



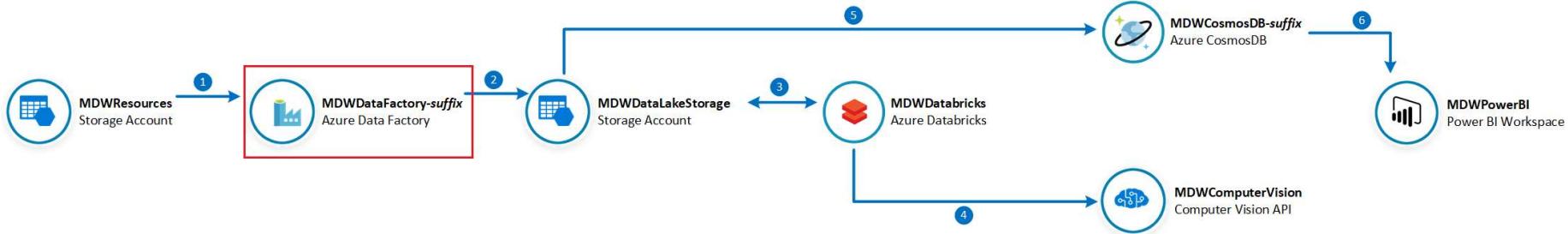
14. On the **New Linked Service (Azure Databricks)** blade, enter the following details:
 - i. **Name:** MDWCosmosDB
 - j. **Connect via integration runtime:** AutoResolveIntegrationRuntime
 - k. **Account selection method:** From Azure subscription
 - l. **Azure subscription:** <select your subscription>
 - m. **Cosmos DB account name:** mdwcosmosdb-suffix
 - n. **Database name:** NYC
15. Click **Test connection** to make sure you entered the correct connection details. You should see a “Connection successful” message above the button.
16. If the connection was successful, then click **Finish**. If you got an error message, please review the connection details above and try again.

← New Linked Service (Azure Cosmos DB (SQL API)) ×

Name *	MDWCosmosDB_
Description	
Connect via integration runtime *	AutoResolveIntegrationRuntime
Connection String	
Azure Key Vault	
Account selection method	<input checked="" type="radio"/> From Azure subscription <input type="radio"/> Enter manually
Azure subscription	Microsoft Azure Internal Consumption (96bd7145-ad7f-445a-9763-862e32480bf1)
Cosmos DB account name *	mdwcosmosdb-zhlnm
Database name *	NYC

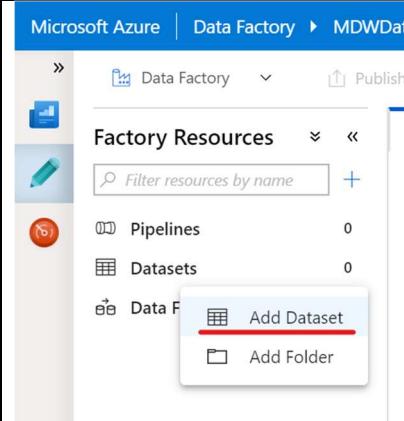
Create Azure Data Factory data sets.

In this section you will create an Azure Data Factory data sets that will be used in the data pipeline.

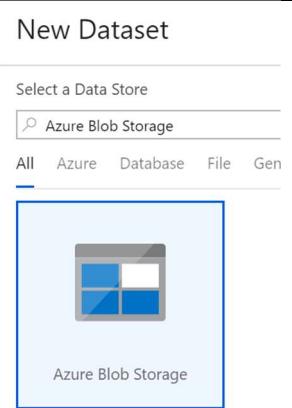


IMPORTANT: Execute these steps on your host computer

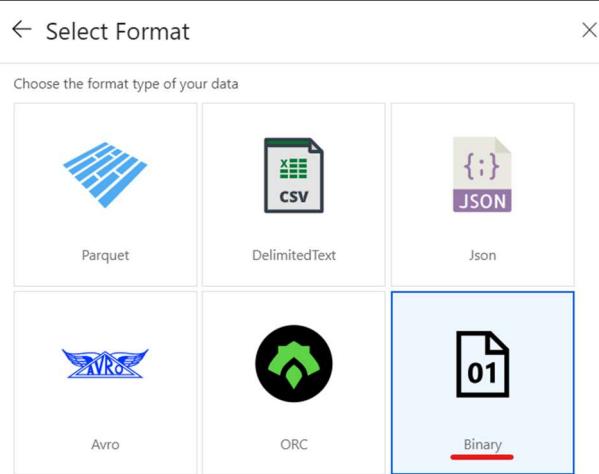
1. Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel. Under **Factory Resources** tab, click the ellipsis (...) next to **Datasets** and then click **Add Dataset** to create a new dataset.



2. Type “Azure Blob Storage” in the search box and select **Azure Blob Storage**. Click **Finish**.



3. On the **Select Format** blade, select **Binary** and click **Continue**.



- On the **New Data Set** tab, enter the following details:
 - General > Name:**
MDWResources_NYCIImages
 - Connection > Linked Service:**
MDWResources
 - Connection > File Path:** nycimages
 - Connection > Binary Copy:** Checked
- Leave remaining fields with default values.

Azure Blob Storage
MDWResources_NYCIImages

General Connection Schema Parameters

Linked service * MDWResources Test connection Edit + New

File path nycimages / File Browse

Compression type None

Start time (UTC) End time (UTC)

Filter by last modified

Binary copy

- Repeat the process to create another dataset, this time referencing the NYCIImages container in your MDWDataLake storage account.
- Type “Azure Blob Storage” in the search box and select **Azure Blob Storage**. Click **Finish**.

New Dataset

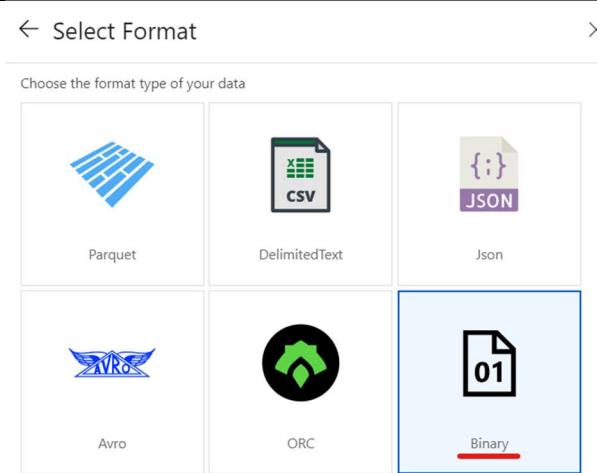
Select a Data Store

Azure Blob Storage

All Azure Database File Gen

Azure Blob Storage

8. On the **Select Format** blade, select **Binary** and click **Continue**.



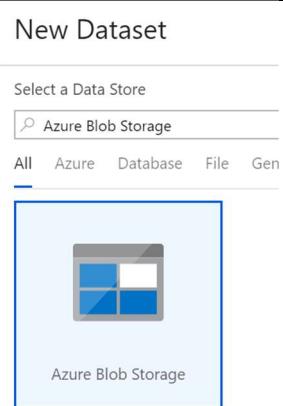
9. On the **New Data Set** tab, enter the following details:

- a. **General > Name:**
MDWDataLake_NYCTaxiData
- b. **Connection > Linked Service:**
MDWDataLake
- c. **Connection > File Path:** nycimages

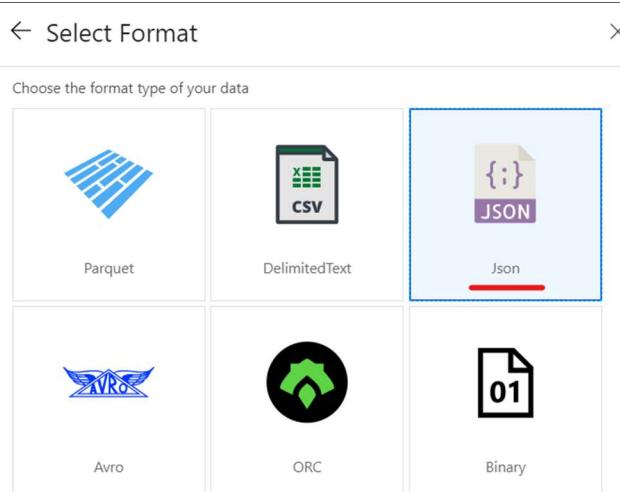
10. Leave remaining fields with default values.

General		Connection		Schema		Parameters	
Linked service *	MDWDataLake	Test connection	<input type="button" value="Edit"/>	<input type="button" value="+ New"/>			
File path	nycimages	/	File	<input type="button" value="Browse"/>			
Compression type	None						
Filter by last modified		Start time (UTC)		End time (UTC)			
Binary copy	<input checked="" type="checkbox"/>						

11. Repeat the process to create another dataset, this time referencing the NYCIImageMetadata container in your MDWDataLake storage account.
12. Type “Azure Blob Storage” in the search box and select **Azure Blob Storage**. Click **Finish**.



13. On the **Select Format** blade, select **JSON** and click **Continue**.



14. On the **New Data Set** tab, enter the following details:

- a. **General > Name:**
MDWDataLake_NYCIImageMetadata
- b. **Connection > Linked Service:**
MDWDataLake
- c. **Connection > File Path:**
nycimagedata
- d. **File format:** JSON format

15. Leave remaining fields with default values.

The screenshot shows the 'New Data Set' dialog in Azure Data Studio. The 'Connection' tab is selected. The 'Linked service' dropdown is set to 'MDWDataLake'. The 'File path' field contains 'nycimagedata'. The 'File format' dropdown is set to 'JSON format'. Other tabs like 'General', 'Schema', and 'Parameters' are visible at the top. A preview of the data is shown below the form.

16. Repeat the process to create another dataset, this time referencing the ImageMetadata collection in your MDWCosmosDB database.

17. Type “Cosmos DB” in the search box and select **Azure Cosmos DB (SQL API)**. Click **Continue**.

New Dataset

The screenshot shows the 'Select a Data Store' search results. A search bar contains 'Cosmos'. Below it, a list of data stores is shown, with 'Azure Cosmos DB (SQL API)' highlighted by a blue dashed border. Other options include 'Azure Cosmos DB (MongoDB API)', 'Azure', 'Database', 'File', 'Generic protocol', 'NoSQL', and 'S'.

18. On the **New Data Set** tab, enter the following details:

- a. **General > Name:**
MDWCosmosDB_NYCIImageMetadata
- b. **Connection > Linked Service:**
MDWCosmosDB
- c. **Collection name:** ImageMetadata

19. Leave remaining fields with default values.

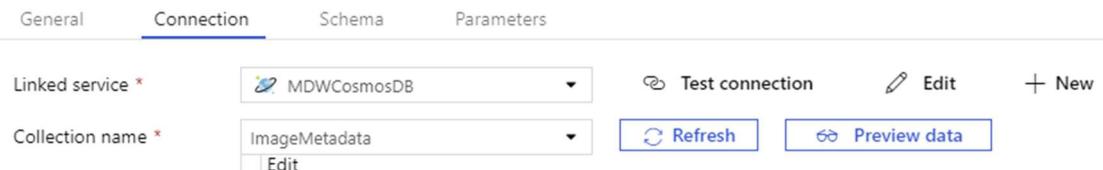


DocumentDB Collection (SQL API)
MDWCosmosDB_NYCIImageMetadata

General Connection Schema Parameters

Linked service * MDWCosmosDB Test connection Edit + New

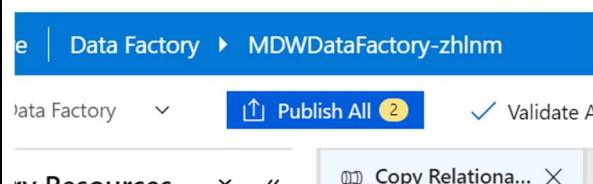
Collection name * ImageMetadata Refresh Preview data Edit



20. Publish your dataset changes by clicking the **Publish all** button.

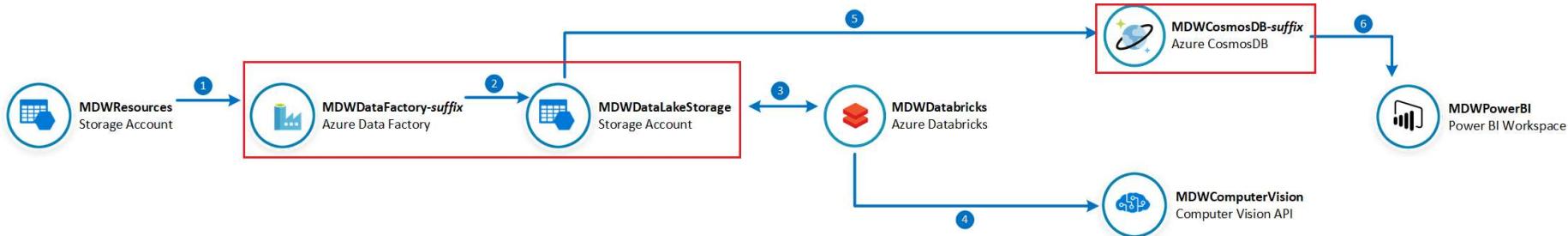
e | Data Factory > MDWDataFactory-zhlnm

Data Factory Publish All 2 Validate A



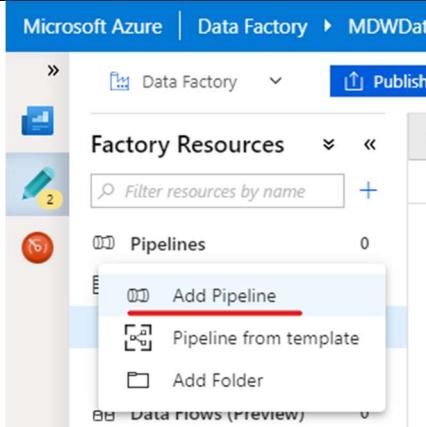
Create Azure Data Factory pipeline to generate and save image metadata to Cosmos DB.

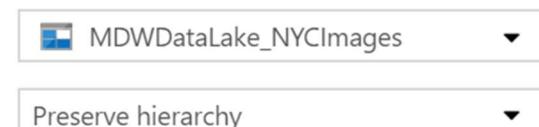
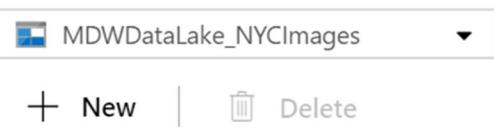
In this section you will create an Azure Data Factory pipeline to copy New York images from MDWResources into your MDWDataLakesuffix storage account. The pipeline will then execute a Databricks notebook for each image and generate a metadata file in the NYCIImageMetadata container. The pipeline finishes by saving the image metadata content in a CosmosDB database.



IMPORTANT: Execute these steps on your host computer

1. Open the **Azure Data Factory** portal and click the **Author** option on the left-hand side panel. Under **Factory Resources** tab, click the ellipsis (...) next to **Pipelines** and then click **Add Pipeline** to create a new dataset.
2. On the New Pipeline tab, enter the following details:
 - a. **General > Name:** Copy NYC Images
 - b. **Variables > + New >**
 - i. **Name:** ImageMetadataContainerUrl
 - ii. **Default Value:**
<https://mdwdatalakesuffix.blob.core.windows.net/nycimages/>
3. Leave remaining fields with default values.

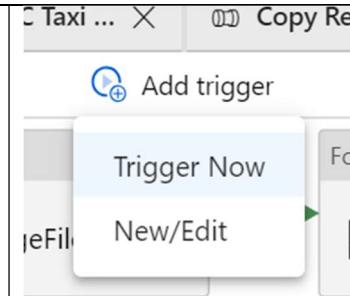


	<p>General Parameters Variables Output</p> <p>+ New Delete</p> <table border="1"> <thead> <tr> <th>NAME</th><th>TYPE</th><th>DEFAULT VALUE</th></tr> </thead> <tbody> <tr> <td>ImageMetadataContainerUr</td><td>String</td><td>https://mdwdatalakezhlnm.blob.core.windows.net</td></tr> </tbody> </table>	NAME	TYPE	DEFAULT VALUE	ImageMetadataContainerUr	String	https://mdwdatalakezhlnm.blob.core.windows.net
NAME	TYPE	DEFAULT VALUE					
ImageMetadataContainerUr	String	https://mdwdatalakezhlnm.blob.core.windows.net					
<p>4. From the Activities panel, type “Copy Data” in the search box. Drag the Copy Data activity on to the design surface. This copy activity will copy image files from MDWResources to MDWDatalake.</p> <p>5. Select the Copy Data activity and enter the following details:</p> <ul style="list-style-type: none"> a. General > Name: CopyImageFiles b. Source > Source dataset: MDWResources_NYCIImages c. Sink > Sink dataset: MDWDatalake_NYCIImages d. Sink > Copy Behavior: Preserve Hierarchy <p>6. Leave remaining fields with default values.</p>	<p>General Source Sink Mapping Settings</p> <p>Source dataset *</p>  <p>General Source Sink Mapping Settings</p> <p>Sink dataset *</p> 						
<p>7. From the Activities panel, type “Get Metadata” in the search box. Drag the Get Metadata activity on to the design surface. This activity will retrieve a list of image files saved in the NYCIImages container by the previous CopyImageFiles activity.</p> <p>8. Select the Get Metadata activity and enter the following details:</p> <ul style="list-style-type: none"> a. General > Name: GetImageFileList b. Source > Source dataset: MDWDatalake_NYCIImages c. Source > Field list: Child Items <p>9. Leave remaining fields with default values.</p>	<p>General Dataset User Properties</p> <p>Dataset *</p>  <p>Field list</p> <p>+ New Delete</p> <p><input type="checkbox"/> ARGUMENT</p> <p>Child Items</p>						

<p>10. Create a Success (green) precedence constraint between CopyImageFiles and GetImageFileList activities. You can do it by dragging the green connector from CopyImageFiles and landing the arrow onto GetImageFileList.</p>	
<p>11. From the Activities panel, type “ForEach” in the search box. Drag the ForEach activity on to the design surface. This ForEach activity will act as a container for other activities that will be executed for each image files returned by the GetImageFileList activity.</p> <p>12. Select the ForEach activity and enter the following details:</p> <ul style="list-style-type: none"> a. General > Name: ForEachImage b. Settings > Items: @activity('GetImageFileList').output.childItems <p>13. Leave remaining fields with default values.</p>	<p>General Settings Activities (0) User Properties</p> <p>Sequential <input type="checkbox"/></p> <p>Batch count <input type="text"/></p> <p>Items <input type="text" value="@activity('GetImageFileList').output.childItems"/></p>
<p>14. Create a Success (green) precedence constraint between GetImageFileList and ForEachImage activities. You can do it by dragging the green connector from GetImageFileList and landing the arrow onto ForEachImage.</p>	
<p>15. Double-click the ForEachImage activity to edit its contents. Note the design context is displayed on the top left-hand side of the design canvas.</p>	
<p>16. From the Activities panel, type “Notebook” in the search box. Drag the Notebook activity on to the design surface. This Notebook activity will pass the image URL as a parameter to the Databricks notebook we created previously.</p> <p>17. Select the Notebook activity and enter the following details:</p>	<p>General Azure Databricks Settings User Properties</p> <p>Databricks Linked Service * <input type="text" value="MDWDatabricks"/></p>

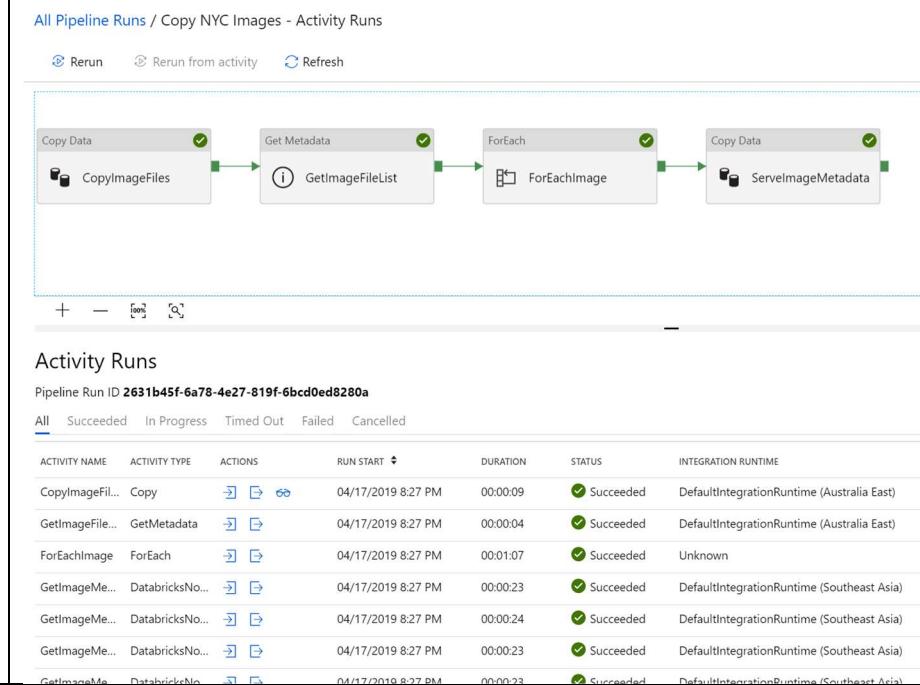
<p>a. General > Name: GetImageMetadata</p> <p>b. Azure Databricks > Databricks Linked Service: MDWDatabricks</p> <p>c. Settings > Notebook path: Click Browse and navigate to <code>/Users/your-user-name/NYCIImageMetadata-Lab</code></p> <p>d. Base Parameters: + New</p> <ul style="list-style-type: none"> i. <code>nyclImageUrl: @concat(variables('ImageMetadataContainerUrl'), item().name)</code> <p>18. Leave remaining fields with default values.</p>	<p>The screenshot shows the 'Settings' tab of a Databricks linked service. Under 'Base Parameters', there is a single entry: 'nyclImageUrl' with the value '@concat(variables('ImageMetadataContainerUrl'), item().name)' highlighted with a blue box.</p>
<p>19. Navigate back to the “Copy NYC Images” pipeline canvas.</p> <p>20. From the Activities panel, type “Copy Data” in the search box. Drag the Copy Data activity on to the design surface. This copy activity will copy image metadata from the JSON files sitting on the NYCIImageMetadata container in MDWDataLake to the ImageMetadata collection on CosmosDB.</p> <p>21. Select the Copy Data activity and enter the following details:</p> <ul style="list-style-type: none"> a. General > Name: ServelImageMetadata b. Source > Source dataset: MDWDataLake_NYCIImageMetadata c. Sink > Sink dataset: MDWCosmosDB_NYCIImageMetadata <p>22. Leave remaining fields with default values.</p>	<p>The screenshot shows the Azure Data Factory pipeline canvas with the following sequence of activities connected by green arrows:</p> <pre> graph LR A[Copy Data] --> B[Get Metadata] B --> C[ForEach] C --> D[ForEachImage] D --> E[Copy Data] </pre> <p>Activity A is labeled 'Copy Data' with a 'CopyImageFiles' icon. Activity B is labeled 'Get Metadata' with a 'GetImageFileList' icon. Activity C is labeled 'ForEach' with a 'ForImage' icon. Activity D is labeled 'ForEachImage' with a 'ForEachImage' icon. Activity E is labeled 'Copy Data' with a 'ServelImageMetadata' icon.</p> <p>The screenshot shows the 'Data Factory' blade with the pipeline name 'MDWDataFactory-zhlnm'. Below it is a 'Publish All' button with a count of 1, indicating changes ready to be published. Other buttons include 'Validate All' and 'Resources'.</p>
<p>23. Create a Success (green) precedence constraint between GetImageFileList and ForEachImage activities. You can do it by dragging the green connector from GetImageFileList and landing the arrow onto ForEachImage.</p> <p>24. Publish your pipeline changes by clicking the Publish all button.</p>	<p>The screenshot shows the 'Data Factory' blade with the pipeline name 'MDWDataFactory-zhlnm'. Below it is a 'Publish All' button with a count of 1, indicating changes ready to be published. Other buttons include 'Validate All' and 'Resources'.</p>

25. To execute the pipeline, click on **Add trigger** menu and then **Trigger Now**.
 26. On the **Pipeline Run** blade, click **Finish**.



27. To monitor the execution of your pipeline, click on the **Monitor** menu on the left-hand side panel.
 28. You should be able to see the **Status** of your pipeline execution on the right-hand side panel.
 29. Click the **View Activity Runs** button for detailed information about each activity execution in the pipeline. The whole execution should last between 7-8 minutes.

All	Succeeded	In Progress	Queued	Failed	Cancelled
Pipeline Name	Actions	Run Start	Duration	Triggered By	Status
Copy NYC Images		04/17/2019, 8:27:35 PM	00:01:43	Manual trigger	Succeeded

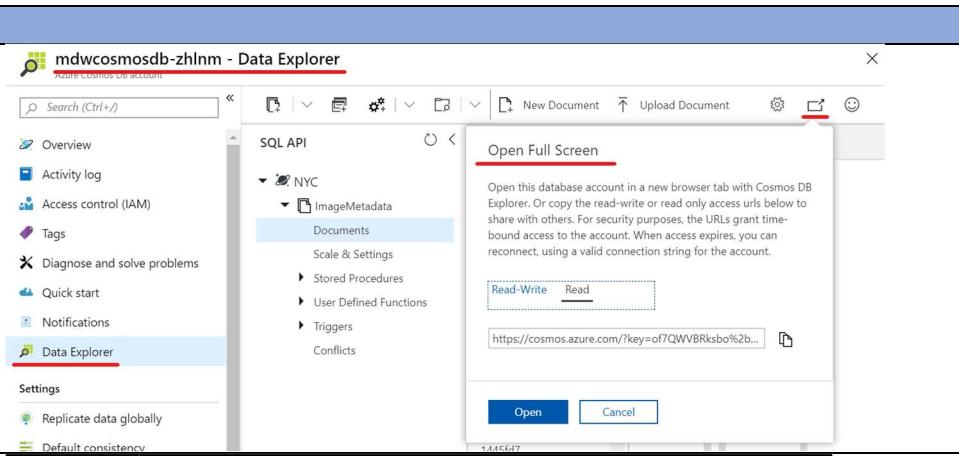


Explore Image Metadata Documents in CosmosDB

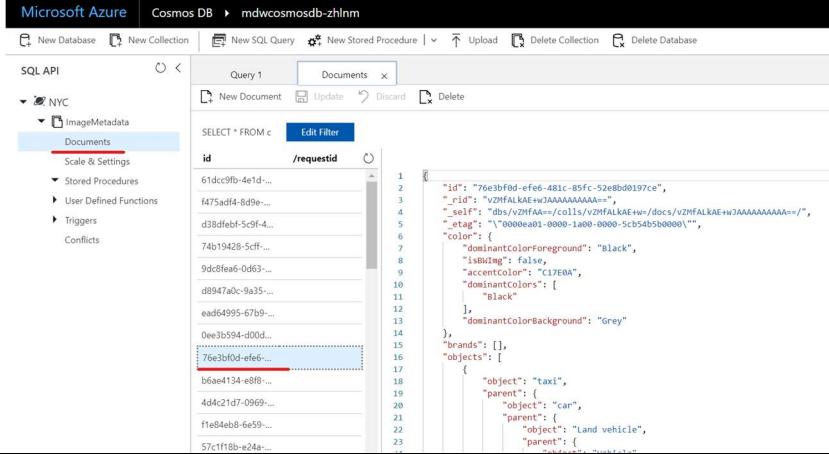
In this section you will explore the image metadata records generated by the Azure Data Factory pipeline in CosmosDB. You will use the Cosmos DB's SQL API to write SQL-like queries and retrieve data based on their criteria.

IMPORTANT: Execute these steps on your host computer

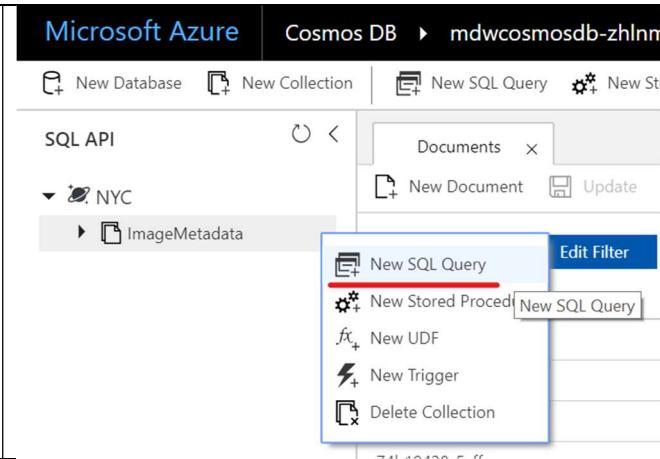
1. In the Azure Portal, go to the lab resource group and locate the CosmosDB account **MDWCosmosDB-suffix**.
2. On the **Data Explorer** panel, click **Open Full Screen** button on the top right-hand side of the screen.
3. On the Open Full Screen pop-up window, click **Open**.



4. On the **Azure Cosmos DB Data Explorer** window, under **NYC > ImageMetadata** click **Documents** to see the full list of documents in the collection.
5. Click the document id to see the content of each document.



- Click the ellipsis (...) next to ImageMetadata collection.
- On the pop-up menu, click New SQL Query to open a new query tab.



- On the New Query 1 window, try the different SQL Commands from the list.
- Click the **Execute Selection** button to execute your query.

SQL Commands

```

SELECT m.id
    , m.imageUrl
FROM ImageMetadata as m

SELECT m.id
    , m.imageUrl
    , tags.name
FROM ImageMetadata as m
    JOIN tags IN m.tags
WHERE tags.name = 'wedding'
  
```

10. Check the results in the **Results** panel.

The screenshot shows the Azure Data Studio interface. In the top right, there is a 'Query 1' tab with a close button. Below it are two buttons: 'Execute Selection' with a play icon and 'Open Query From Disk' with a disk icon. The main area contains a code editor with the following SQL query:

```
1  SELECT m.id  
2      , m.imageUrl  
3  FROM ImageMetadata as m  
4  
5  
6
```

Below the code editor, there are two tabs: 'Results' (which is selected) and 'Query Stats'. The 'Results' tab displays the output of the query. It shows a count of 1 - 94 and a JSON array of two objects. The first object has an 'id' of '61dcc9fb-4e1d-4acf-8b4f-03f7b628b6a1' and an 'imageUrl' of 'https://mdwdatalakezhlnm.blob.core.windows.net/nycimages/NYC22.jpg'. The second object has an 'id' of '547f-4de0-1f8c-2070-b270710700b'.

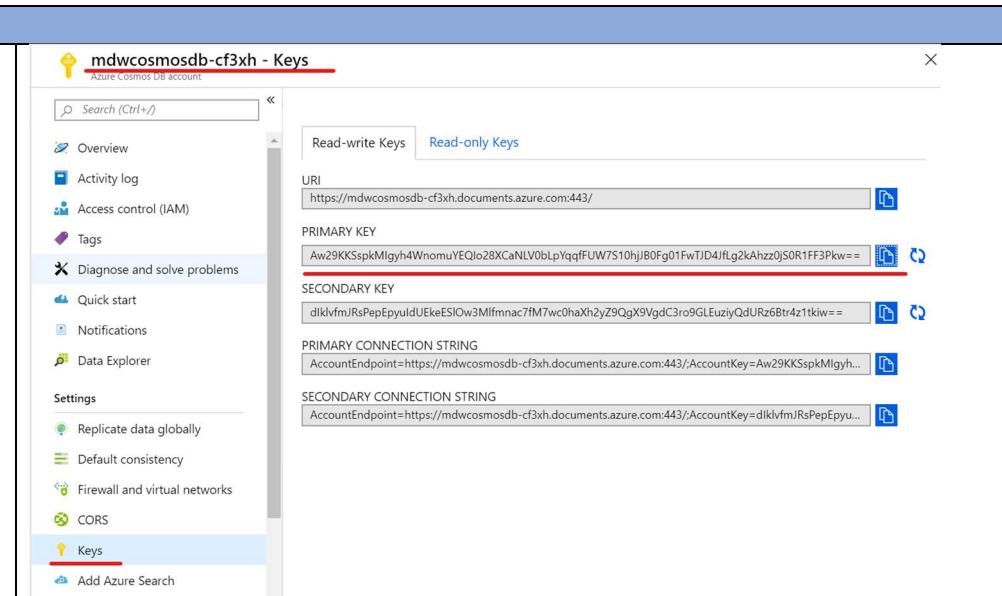
```
[  
  {  
    "id": "61dcc9fb-4e1d-4acf-8b4f-03f7b628b6a1",  
    "imageUrl": "https://mdwdatalakezhlnm.blob.core.windows.net/nycimages/NYC22.jpg"  
  },  
  {  
    "id": "547f-4de0-1f8c-2070-b270710700b"  
  }]
```

Visualize Data with Power BI

In this section you are going to use Power BI to visualize data from Cosmos DB. The Power BI report will use an Import connection to retrieve image metadata from Cosmos DB and visualise images sitting in your data lake.

IMPORTANT: Execute these steps on your host computer

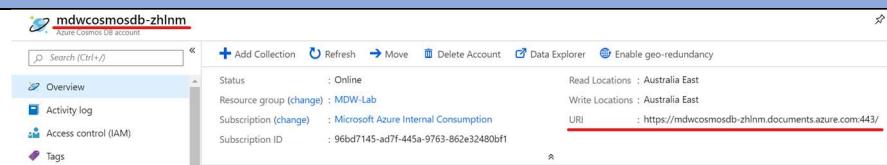
1. Navigate to the Azure Portal and retrieve the **mdwcosmosdb-suffix** access key.
2. Save it to notepad. You will need it in the next step.



The screenshot shows the 'Keys' blade of an Azure Cosmos DB account. The 'Read-write Keys' tab is active. It displays the Primary Key and Secondary Key, both of which are redacted. Below these are the Primary Connection String and Secondary Connection String, also redacted.

IMPORTANT: Steps to be executed on MDWDesktop

3. On MDWDesktop, download the Power BI report from the link <https://aka.ms/MDWLab4> and save it in the Desktop.
4. Open the file MDWLab4.pbit with Power BI Desktop.
5. When prompted to enter the value of the **MDWCosmosDB** parameter, type the full server URI: **https://mdwcosmosdb-suffix.documents.azure.com:443/**
6. Click **Load**.



The screenshot shows the 'Overview' blade of an Azure Cosmos DB account. The 'URI' field is highlighted with a red box.

MDWLab4

Modern Data Warehouse Power BI Template - Lab 4

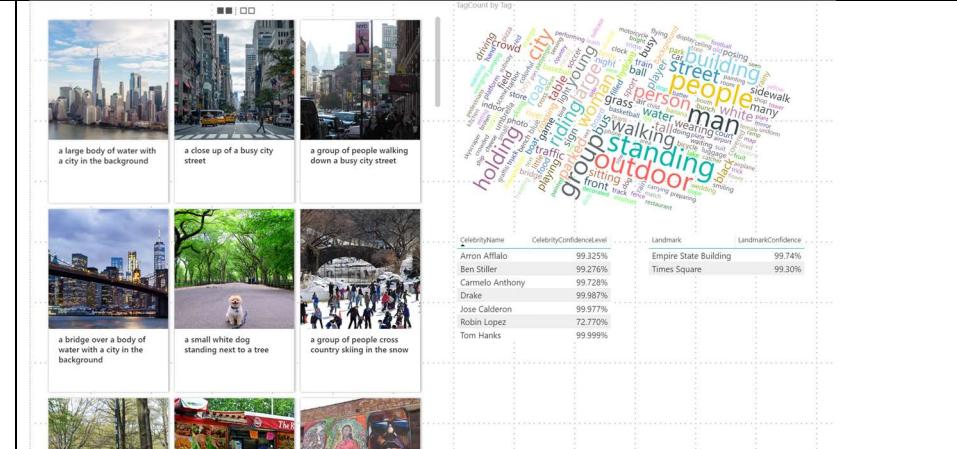
MDWCosmosDB

<https://mdwcosmosdb-zhlnm.documents.azure.com:443/>

7. When prompted for an Account Key, paste the MDWCosmosDB account key you retrieved in the previous exercise.
8. Click **Connect**.



9. Once data finish loading, interact with the report by clicking on the different images displayed and check the accuracy of their associated metadata.
10. Save your work and close Power BI Desktop.



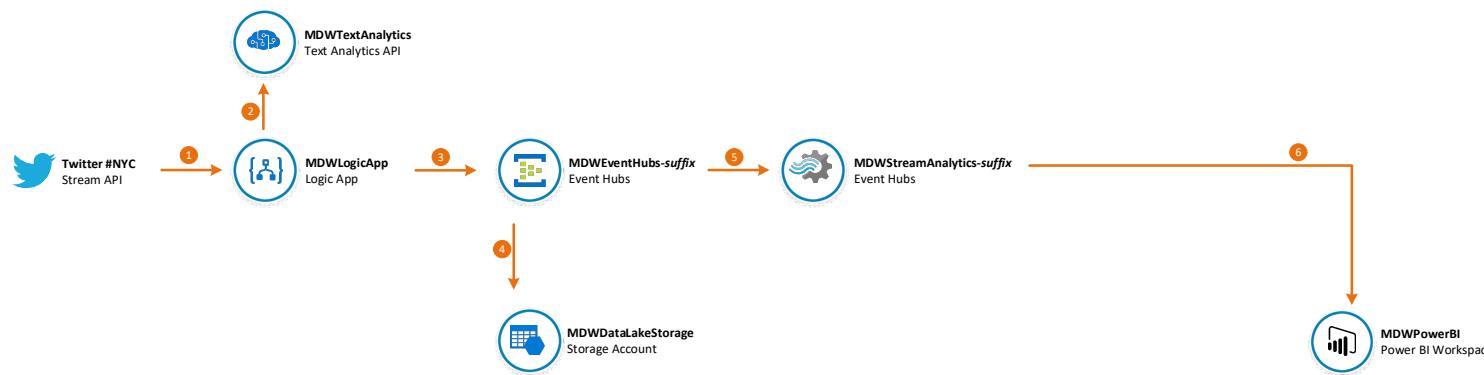
Lab 5: Ingest and Analyse real-time data with Event Hubs and Stream Analytics

In this lab you will use an Azure Logic App to connect to Twitter and generate a stream of messages using the hashtag #NYC. The logic app will invoke Azure Text Analytics Cognitive service to score Tweet sentiment and send messages to Event Hubs. You will then use Stream Analytics to generate the average tweet sentiment on the last 60 seconds and send the results to a real-time dataset in Power BI.

IMPORTANT: This lab requires you have a valid Twitter account. If you don't have one, you can sign up for free following the instructions here: <https://twitter.com/signup>.

IMPORTANT: This lab requires you have a valid Power BI account. If you don't have one you can register for a 60-day trial here: <https://powerbi.microsoft.com/en-us/power-bi-pro/>

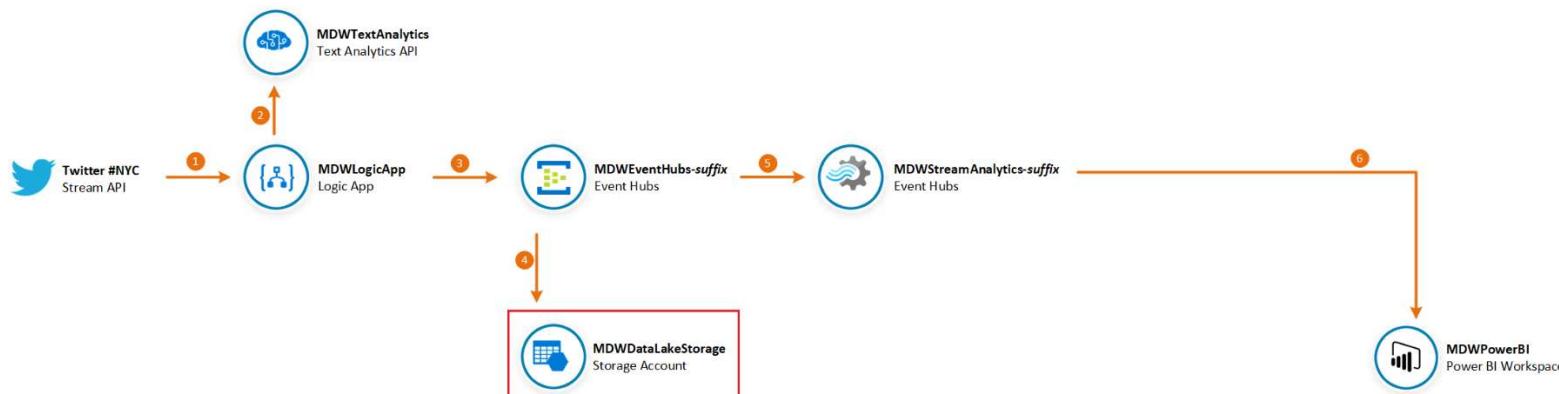
Lab Architecture



- 1 Build an Azure Logic App to invoke the Twitter API and retrieve Tweets with the hashtag #NYC
- 2 For each Tweet, invoke the Azure Text Analytics Cognitive service to detect its sentiment score
- 3 Format and send the Tweet's JSON message to Event Hubs
- 4 Save Tweet messages into your data lake for future analysis (cold path)
- 5 Send stream of Tweet messages to Stream Analytics for real-time analytics (hot path)
- 6 Visualize real-time data generated by Stream Analytics with Power BI

Create NYCTweets Container in Azure Blob Storage

In this section you will create a container in your MDWDataLake that will be used as a repository for the NYC image files. You will copy 30 files from the MDWResources Storage Account into your NYCTaxiData container.



IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, go to the lab resource group and locate the Azure Storage account **mdwdatalakesuffix**.
2. On the **Overview** panel, click **Blobs**.

Screenshot of the Azure Storage account overview page for **mdwdatalakezhlnm**:

Dashboard > MDW-Lab > mdwdatalakezhlnm

mdwdatalakezhlnm
Storage account

Overview (selected)

- Activity log
- Access control (IAM)
- Tags
- Diagnose and solve problems
- Events
- Storage Explorer (preview)

Settings

- Access keys
- Geo-replication
- CORS

Services

- Blobs**: REST-based object storage for unstructured data
[Learn more](#)

3. On the **mdwdalalakesuffix – Blobs** blade, click **+ Container**.

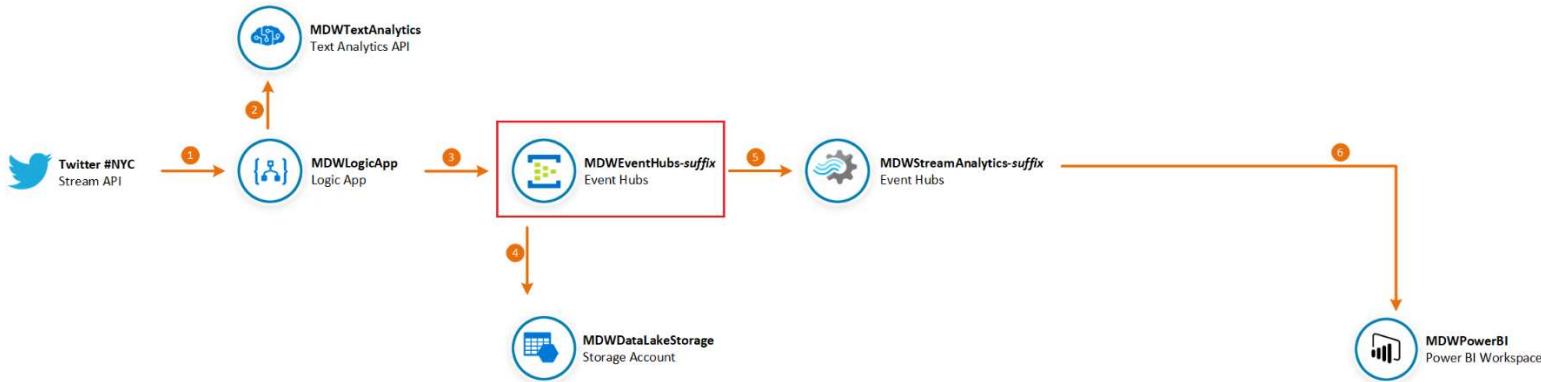
The screenshot shows the 'Blobs' blade for the storage account 'mdwdatalakezhlnm'. At the top right, there is a red box around the '+ Container' button. Below it, the storage account name is displayed. To the left is a sidebar with links: Overview, Activity log, Access control (IAM), Tags, and Diagnose and solve problems. On the right, there is a search bar labeled 'Search containers by prefix' and a table with one row containing the text 'NAME' and 'polybase'.

4. On the **New container** blade, enter the following details:
a. **Name:** nyctweets
b. **Public access level:** Private (no anonymous access)
5. Click **OK** to create the new container.

The screenshot shows the 'New container' dialog box. It has fields for 'Name' (containing 'nyctweets') and 'Public access level' (set to 'Private (no anonymous access)'). At the bottom are 'OK' and 'Cancel' buttons.

Create and Configure Event Hubs

In this section you will prepare Event Hubs to ingest Twitter data collected by the Logic App and save incoming messages to your Data Lake storage account.



IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, go to the lab resource group and locate the Event Hubs resource **MDWEVENTHUBS-suffix**.
2. On the **Event Hubs** panel, click **+ Event Hub** button to create a new event hub.

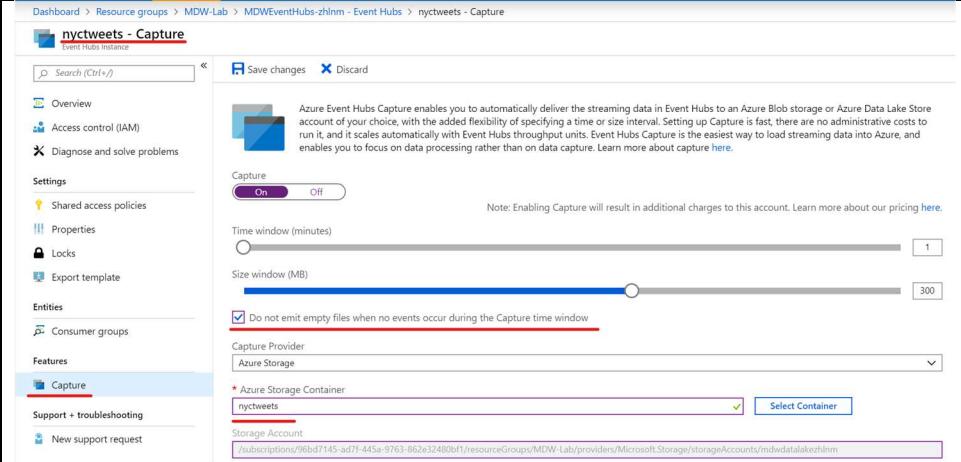
Screenshot of the Azure Event Hubs blade:

- The URL is [Dashboard > Resource groups > MDW-Lab > MDWEVENTHUBS-zhlnm - Event Hubs](#).
- The main title is **MDWEVENTHUBS-zhlnm - Event Hubs**.
- The left sidebar includes links: Overview, Access control (IAM), Tags, Diagnose and solve problems, Events, Settings (Shared access policies, Geo-Recovery, Properties, Locks, Export template), Entities, and **Event Hubs**.
- The right pane shows a table with a single row: **NAME** (no Event Hubs yet).
- Buttons at the top right include **Event Hub** (highlighted in red) and **Refresh**.

3. On the **Create Event Hub** blade, type “NYCTweets” in the **Name** field and set the **Partition Count** to 8.
4. Click **Create**.

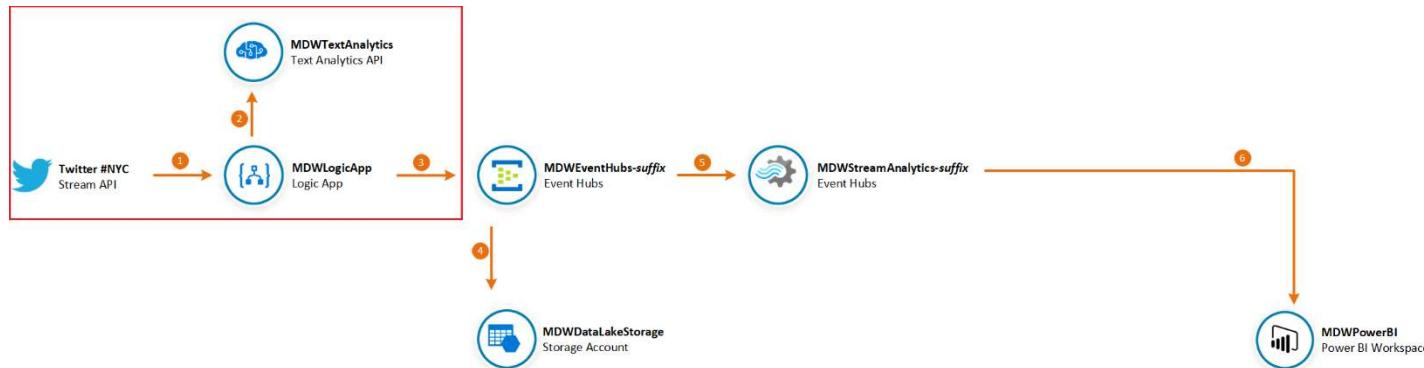


5. Once the NYCTweets event hub is created, click **Capture** on the left-hand side menu.
6. Enter the following details:
 - a. **Capture:** On
 - b. **Time window (minutes):** 1
 - c. **Do not emit empty files when no events occur during the Capture time window:** Checked.
 - d. **Capture Provider:** Azure Storage
 - e. **Azure Storage Container:** <select the nyctweets container in your MDWDataLakesuffix storage account.
7. Leave remaining fields with their default values.
8. Click **Save Changes**.



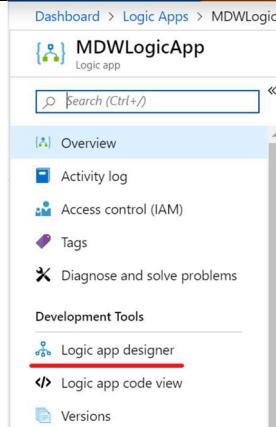
Create Azure Logic App to Read #NYC Tweets and post them to Event Hubs

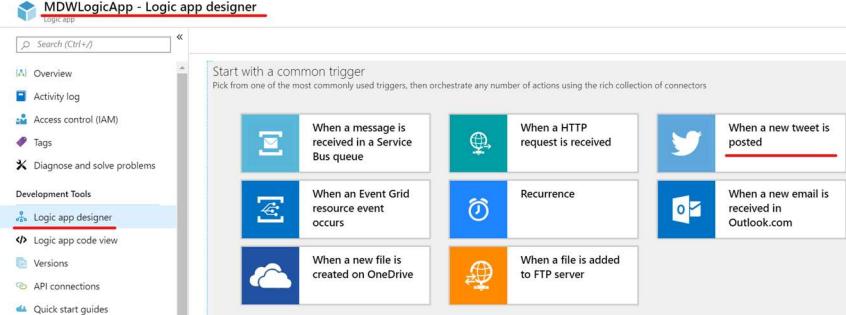
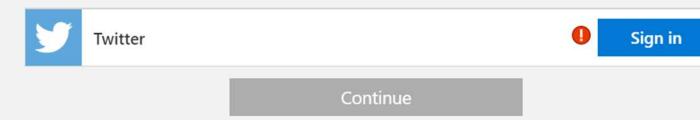
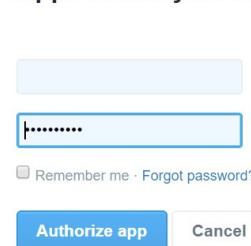
In this section you will create a Logic App to invoke the Twitter API and retrieve tweets for the hashtag #NYC. Tweets will then be formatted into a JSON message and sent to Event Hubs for processing.

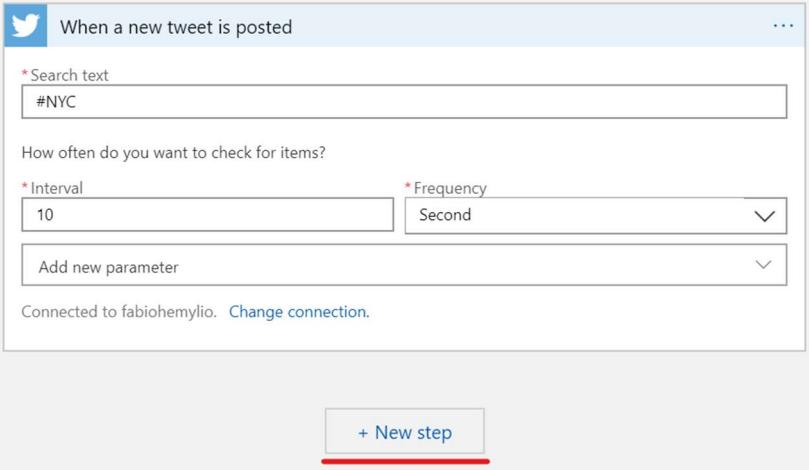
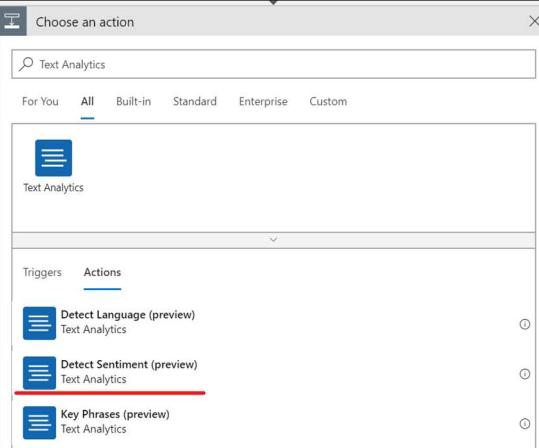


IMPORTANT: Execute these steps on your host computer

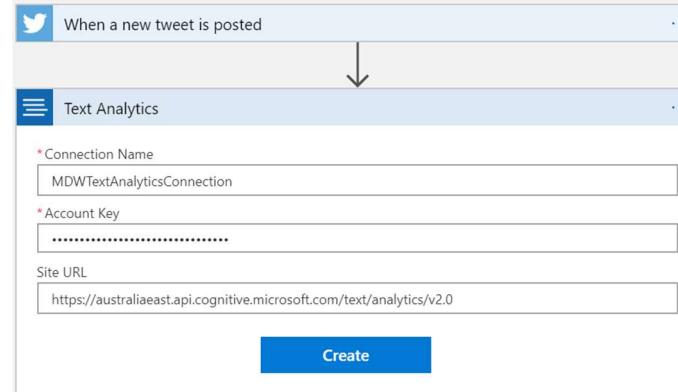
1. In the Azure Portal, go to the lab resource group and locate the Logic App resource **MDWLogicApp**.
2. On the **MDWLogicApp** menu, click **Logic app designer** to open the design blade.



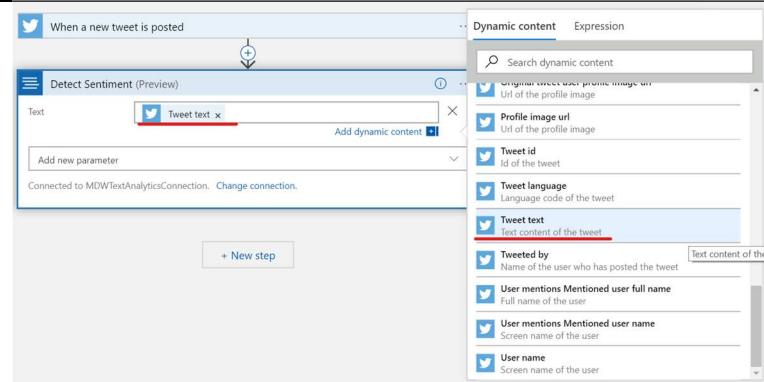
<p>3. On the Logic app designer blade, scroll down to the section Start with a common trigger.</p> <p>4. Click When a new tweet is posted.</p>	
<p>5. On the design surface you will see the Twitter connector. Click Sign in and use your personal Twitter account to authenticate.</p>	
<p>6. On the log on screen review the permissions that will be granted to logic apps. If you agree with the permissions granted then enter your credentials and click Authorise app.</p> <p>7. On the Confirmation Required page, click Allow access to proceed with the authentication process and to grant the right permissions to Logic Apps.</p>	<p>Authorize Microsoft Azure Logic Apps to use your account?</p>  <p>This application will be able to:</p> <ul style="list-style-type: none"> • Read Tweets from your timeline. • See who you follow, and follow new people. • Update your profile. • Post Tweets for you. • Access your direct messages. <p>Will not be able to:</p> <ul style="list-style-type: none"> • See your email address. • See your Twitter password.

<p>8. Once authenticated with Twitter you will notice the green tick next to your user name.</p> <p>9. Click Continue to configure the Twitter connector.</p>	<p>This logic app will connect to:</p> 
<p>10. On the When a new tweet is posted properties, enter the follow details:</p> <ul style="list-style-type: none"> a. Search text: #NYC b. Interval: 10 c. Frequency: Seconds <p>11. Leave remaining fields with their default values.</p> <p>12. Click + New step to create a new Logic App task.</p>	
<p>13. On the Choose an action box, type “Text Analytics” in the search field. Select Detect Sentiment (preview) from the Actions tab.</p>	

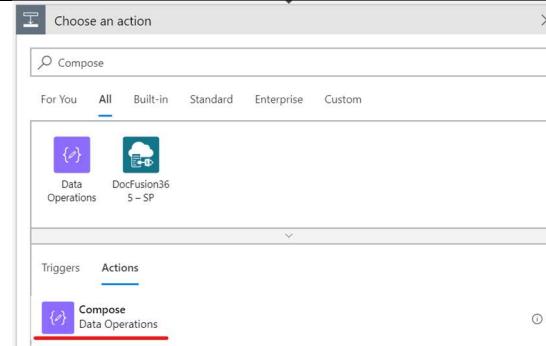
14. On the Text Analytics properties, type “MDWTextAnalyticsConnection” in the Connection Name field.
15. Open the **Azure Portal** in a new browser tab and copy Endpoint and Key values for **MDWTextAnalytics**.
16. Paste the values in the respective fields in the **Text Analytics** properties.
17. Click **Create**.



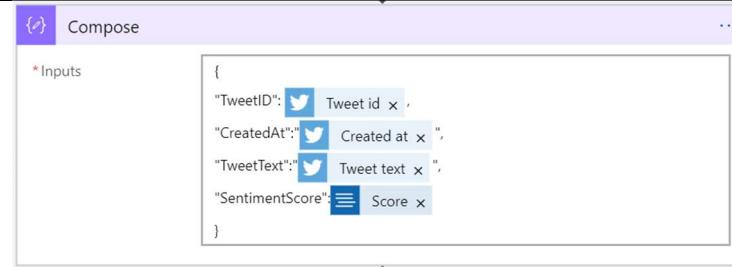
18. In the **Detect Sentiment (preview)** properties, click **Add new parameter**.
19. Select **Text** check box.
20. In the **Text** parameter box, select **Tweet Text** from the **Dynamic content** list.
21. Click **+ New step** to create a new task.



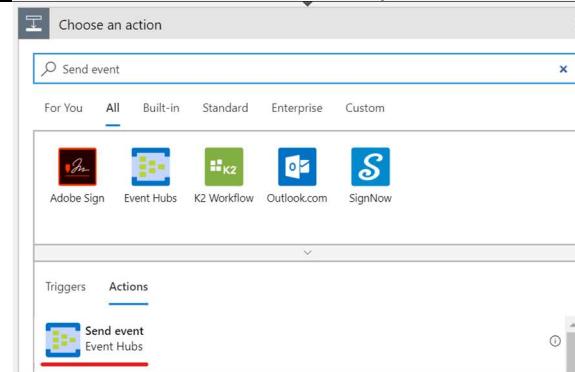
22. Type “Compose” in the search box and select the **Compose Data Operation**.



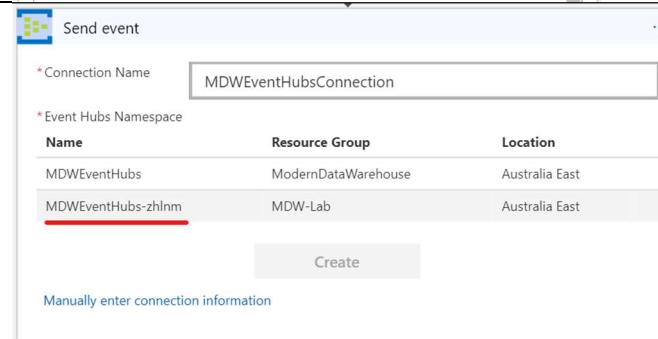
23. On the Compose properties, build a new JSON message using data elements returned by the previous tasks. Your JSON message should look like this.
24. Click **+ New step** to create a new task.

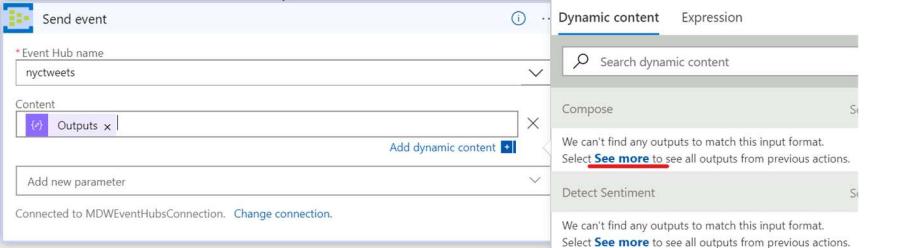
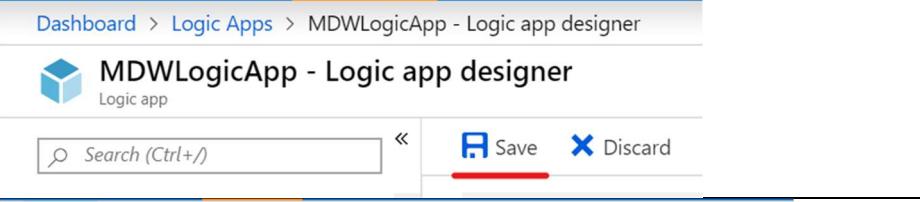
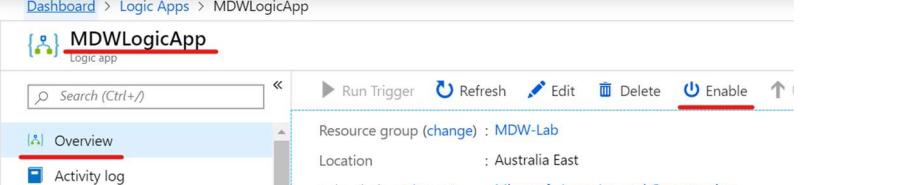
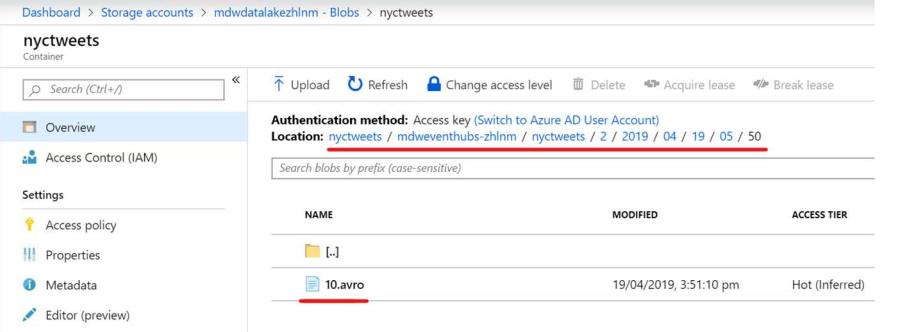


25. On the **Choose an action** box, type “Send event” in the search field. Select **Send Event** from the **Actions** tab.



26. On the **Send event** properties, type “MDWEVENTHUBSConnection” in the **Connection Name** field.
 27. Select **MDWEVENTHUBS-suffix** from the list of **Event Hubs Namespaces**.
 28. Select the default access policy **RootManageSharedAccessKey**.
 29. Click **Create**.



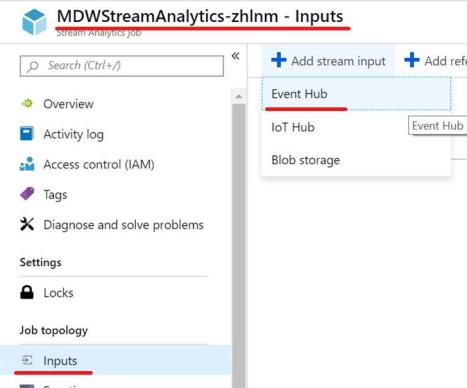
<p>30. On the Send event properties, select nyctweets in the Event Hub name field.</p> <p>31. Click Add new parameter and select Content.</p> <p>32. Click on the Content field. From the Dynamic content pop-up window, select click the See more link under Compose to display the Outputs field.</p> <p>33. Select Outputs as the value for the Content field.</p>							
<p>34. Click the Save button to save your Logic App.</p>							
<p>35. On the Overview panel, click Enable.</p>							
<p>36. In the Azure Portal, navigate to the MDWDataLakesuffix storage account.</p> <p>37. Wait a couple of minutes and you should be able to see new files being created in the nyctweets container.</p>	 <table border="1"> <thead> <tr> <th>NAME</th> <th>MODIFIED</th> <th>ACCESS TIER</th> </tr> </thead> <tbody> <tr> <td>10.avro</td> <td>19/04/2019, 3:51:10 pm</td> <td>Hot (Inferred)</td> </tr> </tbody> </table>	NAME	MODIFIED	ACCESS TIER	10.avro	19/04/2019, 3:51:10 pm	Hot (Inferred)
NAME	MODIFIED	ACCESS TIER					
10.avro	19/04/2019, 3:51:10 pm	Hot (Inferred)					

Create and Configure Stream Analytics

In this section you will configure Stream Analytics to perform analytic queries on streaming data sent by Event Hubs and generate outputs to Power BI.

IMPORTANT: Execute these steps on your host computer

1. In the Azure Portal, go to the lab resource group and locate the Event Hubs resource **MDWStreamAnalytics-suffix**.
2. On the **Inputs** panel, click **+ Add stream input** button and select **Event Hub** to create a new stream input.

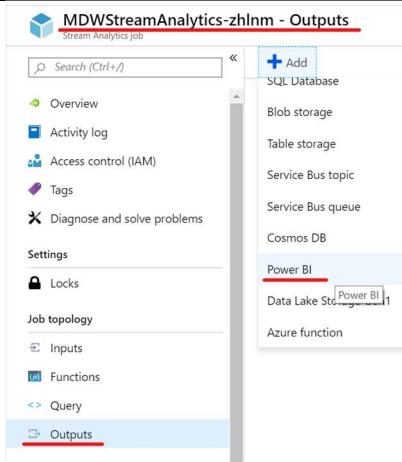


3. On the **Event Hub New input** blade enter the following details:
 - a. **Input alias:** MDWEVENTHUBS
 - b. **Event Hub namespace:** MDWEVENTHUBS-suffix
 - c. **Event hub name > Use existing:** nyctweets
4. Leave remaining fields with their default values.

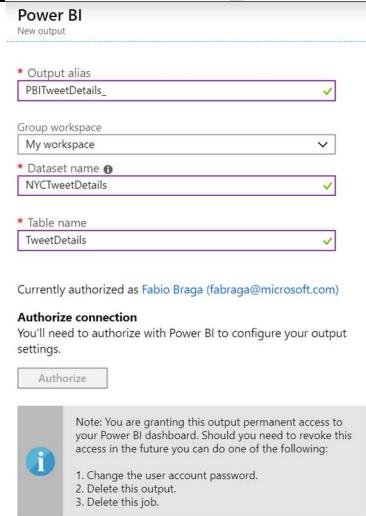
The screenshot shows the 'Event Hub New input' blade with the following configuration:

- Input alias:** MDWEVENTHUBS
- Provide Event Hub settings manually:** Unselected
- Select Event Hub from your subscriptions:** Selected
- Subscription:** Microsoft Azure Internal Consumption
- Event Hub namespace:** MDWEVENTHUBS-zhlnm
- Event Hub name:** nyctweets (radio button selected for 'Use existing')
- Event Hub policy name:** RootManageSharedAccessKey
- Event Hub policy key:** (redacted)
- Event Hub consumer group:** (redacted)
- Event serialization format:** JSON

- On the **Inputs** panel, click **+ Add** button and select **Power BI** to create a new stream output.



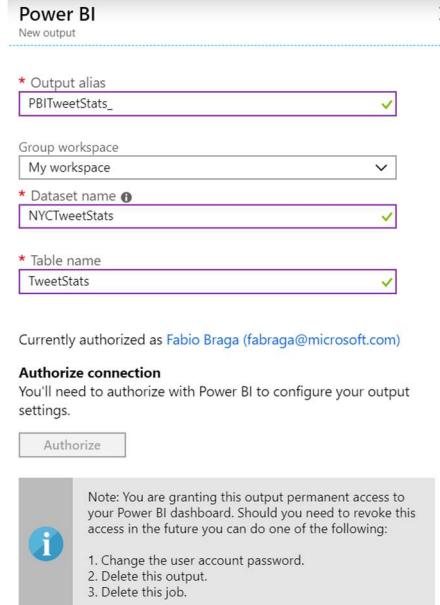
- On the **Power BI New Output** blade, click **Authorize** to authenticate with Power BI.
- Once authenticated, enter the following details:
 - Output alias:** PBITweetDetails
 - Group Workspace:** My Workspace
 - Dataset name:** NYCTweetDetails
 - Table name:** TweetStats
- Leave remaining fields with their default values.
- Click **Save**.



10. Repeat the process to create another Power BI Output. This time enter the following details:

- a. **Output alias:** PBITweetStats
- b. **Group Workspace:** My Workspace
- c. **Dataset name:** NYCTweetStats
- d. **Table name:** TweetStats

11. Click **Save**.



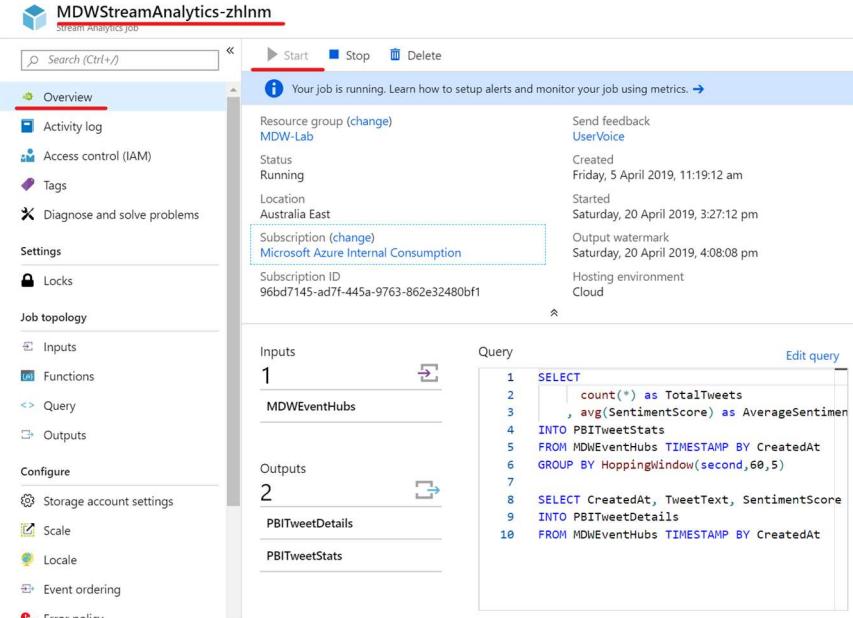
12. On the **Query** panel, note the inputs and outputs you created in the previous steps.

```
1 SELECT
2 | count(*) as TotalTweets
3 , avg(SentimentScore) as AverageSentiment
4 INTO PBITweetStats
5 FROM MDWEventHubs TIMESTAMP BY CreatedAt
6 GROUP BY HoppingWindow(second,60,5)
7 |
8 SELECT CreatedAt, TweetText, SentimentScore
9 INTO PBITweetDetails
10 FROM MDWEventHubs TIMESTAMP BY CreatedAt
```

13. Enter the following SQL commands in the query window.

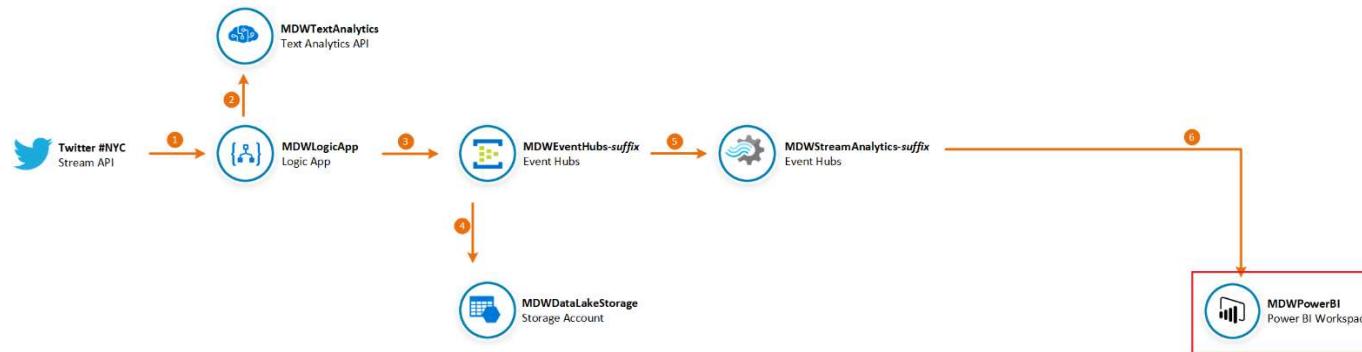
14. Click **Save**.

```
SELECT
    count(*) as TotalTweets
    , avg(SentimentScore) as AverageSentiment
INTO PBITweetStats
```

	<pre> FROM MDWEVENTHUBS TIMESTAMP BY CreatedAt GROUP BY HoppingWindow(second,60,5) SELECT CreatedAt, TweetText, SentimentScore INTO PBITweetDetails FROM MDWEVENTHUBS TIMESTAMP BY CreatedAt </pre>
15. On the Overview panel, click Start to start the Stream Analytics job.	 <p>The screenshot shows the Azure Stream Analytics Overview page for a job named "MDWStreamAnalytics-zhlnm". The job is running and was created on Friday, 5 April 2019, at 11:19:12 am. It was started on Saturday, 20 April 2019, at 3:27:12 pm. The output watermark is Saturday, 20 April 2019, at 4:08:08 pm, and it's hosted in the Cloud environment. The job has one input, "MDWEVENTHUBS", and two outputs: "PBITweetDetails" and "PBITweetStats". The "Query" section shows the T-SQL code for both outputs:</p> <pre> 1 SELECT 2 count(*) as TotalTweets 3 , avg(SentimentScore) as AverageSentiment 4 INTO PBITweetStats 5 FROM MDWEVENTHUBS TIMESTAMP BY CreatedAt 6 GROUP BY HoppingWindow(second,60,5) 7 8 SELECT CreatedAt, TweetText, SentimentScore 9 INTO PBITweetDetails 10 FROM MDWEVENTHUBS TIMESTAMP BY CreatedAt </pre>

Create Power BI Dashboard to Visualise Real-Time Data

In this section you will log on to the Power BI portal to create a dashboard to visualize real-time Tweeter statistics data sent by Stream Analytics.

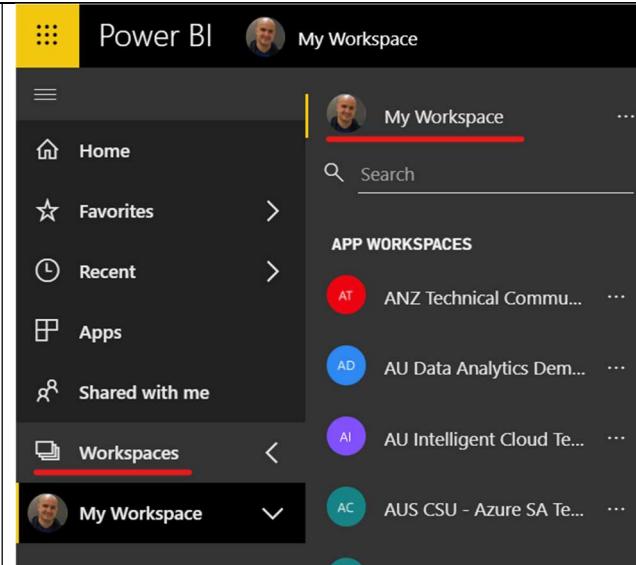


IMPORTANT: Steps to be executed on the Student's Computer

1. Open a new browser tab and navigate to <https://www.powerbi.com>
2. Enter your credentials to authenticate with the Power BI service.

A screenshot of the Microsoft Power BI website homepage. The top navigation bar includes the Microsoft logo, the "Power BI" text, and links for "Sign in" (underlined), "Sign up free", and a menu icon. The main banner features the text "Business intelligence like never before" and "Go from data to insights in minutes. Any data, any way, anywhere. And all in one view." Below the banner is a yellow button labeled "START FREE >".

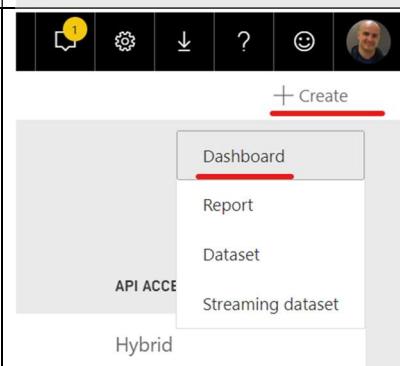
3. Once authenticated, open the **Workspaces** menu and click **My Workspace** at the top of the **Workspaces** list.



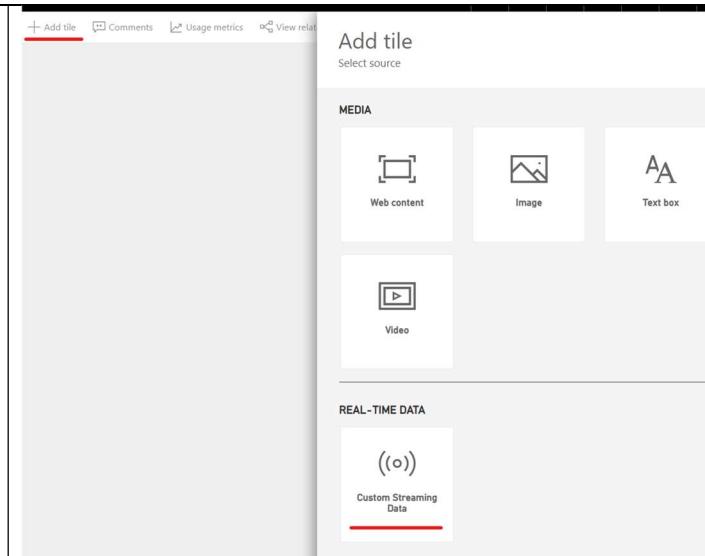
4. Navigate to the **Datasets** tab and verify that two datasets have been created by Stream Analytics: **NYCTweetDetails** and **NYCTweetStats**.

A screenshot of the Power BI Datasets tab. At the top, there is a search bar labeled 'Search content...'. Below the search bar, there are tabs for 'Dashboards', 'Reports', 'Workbooks', and 'Datasets', with 'Datasets' being the active tab. The main area displays a list of datasets with columns for 'NAME' and 'ACTIONS'. Two datasets are listed: 'NYCTweetDetails' and 'NYCTweetStats'. Each dataset row has a small icon, a name, and a set of icons for managing the dataset.

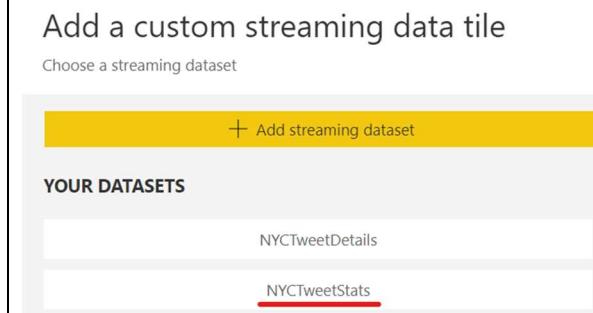
5. On the top right-hand side corner click **+ Create** and then click **Dashboard** from the dropdown menu to create a new dashboard.
 6. Type “NYC Tweet Stats” in the **Dashboard name** field and click **Create**.



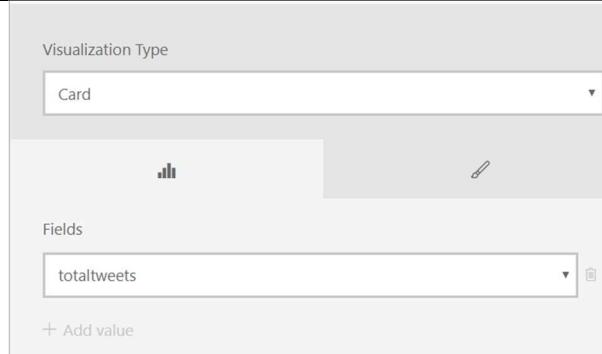
7. Click **+ Add tile** button on the toolbar.
8. On the **Add tile** blade, select **Custom Streaming Data** under the **Real-Time Data** section.
9. Click **Next**.



10. On the Add a custom streaming data tile blade, select the **NYCTweetStats** dataset.
11. Click **Next**.



12. On the **Visualization Type** field select **Card**.
13. On the **Fields** field select **totaltweets**.
14. Click **Next**.



15. On the Tile details blade, enter the following details:
a. **Title:** Total Tweets
b. **Subtitle:** in the last minute.
16. Leave remaining fields with their default values.
17. Click **Apply**.

Tile details

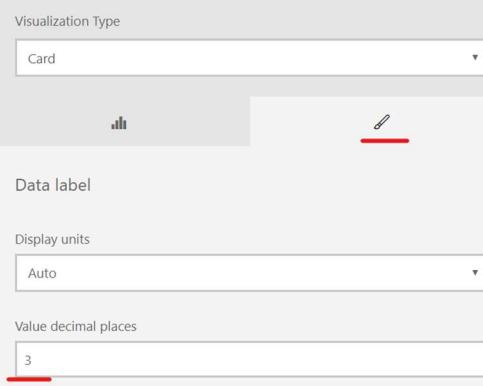
The screenshot shows the "Tile details" blade. It includes a section for "Required Details" with a checked checkbox for "Display title and subtitle". Below this, there are two input fields: "Title" containing "Total Tweets" and "Subtitle" containing "in the last minute".

18. Repeat the process to create another tile, this time to display the average sentiment score. Use the following details:

- a. **Dataset:** NYCTweetStats
- b. **Visualization Type:** Card
- c. **Fields:** averagesentiment
- d. **Data Label > Value decimal places:** 3
- e. **Details > Display Title and Subtitle:** Checked
- f. **Details > Title:** Average Sentiment Score
- g. **Details > Subtitle:** in the last minute

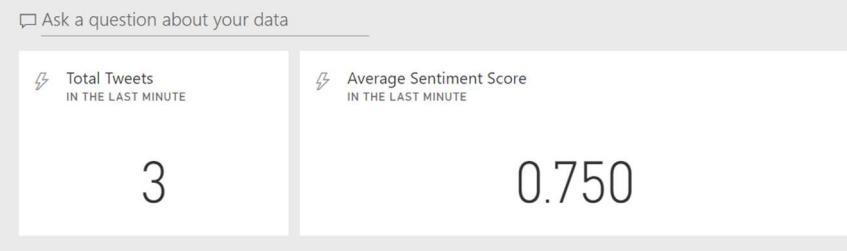
Add a custom streaming data tile

Choose a streaming dataset > Visualization design



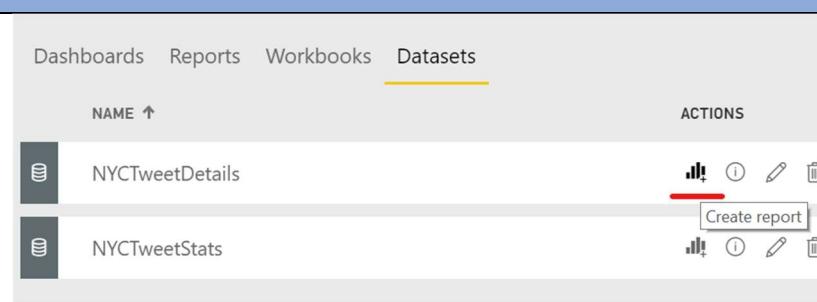
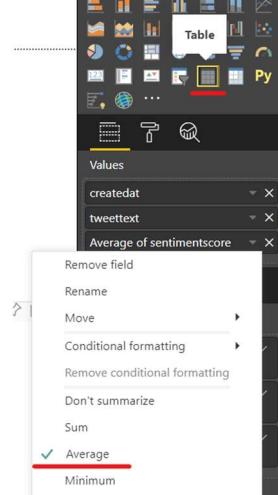
19. You should be able to see the values for both tiles changing every few seconds.

20. Try logging on to Tweeter and posting different messages with the hashtag **#NYC**. Go back to your Power BI dashboard and see the effect your messages had on the average sentiment score.

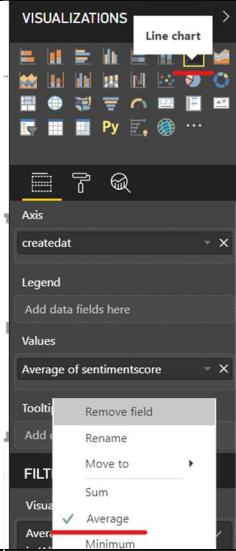


Create Power BI Report to Visualise Real-Time Data

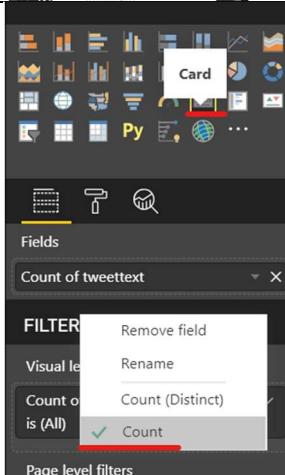
In this section you will log on to the Power BI portal to create a dashboard to visualize real-time Tweeter statistics data sent by Stream Analytics.

IMPORTANT: Steps to be executed on the Student's Computer	
1. On the Power BI portal navigate to My Workspace and click on the Datasets tab again. 2. Click the Create Report icon under Actions for the NYCTweetDetails dataset.	
3. On the report canvas add a Table visual and add the following dataset fields to the Values property: a. createdat b. tweettext c. sentimentscore 4. Change the default aggregation function for sentimentscore to Average .	

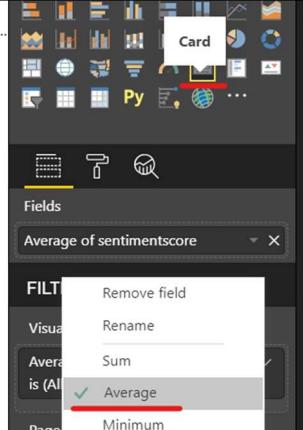
5. On the report canvas add a **Line chart** visual and add the following dataset fields:
 - a. Axis: createdat
 - b. Values: sentimentscore
6. Change the default aggregation function for **sentimentscore** to **Average**.



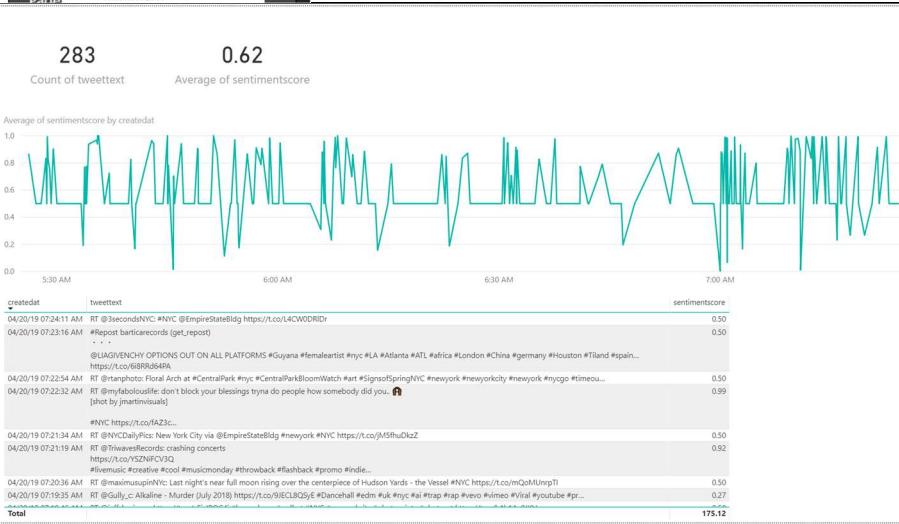
7. On the report canvas add a **Card** visual and add the following dataset fields to its **Fields** property:
 - a. tweettext
8. Change the default aggregation function for **tweettext** to **Count**.



9. On the report canvas add a **Card** visual and add the following dataset fields to its **Fields** property:
 - b. sentimentscore
10. Change the default aggregation function for **sentimentscore** to **Average**.



11. Place your visuals in a way you can see all of them.
12. Click the **Save** button on the toolbar to save your report.
13. Type “NYC Tweet Details” in the **Report name** field.



References

© 2015 Microsoft Corporation. All rights reserved.

By using this demo/lab, you agree to the following terms:

The technology/functionality described in this demo/lab is provided by Microsoft Corporation for purposes of obtaining your feedback and to provide you with a learning experience. You may only use the demo/lab to evaluate such technology features and functionality and provide feedback to Microsoft. You may not use it for any other purpose. You may not modify, copy, distribute, transmit, display, perform, reproduce, publish, license, create derivative works from, transfer, or sell this demo/lab or any portion thereof.

COPYING OR REPRODUCTION OF THE DEMO/LAB (OR ANY PORTION OF IT) TO ANY OTHER SERVER OR LOCATION FOR FURTHER REPRODUCTION OR REDISTRIBUTION IS EXPRESSLY PROHIBITED.

THIS DEMO/LAB PROVIDES CERTAIN SOFTWARE TECHNOLOGY/PRODUCT FEATURES AND FUNCTIONALITY, INCLUDING POTENTIAL NEW FEATURES AND CONCEPTS, IN A SIMULATED ENVIRONMENT WITHOUT COMPLEX SET-UP OR INSTALLATION FOR THE PURPOSE DESCRIBED ABOVE. THE TECHNOLOGY/CONCEPTS REPRESENTED IN THIS DEMO/LAB MAY NOT REPRESENT FULL FEATURE FUNCTIONALITY AND MAY NOT WORK THE WAY A FINAL VERSION MAY WORK. WE ALSO MAY NOT RELEASE A FINAL VERSION OF SUCH FEATURES OR CONCEPTS. YOUR EXPERIENCE WITH USING SUCH FEATURES AND FUNCITONALITY IN A PHYSICAL ENVIRONMENT MAY ALSO BE DIFFERENT.

Feedback

If you give feedback about the technology features, functionality and/or concepts described in this demo/lab to Microsoft, you give to Microsoft, without charge, the right to use, share and commercialize your feedback in any way and for any purpose. You also give to third parties, without charge, any patent rights needed for their products, technologies and services to use or interface with any specific parts of a Microsoft software or service that includes the feedback. You will not give feedback that is subject to a license that requires Microsoft to license its software or documentation to third parties because we include your feedback in them. These rights survive this agreement.

MICROSOFT CORPORATION HEREBY DISCLAIMS ALL WARRANTIES AND CONDITIONS WITH REGARD TO THE DEMO/LAB, INCLUDING ALL WARRANTIES AND CONDITIONS OF MERCHANTABILITY, WHETHER EXPRESS, IMPLIED OR STATUTORY, FITNESS FOR A PARTICULAR PURPOSE, TITLE AND NON-INFRINGEMENT. MICROSOFT DOES NOT MAKE ANY ASSURANCES OR REPRESENTATIONS WITH REGARD TO THE ACCURACY OF THE RESULTS, OUTPUT THAT DERIVES FROM USE OF DEMO/ LAB, OR SUITABILITY OF THE INFORMATION CONTAINED IN THE DEMO/LAB FOR ANY PURPOSE.

Disclaimer

This demo/lab contains only a portion of new features and enhancements in Microsoft Azure data services. Some of the features might change in future releases of the product. In this demo/lab, you will learn about some, but not all, new features.