

Chapter 1: Creating an Azure Databricks Service

[Home](#) > [Create a resource](#) >

Azure Databricks ⚙ ...

Microsoft



Azure Databricks ♡ Add to Favorites

Microsoft

★★★★★ 4.3 (156 ratings)

[Create](#)

[Overview](#) [Plans](#) [Usage Information + Support](#) [Reviews](#)

Fast, easy, and collaborative Apache Spark-based analytics platform

Accelerate innovation by enabling data science with a high-performance analytics platform that's optimized for Azure.

Drive innovation and increase productivity

Create an Azure Databricks workspace

Basics Networking Advanced Tags Review + create

Project Details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * ⓘ

Resource group * ⓘ

[Create new](#)

Instance Details

Workspace name *

Region *

Pricing Tier * ⓘ

A resource group is a container that holds related resources for an Azure solution.

Name *

CookbookRG 

[OK](#)

[Cancel](#)

[Review + create](#)

[< Previous](#)

[Next : Networking >](#)

Create an Azure Databricks workspace

...

Basics Networking Advanced Tags Review + create

Project Details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * ⓘ

Resource group * ⓘ

 ⓘ

[Create new](#)

Instance Details

Workspace name *

 ⓘ

Region *

 ⓘ

Pricing Tier * ⓘ

 ⓘ

Standard (Apache Spark, Secure with Azure AD)

Premium (+ Role-based access controls)

Trial (Premium - 14-Days Free DBUs)

[Review + create](#)

[< Previous](#)

[Home](#) > [Create a resource](#) > [Azure Databricks](#) >

Create an Azure Databricks workspace

...

Basics

Networking

Advanced

Tags

Review + create

Deploy Azure Databricks workspace with Secure Cluster Connectivity (No Public IP) ⓘ

Yes No

Deploy Azure Databricks workspace in your own Virtual Network (VNet)

Yes No

[Review + create](#)

[< Previous](#)

[Next : Advanced >](#)

Create an Azure Databricks workspace

...

 Validation Succeeded

Basics Networking Advanced Tags **Review + create**

Summary

Basics

Workspace name	BigDataWorkspace
Subscription	[REDACTED]
Resource group	CookbookRG
Region	East US
Pricing Tier	trial

Networking

Deploy Azure Databricks workspace with Secure Cluster Connectivity (No Public IP)	No
Deploy Azure Databricks workspace in your own Virtual Network (VNet)	No

Create

< Previous

Download a template for automation

BigDataWorkspace Azure Databricks Service

Search (Ctrl+ /) Delete

Overview

- Activity log
- Access control (IAM)
- Tags

Settings

- Virtual Network Peerings
- Encryption
- Properties
- Locks

Automation

- Tasks (preview)
- Export template

Support + troubleshooting

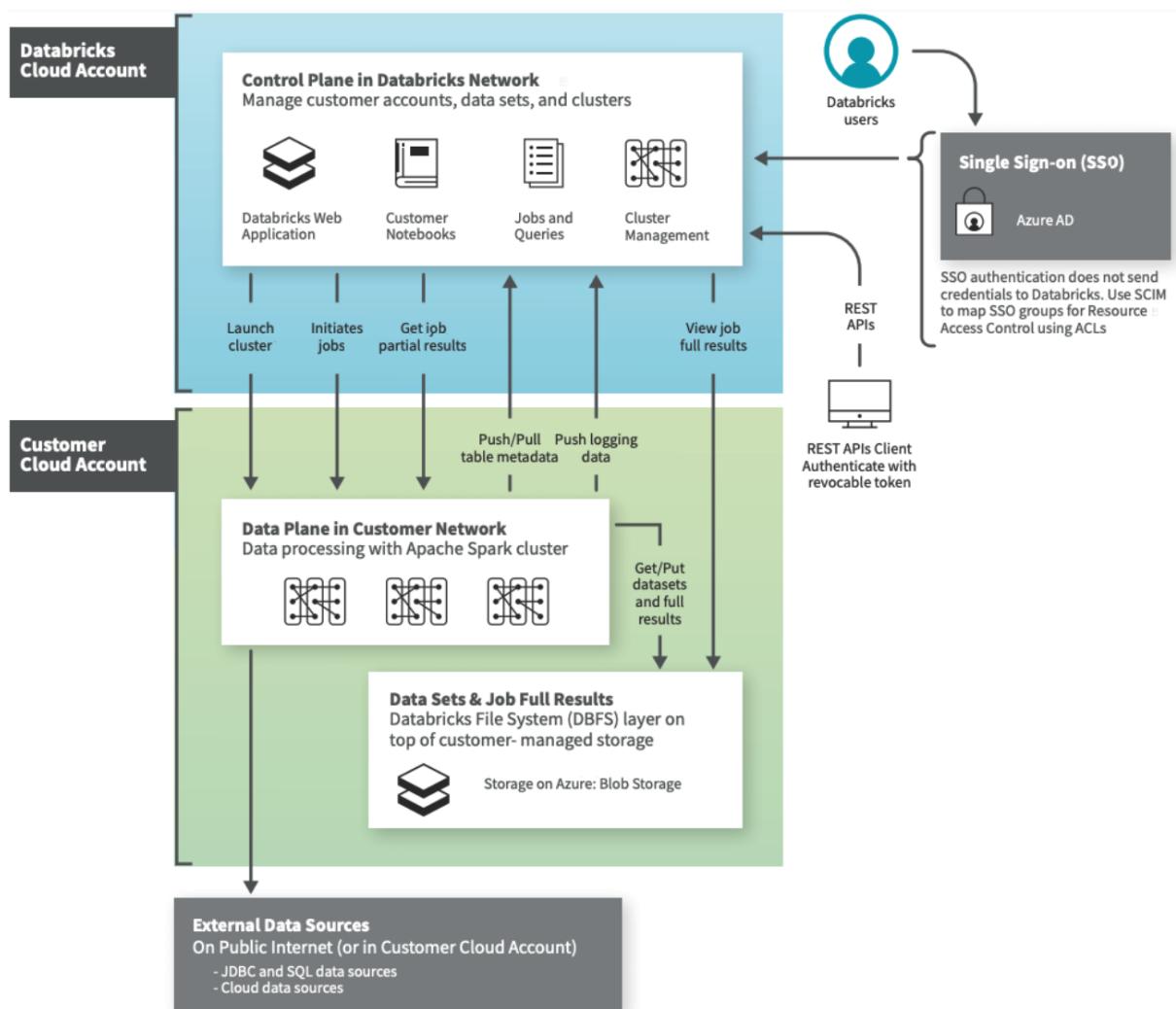
- New support request

Status: Active
Resource group: CookbookRG
Location: East US
Subscription ID: [REDACTED]
Tags (change): Click here to add tags

Managed Resource Group: databricks-rg-BigDataWorkspace-g3eexhrn3ipw
URL: https://adb-[REDACTED].azuredatabricks.net
Pricing Tier: Trial (Premium - 14-Days Free DBUs)



Launch Workspace



databricks-rg-BigDataWorkspace-g3eexhzrn3ipw

Resource group

Search (Ctrl+ /) Create Edit columns Delete resource group Refresh Export to CSV Open query Assign tag

Overview Activity log Access control (IAM) Tags Events

Subscription (change) : Deployments : 2 Success

Subscription ID : Location : East US

Tags (change) : Click here to add tags

Filter for any field... Type == all X Location == all X Add filter

Showing 1 to 3 of 3 records. Show hidden types

Name	Type
dbstoragefnppxdzsotio	Storage account
workers-sg	Network security group
workers-vnet	Virtual network

Deployments Security Policies Properties



Microsoft Azure (Preview)

PowerShell Bash

Cloud Shell. Succeeded.

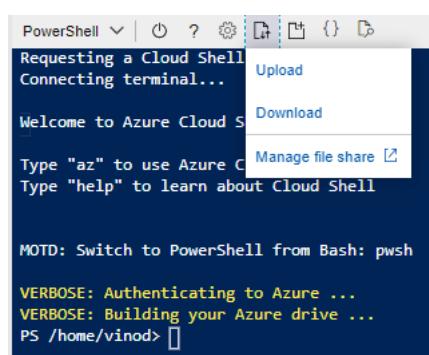
Welcome to Azure Cloud Shell

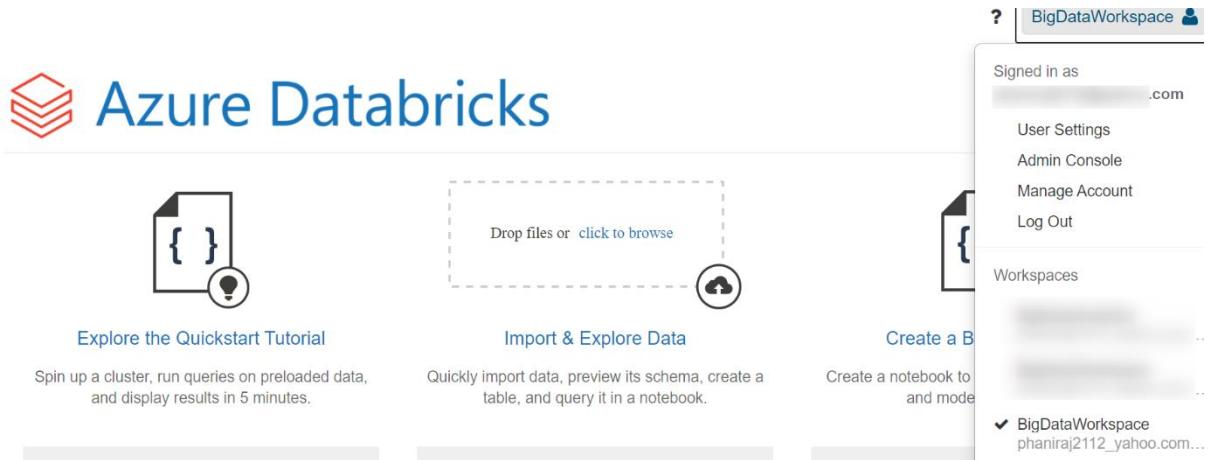
Type "az" to use Azure CLI
Type "help" to learn about Cloud Shell

Your Cloud Shell session will be ephemeral so no files or system changes will persist.

MOTD: Download files or directories from the cloudshell: Export-File

VERBOSE: Authenticating to Azure ...
VERBOSE: Building your Azure drive ...
PS /home/vinod> []



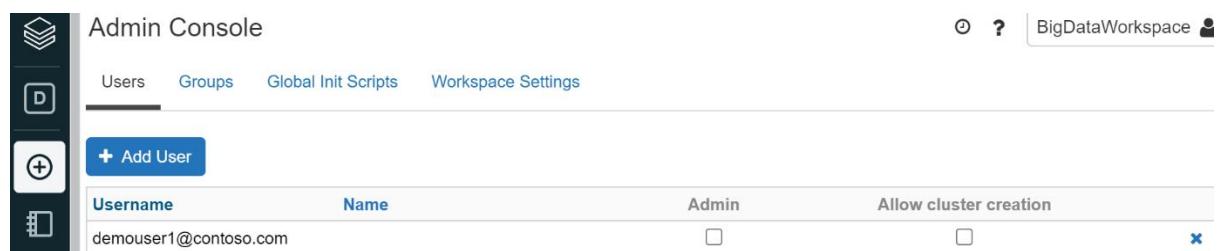


The screenshot shows the Azure Databricks home page. At the top right, there is a user menu with options like "Signed in as", "User Settings", "Admin Console", "Manage Account", and "Log Out". Below the menu, there are three main sections: "Explore the Quickstart Tutorial", "Import & Explore Data", and "Create a Notebook". Each section has a brief description and a corresponding icon.

Explore the Quickstart Tutorial
Spin up a cluster, run queries on preloaded data, and display results in 5 minutes.

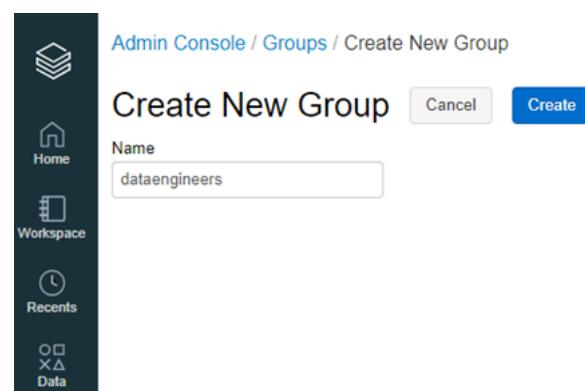
Import & Explore Data
Quickly import data, preview its schema, create a table, and query it in a notebook.

Create a Notebook
Create a notebook to and mode

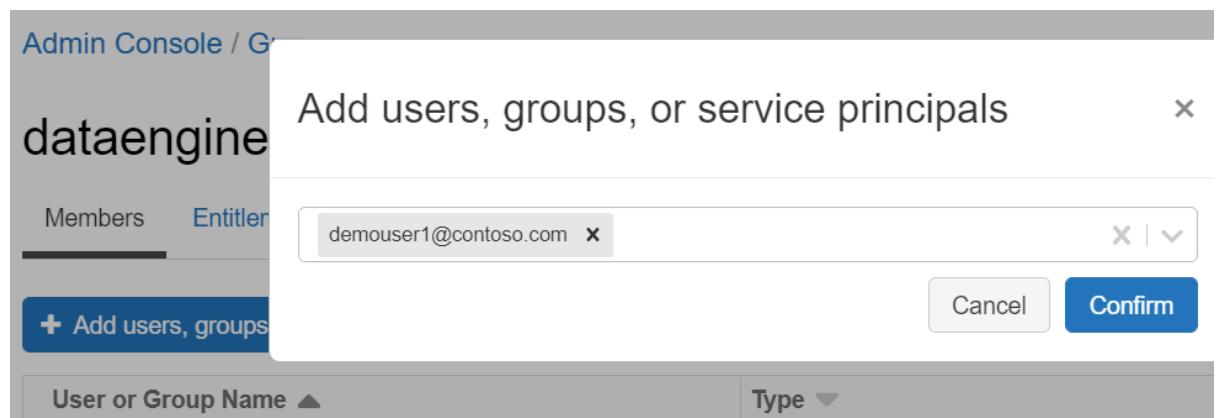


The screenshot shows the "Admin Console" interface, specifically the "Users" tab. It includes a sidebar with icons for Home, Workspace, Recents, and Data. The main area displays a table with one user entry:

Username	Name	Admin	Allow cluster creation
demouser1@contoso.com		<input type="checkbox"/>	<input type="checkbox"/>



The screenshot shows a modal dialog titled "Create New Group" with a "Cancel" button and a "Create" button. On the left is a sidebar with Home, Workspace, Recents, and Data icons. The main area has a "Name" field containing "dataengineers".



The screenshot shows a modal dialog titled "Add users, groups, or service principals" with a "Cancel" button and a "Confirm" button. On the left is a sidebar with Members and Entitlements tabs, and a "dataengineers" group selected. The main area has a search bar with "demouser1@contoso.com" and a "User or Group Name" input field.

Admin Console



Users Groups Global Init Scripts

Workspace Settings

Filter

Access Control

- > Workspace Access Control: Enabled
- > Cluster, Pool and Jobs Access Control: Enabled
- > Table Access Control: Disabled
- > Personal Access Tokens: Enabled Permission Settings
- > Workspace Visibility Control: Enabled
- > Cluster Visibility Control: Enabled
- > Job Visibility Control: Enabled

Create Cluster

New Cluster [Cancel](#) [Create Cluster](#)

Cluster Name
DevCluster

Cluster Mode ?
Standard

Databricks Runtime Version ? [Learn more](#)
Runtime: 7.6 (Scala 2.12, Spark 3.0.1)

Autopilot Options

Enable autoscaling ?
 Terminate after minutes of inactivity ?

Worker Type ? **Workers**

Standard_DS3_v2 14 GB Memory, 4 Cores, 0.75 DBU Spot instances ?

New Configure separate pools for workers and drivers for flexibility. [Learn more](#)

Driver Type

Same as worker 14 GB Memory, 4 Cores, 0.75 DBU



[Clusters /](#)

DevCluster [Edit](#) [Permissions](#) [Clone](#) [Restart](#) [Terminate](#) [Delete](#)

[Configuration](#) [Notebooks](#) [Libraries](#) [Event Log](#) [Spark UI](#) [Driver Logs](#) [Metrics](#) [Apps](#) [Spark Cluster UI - Master](#) ▾

Unrestricted

Cluster Mode ?
Standard

Databricks Runtime Version
7.6 (includes Apache Spark 3.0.1, Scala 2.12)

Autopilot Options

Enable autoscaling ?
 Terminate after minutes of inactivity ?

Worker Type ? **Workers** Current

Standard_DS3_v2 14 GB Memory, 4 Cores, 0.75 DBU 2 Spot instances ?

Driver Type

Standard_DS3_v2 14 GB Memory, 4 Cores, 0.75 DBU



 **databricks-rg-BigDataWorkspace-g3eexhzrn3ipw** X ...

Resource group

« »

+ Create Edit columns Delete resource group Refresh Export to CSV Open query Assign

Overview

Activity log

Access control (IAM)

Tags

Events

Settings

Deployments

Security

Policies

Properties

Locks

Cost Management

Cost analysis

Cost alerts (preview)

Essentials

Subscription ([change](#)) [REDACTED] Deployments : 2 Su

Subscription ID Location : East

Tags ([change](#)) : [Click here to add tags](#)

Type = **all** X Location = **all** X + Add filter

Showing 1 to 24 of 24 records. Show hidden types ①

<input type="checkbox"/> Name ↑↓	Type ↑↓
<input type="checkbox"/> workers-vnet	Virtual network
<input type="checkbox"/> a2bb2b656e4a4d4fa44dbbc3392208bc	Virtual machine
<input type="checkbox"/> cde9a7fc591448b8890d2b430a12f167	Virtual machine
<input type="checkbox"/> f286449f997e490399d465090f0cd710	Virtual machine
<input type="checkbox"/> dbstoragefnppxdzsotioi	Storage account
<input type="checkbox"/> a2bb2b656e4a4d4fa44dbbc3392208bc-publicIP	Public IP address
<input type="checkbox"/> cde9a7fc591448b8890d2b430a12f167-publicIP	Public IP address

Microsoft Azure | Databricks

Workspace

Workspace

Shared

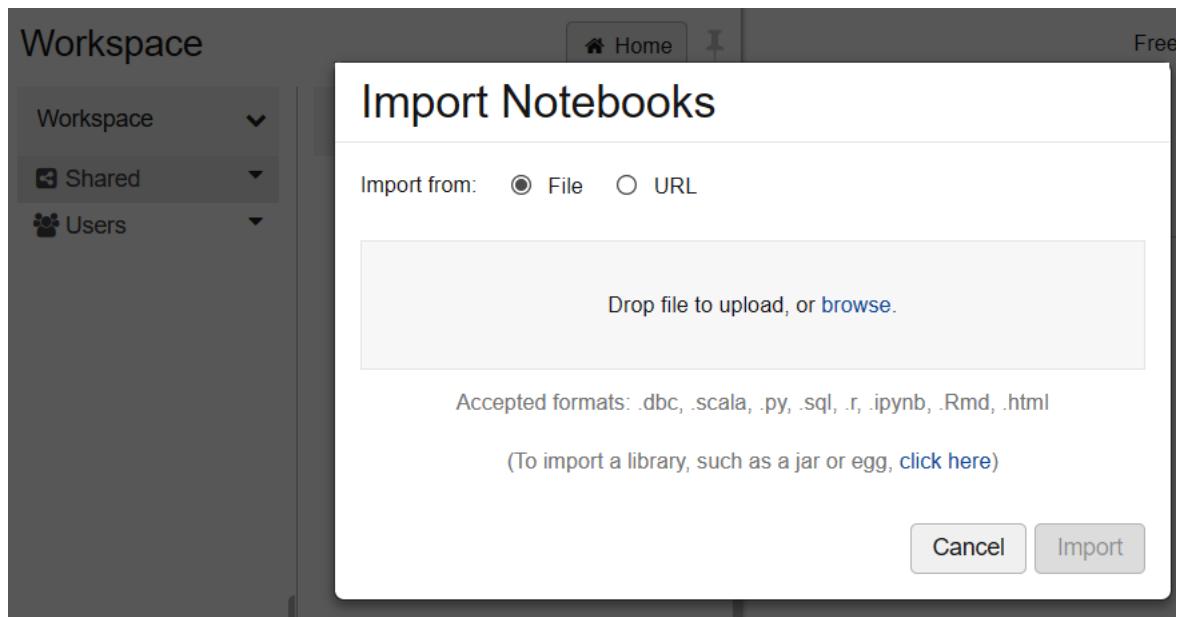
- Create
- Clone
- Import
- Export
- Copy Link Address

Home

Workspace

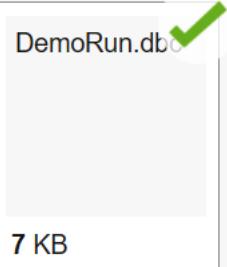
Recents

Data



Import Notebooks

Import from: File URL



Accepted formats: .dbc, .scala, .py, .sql, .r, .ipynb, .Rmd, .html

(To import a library, such as a jar or egg, [click here](#))

[Cancel](#) [Import](#)

Workspace

The workspace sidebar shows a "Shared" dropdown menu with "DemoRun" listed under "Shared". Other options in the "Shared" menu include "Shared" and "Users". The "Workspace" dropdown menu is also visible on the far left.

Jobs / Create NEW

CancelCreate[Runs](#) [Configuration](#)

Schedule Type Manual (Paused)
 Scheduled

Schedule ?

Every | at | : | (

Show Cron Syntax

Task

Type *

 |

Cluster * ?

 |

Parameters ?

 |

Add

[Advanced options >](#)

Configure New Cluster

2 Workers: 28 GB Memory, 8 Cores, 1.5 DBU

1 Driver: 14 GB Memory, 4 Cores, 0.75 DBU ?

Cluster Mode ?

 | |

Databricks Runtime Version ?

[Learn more](#) |

Autopilot Options

Enable autoscaling ?

Worker Type ?

Workers

 | Spot instances ?

New Configure separate pools for workers and drivers for flexibility. [Learn more](#)

Driver Type

 | CancelConfirm

Jobs / DailyJob NEW

Free trial ends in 13 days. Upgrade to Premium in Azure Portal ? BigDataWorkspace

DailyJob

Run Now More ...

Runs Configuration

ID: 145 Creator: [REDACTED].com Run As: [REDACTED].com Schedule: Paused - At 07:39 PM (UTC) Task: Notebook at /Shared/DemoRun

Active Runs

Run Start Time Run ID Launched Duration Spark Status

Run Now / Run Now With Different Parameters

0 - 0 < > 20 / Page

Completed Runs (past 60 days)

Latest Successful Run (Refreshes Automatically)

Run Start Time Run ID Launched Duration Spark Status

[View Details](#) Jul 9 2021, 19:47 PM IST 1 Manually 4m 40s [Spark UI / Logs / Metrics](#) Succeeded - [Delete](#)

The screenshot shows the Databricks Data workspace interface. On the left, the navigation sidebar includes options like Data Science & ML, Create, Workspace, Repos, Recents, Search, and Data (which has 1 notification). The main area is titled "Data" and shows the "Databases" section with "default" selected (2 entries). To the right, the "Tables" section shows "mnmdata" with 3 entries. A green circle with the number "1" is overlaid on the Data button in the sidebar, and another green circle with the number "3" is overlaid on the mnmdata entry in the tables list.

The screenshot shows the 'User Settings' page with a sidebar on the left containing icons for Home, Databricks, Create, Notebook, and Help. The main content area has a title 'User Settings' and three tabs: 'Access Tokens' (selected), 'Git Integration', and 'Notebook Settings'. A sub-section titled 'Personal access tokens can be used for secure authentication' contains a blue button labeled 'Generate New Token'. Below this, there are two columns: 'Comment' and 'Creation ↑'. A message states 'No tokens exist.'

User Settings

Access Tokens Git Integration Notebook Settings

Personal access tokens can be used for secure authentication

Generate New Token

Comment Creation ↑

No tokens exist.

Generate New Token ×

Comment

Used In Power BI report

Lifetime (days) ?

90

Cancel

Generate



Configuration Notebooks Libraries Event Log Spark UI Driver Logs Metrics

Driver Type

Standard_DS3_v2 14 GB Memory, 4 Cores, 0.75 DBU

▼ Advanced Options

Azure Data Lake Storage Credential Passthrough ?

Enable credential passthrough for user-level data access

Spark Tags Logging Init Scripts JDBC/ODBC Permissions

Server Hostname

adb-34 .7.azuredatabricks.net

Port

443

Protocol

HTTPS

HTTP Path

sql/protocolv1/o/ /0/ 9

JDBC URL ?

Get Data

Search

All
File
Database
Power Platform
Azure
Online Services
Other

Azure

- Azure SQL database
- Azure Synapse Analytics (SQL DW)
- Azure Analysis Services database
- Azure Database for PostgreSQL
- Azure Blob Storage
- Azure Table Storage
- Azure Cosmos DB
- Azure Data Explorer (Kusto)
- Azure Data Lake Storage Gen2
- Azure Data Lake Storage Gen1
- Azure HDInsight (HDFS)
- Azure HDInsight Spark
- HDInsight Interactive Query
- Azure Cost Management
- Azure Databricks**
- Azure Databricks Insights (Beta)

Certified Connectors | Template Apps Connect Cancel

Azure Databricks

Server Hostname ⓘ

HTTP Path ⓘ

Advanced Options (optional)

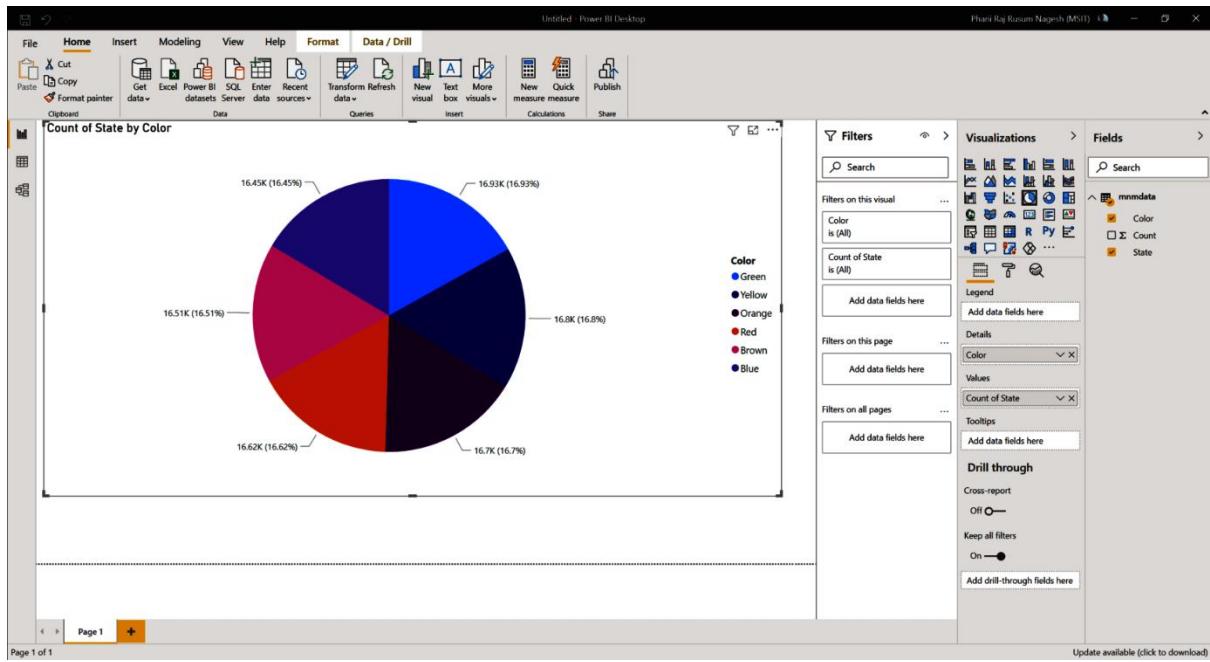
Database (optional) ⓘ

Batch Size (rows) (optional) ⓘ

Data Connectivity mode ⓘ

Import
 DirectQuery

OK Cancel



Chapter 2: Reading and Writing Data from and to Various Azure Services and File Formats

<input type="checkbox"/> Name ↑↓	Type ↑↓	Location ↑↓
<input type="checkbox"/> cookbookadlsgen2storage	Storage account	East US
<input type="checkbox"/> cookbookblobstorage	Storage account	East US

Home > demo

The screenshot shows the Azure Active Directory portal. In the left sidebar, under the 'Azure Active Directory' section, the 'App registrations' item is selected and highlighted with a red box and a green '1' marker. In the main content area, there is a 'New registration' button, which is also highlighted with a red box and a green '2' marker. A callout box contains the text: 'Try out the new App registrations'. Another callout box contains an informational message: 'Starting June 30th, 2020 we will no longer provide feature updates. Ap and Microsoft Graph. Learn more'.

demo | App registrations

Azure Active Directory

Administrative units

Enterprise applications

Devices

App registrations

Identity Governance

Application proxy

New registration

Endpoints

Try out the new App registrations

i Starting June 30th, 2020 we will no longer provide feature updates. Ap and Microsoft Graph. [Learn more](#)

Register an application

...

* Name

The user-facing display name for this application (this can be changed later).

ADLSGen2App

Supported account types

Who can use this application or access this API?

- Accounts in this organizational directory only (demo only - Single tenant)
- Accounts in any organizational directory (Any Azure AD directory - Multitenant)
- Accounts in any organizational directory (Any Azure AD directory - Multitenant)
- Personal Microsoft accounts only

By proceeding, you agree to the Microsoft Platform Policies [↗](#)

Register

 demo | App registrations ⚡ ...
Azure Active Directory

New registration Endpoints Troubleshooting Refresh Downl

Try out the new App registrations search preview! Click to enable the preview. →

Starting June 30th, 2020 we will no longer add any new features to Azure Active Directory. We will continue to provide technical support and security updates but we will no longer provide Microsoft Authentication Library (MSAL) and Microsoft Graph. [Learn more](#)

All applications **Owned applications** Deleted applications (Preview)

Start typing a name or Application ID to filter these results

Display name	Application
ADLSGen2App	06d3c701-2637

Administrative units
Enterprise applications
Devices
App registrations
Identity Governance
Application proxy
Licenses
Azure AD Connect
Custom domain names
Mobility (MDM and MAM)
Password reset
Company branding

ADLSGen2App

Search (Ctrl+ /) Overview 1 Quickstart Integration assistant

Delete Endpoints Preview features

Got a second? We would love your feedback on Microsoft

Manage

- Branding
- Authentication 2
- Certificates & secrets
- Token configuration

Essentials

Display name ADLSGen2App

Application (client) ID 06d3c701-2637-48fb-869e-19e3aa3587ff

Object ID 83763e3b-84c3-44b2-95d1-a1a0ecbffb57

Directory (tenant) ID 82 f0c82a0

Home > ADLSGen2App

ADLSGen2App | Certificates & secrets

Search (Ctrl+ /) Overview Quickstart Integration assistant

Manage

- Branding
- Authentication
- Certificates & secrets 1
- Token configuration
- API permissions
- Expose an API

Add a client secret

Got feedback?

Upload certificate

Thumbprint

No certificates have been added for

Client secrets

A secret string that the application uses

New client secret

Description

Add Cancel

Client secrets

A secret string that the application uses to prove its identity when requesting a token. Also can be referred to as application password.

+ New client secret

Description	Expires	Value	Secret ID	
ADLSGen2App SPN Secret	1/12/2022	H...	63e4d48d-cbb2-4a86-b52d-...	

cookbookadlsgen2storage | Access Control

Storage account

Search (Ctrl+ /) < 2 + Add Download role assignments ⚙

- Overview
- Activity log
- Tags
- Diagnose and solve problems
- Access Control (IAM) 1
- Data migration
- Events

Add role assignment 3

Add role assignment (Preview)

Add co-administrator

View my level of access to this resource.

View my access

Check access
Review the level of access a user, group, service or managed identity has to this resource. [Learn](#)

Add role assignment

Role ⓘ Storage Blob Data Contributor

Assign access to ⓘ User, group, or service principal

Select ⓘ adls

ADLSGen2App

cookbookadlsgen2storage | Access keys

Storage account

Search (Ctrl+ /)

<

2

Show keys



Set rotation reminder



Refresh

 Tables

Security + networking

 Networking

 Access keys 1

 Shared access signature

 Encryption

 Security

Data management

 Geo-replication

Access keys authenticate your applications' requests to this storage account, Key Vault, and replace them often with new keys. The two keys

Remember to update the keys with any Azure resources and apps

Storage account name

cookbookadlsgen2storage

key1

 Rotate key

Key 3

cookbookadlsgen2storage | Shared access signature

Storage account

Search (Ctrl+ /)

<

An account-level SAS can delegate access to multiple storage services (i.e. blob, file, queue, table). account-level SAS.

 Tables

Security + networking

 Networking

 Access keys

 Shared access signature 1

 Encryption

 Security

Data management

 Geo-replication

Learn more

Allowed services ⓘ

Blob File Queue Table

Allowed resource types ⓘ

Service Container Object

Allowed permissions ⓘ

Read Write Delete List Add Create Update Process

Blob versioning permissions ⓘ

Enables deletion of versions

Generate SAS and connection string

Connection string

SAS token

Blob service SAS URL

File service SAS URL

Queue service SAS URL

Table service SAS URL

Microsoft Azure | Databricks

Clusters / DevCluster_Spark2.

Configuration Notebooks Libraries 2

1 Name

Install Library

Library Source 3

Maven Upload DBFS/ADLS PyPI CRAN Workspace

Coordinates 4

com.microsoft.azure:spark-mssql-connector

Invalid Coordinates. Should match "groupId:artifactId:version"

Repository ?

Optional

Exclusions

Dependencies to exclude (log4j:log4j:junit:junit)

← Search Packages ×

Maven Central

Group Id	Artifact Id	Releases	Options
com.microsoft.azure	spark-mssql-connector	1.0.1 <input type="button" value="▼"/>	<input type="button" value="Select"/>

● DevCluster_Spark2.x

[Edit](#)[Permissions](#)[Clone](#)[Restart](#)[Configuration](#) [Notebooks \(0\)](#) [Libraries](#) [Event Log](#) [Spark UI](#) [Driver Logs](#) [Metrics](#) [Apps](#) [Spark](#)[Uninstall](#)[Install New](#)

<input type="checkbox"/>	Name	Type	Status	Source
<input type="checkbox"/>	com.microsoft.azure:spark-mssql-connector:1.0.1	Maven	● Installed	

SQLQuery1.sql - az...icrosoft.com (126)* X

```
SELECT TOP 10 * FROM [dbo].[CustomerTable]
```

121 % ▶

Results Messages

	C_CUSTKEY	C_NAME	C_ADDRESS	C_NATIONKEY
1	23172	Customer#000023172	Y,GqQajAoOed9GXqdHZAuMObNEt19WJH	0
2	39243	Customer#000039243	20JFX538mRg2wrHjcHUEg g4	0
3	66522	Customer#000066522	1bJiQmqr7mjooRovYIX	0
4	35373	Customer#000035373	5u20XRFD8JhBaLtOL3Kjh625GsK	0
5	30207	Customer#000030207	RC,iLpNQ1cCo2kSYncZJXlze82OTrOTMxwelSO	0
6	42862	Customer#000042862	IP26TV1sfWipR9jZECD,QkXVn5o	0
7	42916	Customer#000042916	4yqldhxAihy	0
8	27274	Customer#000027274	G2satKy0awkzwk qwoL5nk O WDr6yDwm47iLdO7	0
9	39478	Customer#000039478	Lxp2V2TIYbaCh	0
10	50677	Customer#000050677	,XzNNjEz,W4MKEKBHVZWW	0

New Container | Enable Azure Synapse Link | New Notebook | [Create new](#)

SQL API | DATA | Sales | NOTEBOOKS | [Gallery](#) | [My Notebooks](#)

Welcome to

Globally distributed, multi-mode

Start with Sample | New Container

New Container

* Database id Create new Use existing
Sales

* Container id Customer

* Partition key /C_MKTSEGMENT

Provision dedicated throughput

Unique keys

OK

Common Tasks | Recents

Clusters /

cosmosdbcluster ● !

[Edit](#)[Permissions](#)[Clone](#)[Configuration](#)[Notebooks](#)[Libraries](#)[Event Log](#)[Spark UI](#)[Driver Logs](#)

Unrestricted

Cluster Mode ?

Standard

Databricks Runtime Version

6.4 Extended Support (includes Apache Spark 2.4.5, Scala 2.11)

Autopilot Options

 Enable autoscaling ? Terminate after 60 minutes of inactivity ?

Worker Type ?

Workers Current

Standard_DS3_v2

14 GB Memory, 4 Cores, 0.75 DBU

1

1

Driver Type

Standard_DS3_v2

14 GB Memory, 4 Cores, 0.75 DBU

Clusters /

cosmosdbcluster ● !

[Configuration](#)[Notebooks](#)[Libraries](#)[Uninstall](#)[Install New](#) Name

Install Library

×

Library Source

[Upload](#)[DBFS/ADLS](#)[PyPI](#)[Maven](#)[CRAN](#)[Workspace](#)

Library Type

[Jar](#)[Python Egg](#)[Python Whl](#)3azure-
cosmosdb-
spark_2.4.0_2.1
3.6.14-uber.jar[Remove file](#)4[Cancel](#)[Install](#)

Clusters /

● cosmosdbcluster 🌐

[Edit](#)[Permissions](#)[Clone](#)[Configuration](#) [Notebooks](#) [Libraries](#) **Event Log** [Spark UI](#) [Driver Logs](#) [Metrics](#)[Uninstall](#)[Install New](#)

<input type="checkbox"/>	Name	Type	Status
<input type="checkbox"/>	azure_cosmosdb_spark_2_4_0_2_11_3_6_14_...	JAR	● Installed
<input type="checkbox"/>	[REDACTED]		● [REDACTED]

```
1 df_json = spark.read.option("multiline","true")
2 .json("dbfs:/mnt/SensorData/JsonData/SimpleJsonData/")
3 display(df_json)
```

▶ (3) Spark Jobs

▼ df_json: pyspark.sql.dataframe.DataFrame

```
Fuel: string
Transmission: string
_id: string
about: string
address: string
color: string
cost: string
currentowner: string
eventtime: string
index: long
latitude: double
longitude: double
phone: string
seatingcapacity: string
sellingcompany: string
```

```
1 df_json = spark.read.option("multiline","true")
2 .json("dbfs:/mnt/SensorData/JsonData/JsonData/")
3 display(df_json)
```

▶ (3) Spark Jobs

▶ 📄 df_json: pyspark.sql.dataframe.DataFrame = [Fuel: string, Transmission: string ...

	owners
1	▼ array
	▼ 0:
	name: "Burton Chase"
	phone: "+1 (896) 412-2343"
	▼ 1:
	name: "Williams Decker"
	phone: "+1 (857) 575-3797"
	▶ [{"name": "Galloway Woods", "phone": "+1 (868) 450-3041"}, {"name": "Chai

```
1 from pyspark.sql.functions import explode
2 data_df = df_json.select("_id",
3     explode("owners").alias("vehicleOwnersExplode"))
4 .select("_id", "vehicleOwnersExplode.*")
5 display(data_df)
```

▶ (2) Spark Jobs

▶ 📄 data_df: pyspark.sql.dataframe.DataFrame

 _id: string
 name: string
 phone: string

	_id	name	phone
1	60f623d0f2861584c56e2660	Burton Chase	+1 (896) 412-2343
2	60f623d0f2861584c56e2660	Williams Decker	+1 (857) 575-3797

```

1 jsonDF = data_df.withColumn("jsonCol",
2                               to_json(struct([data_df[x] for x in data_df.columns])))
3                               .select("jsonCol")
4 display(jsonDF)

```

▶ (2) Spark Jobs

▶ 📄 jsonDF: pyspark.sql.dataframe.DataFrame = [jsonCol: string]

	jsonCol
1	{"_id": "60f623d0f2861584c56e2660", "name": "Burton Chase", "phone": "+1 (896) 412-2343"}
2	{"_id": "60f623d0f2861584c56e2660", "name": "Williams Decker", "phone": "+1 (857) 575-3797"}
3	{"_id": "60f623d0b164af3c012566e9", "name": "Galloway Woods", "phone": "+1 (868) 450-3041"}
4	{"_id": "60f623d0b164af3c012566e9", "name": "Chaney Martin", "phone": "+1 (996) 430-3875"}
5	{"_id": "60f623d0722a2004debafe2a", "name": "Susanna Weber", "phone": "+1 (919) 431-2833"}

Component	Version
Apache Spark	2.4.x, 2.3.x, 2.2.x, and 2.1.x
Scala	2.11
Azure Databricks runtime version	> 3.4

Connector	Maven Coordinate	Scala Version
Spark 2.4.x compatible connector	com.microsoft.azure:spark-mssql-connector:1.0.2	2.11
Spark 3.0.x compatible connector	com.microsoft.azure:spark-mssql-connector_2.12:1.1.0	2.12

Chapter 3: Understanding Spark Query Execution

-1.Introduction to Jobs, Stages and Tasks (Python)

The screenshot shows the Apache Spark UI interface. At the top, there's a code editor window with the following Python code:

```
# Reading customer csv files in a dataframe
df_cust= spark.read.format("csv").option(True).load("dbfs:/mnt/Gen2/Customer/cs")
Command took 1.52 seconds -- by phanir@microsoft.com
display(df_cust.limit(10))
```

Below the code editor, a sidebar shows the current job: Job 7 (1/1).

The main area displays the "Completed Queries (35)" table:

ID	Description	Submitted	Duration	Job IDs
34	display(df_cust_agg)	2021/07/15 15:16:55	0.7 s	[9]
33	display(df_cust_agg)	2021/07/15 15:16:52	3 s	[8]
32	display(df_cust.limit(10))	2021/07/15 15:16:51	0.2 s	[7]
31	# Reading customer csv files in a dataframe df_...	2021/07/15	0.2 s	[5]

Below the table, a section titled "Details for Query 3" provides the following information:

- Submitted Time:** 2021/02/14 05:14:05
- Duration:** 0.3 s
- Succeeded Jobs:** 1

Expand all the details in the query plan visualization

A detailed view of the "Scan csv +details" stage is shown, highlighting the following values:

WholeStageCodegen	
number of files read	10
filesystem read data size total (min, med, max)	64.0 KB (64.0 KB, 6
filesystem read data size (sampled) total (min, med, max)	128.0 KB (128.0 KB
filesystem read time (sampled) total (min, med, max)	107 ms (107 ms, 10
metadata time	0
size of files read total (min, med, max)	23.1 MB (23.1 MB, 2
rows output	10

At the bottom of the stage details, there is a "CollectLimit" button.

```
1 | display(df_cust.limit(10))
```

▼ (1) Spark Jobs

▼ Job 1 [View](#) (Stages: 1/1)

Stage 1: 1/1 [i](#)

1

Summary Metrics for 1 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	0.2 s	0.2 s	0.2 s	0.2 s	0.2 s
GC Time	0 ms	0 ms	0 ms	0 ms	0 ms
Input Size / Records	64.0 KB / 10	64.0 KB / 10	64.0 KB / 10	64.0 KB / 10	64.0 KB / 10

▼ Aggregated Metrics by Executor

Executor ID	Address	Task Time	Total Tasks	Failed Tasks	Killed Tasks	Succeeded Tasks	Input Size / Records	Blacklisted
0 stdout stderr	10.139.64.5:42133	0.3 s	1	0	0	1	64.0 KB / 10	false

▼ Tasks (1)

Index	ID	Attempt	Status	Locality Level	Executor ID	Host	Launch Time	Duration	GC Time	Input Size / Records	Errors
0	1	0	SUCCESS	PROCESS_LOCAL	0	10.139.64.5 stdout stderr	2021/02/14 05:14:05	0.2 s		64.0 KB / 10	

```
1 | display(df_cust_agg)
```

▼ (1) Spark Jobs

▼ Job 2 [View](#) (Stages: 2/2)

Stage 2: 4/4 [i](#)

Stage 3: 200/200 [i](#)

	C_MKTSEGMENT	sum_acctbal	avg_acctbal	max_bonus
1	BUILDING	135888621.94	4508.28153208148	9999.99
2	HOUSEHOLD	135873341.17	4500.756605717316	9999.23
3	AUTOMOBILE	133866847.09	4499.423470354938	9999.96
4	MACHINERY	134438861.67	4488.92656415906	9999.64

Jobs Stages Storage Environment Executors SQL

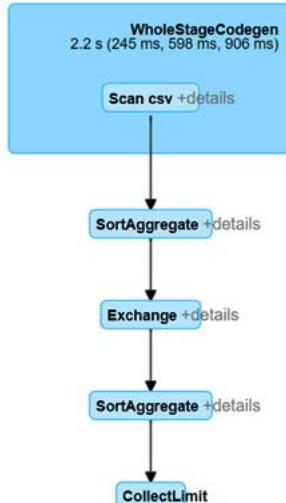
Details for Query 6

Submitted Time: 2021/02/14 05:36:44

Duration: 1 s

Succeeded Jobs: 4

Expand all the details in the query plan visualization



```
display(df_cust.limit(10))
```

▼ (1) Spark Jobs
▶ Job 4 View (Stages: 1/1)

	C_CUSTKEY	C_NAME	C_ADDRESS
1	35165	Customer#000035165	eNQSVdTId
2	30597	Customer#000030597	S9s1dDut80
3	42279	Customer#000042279	ABCvDnNa3
4	42578	Customer#000042578	i6VNaE7iSz
5	37854	Customer#000037854	dL6LCTLpY
6	40053	Customer#000040053	qh8Q6gaffF

Showing all 10 rows.



Jobs	Stages	Storage	Environment	Executors	SQL	JDBC/ODBC Server
Structured Streaming						
SQL						
Completed Queries: 403						
Completed Queries (403)						
Page: 1 2 3 4 5 > 5 Pages. Jump to 1 . Show 100 items in a page. Go						
ID	Description	Submitted	Duration	Job IDs		
402	display(df_cust_agg)	2021/07/15 16:52:10	0.8 s	[7][8]	+details	
401	display(df_cust_agg)	2021/07/15 16:52:08	2 s	[5][6]	+details	
400	display(df_cust.limit(10))	2021/07/15 16:52:07	0.3 s	[4]	+details	

1 2 3 4 5 > 5 Pages. Jump to 1 . Show 100 items in a page. Go

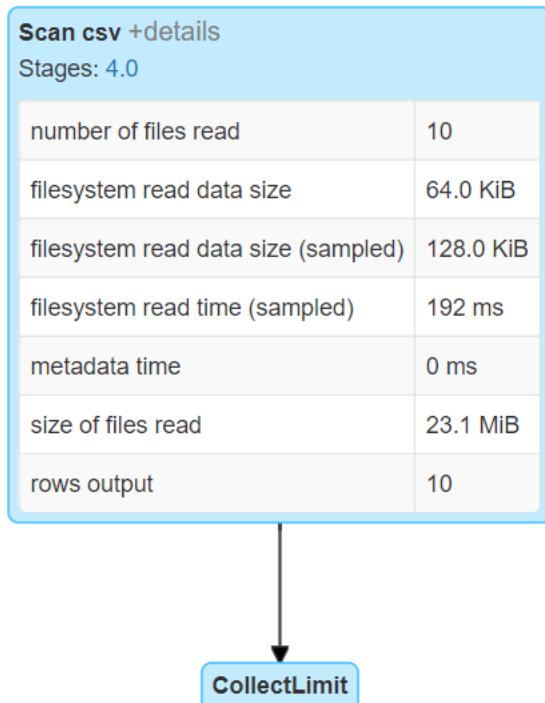
Completed Queries (403)

Page: 1 2 3 4 5 > 5 Pages. Jump to 1 . Show 100 items in a page. Go

ID	Description	Submitted	Duration	Job IDs
402	display(df_cust_agg)	2021/07/15 16:52:10	0.8 s	[7][8]
401	display(df_cust_agg)	2021/07/15 16:52:08	2 s	[5][6]
400	display(df_cust.limit(10))	2021/07/15 16:52:07	0.3 s	[4]

Succeeded Jobs:

Expand all the details in the query plan visualization



DevCluster

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

Completed Jobs (9)

Page: 1

1 Pages. Jump to 1 . Show 100 items in a page.

Job Id (Job Group) ▾	Description	Submitted	Duration	Stages: Succeeded/Total
8 (4730382624081364794_6058216927087136628_7fc7c5f176de4503afe796a8206810ee)	display(df_cust_agg) collectResult at OutputAggregator.scala:194	2021/07/15 16:52:10	0.2 s	1/1 (1 skipped)
7 (4730382624081364794_6058216927087136628_7fc7c5f176de4503afe796a8206810ee)	display(df_cust_agg) collectResult at OutputAggregator.scala:194	2021/07/15 16:52:10	0.5 s	1/1
6 (4730382624081364794_7166894886713111235_7fc7c5f176de4503afe796a8206810ee)	display(df_cust_agg) collectResult at OutputAggregator.scala:194	2021/07/15 16:52:09	0.4 s	1/1 (1 skipped)

DevCluster

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

Sort By

Page: 1 2 >

2 Pages. Jump to 1 . Show 100 items in a page. Go

ID	Description	Submitted	Duration ▾	Job IDs
196	# Reading customer csv files in a dataframe df_... 3	2021/07/15 18:26:38	7 s	[0]
198	display(df_cust_agg)	2021/07/15 18:26:55	2 s	[3][4]
190	show databases	2021/07/15 18:26:34	0.8 s	

Associated SQL Query: 196
 Job Group: 8287089703055764537_9154657314170381234_d0ddf8fc1a4c456f932860dfd7316f42

Completed Stages: 1

- ▶ Event Timeline
- ▶ DAG Visualization
- ▼ Completed Stages (1)

Page: 1 1 Pages. Jump to 1 . Show

Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input
0	8287089703055764537	# Reading customer csv files in a dataframe df_c... load at NativeMethodAccessorImpl.java:0 +details	2021/07/15 18:26:38	7 s	1/1	64.0 KiB

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

▶ Aggregated Metrics by Executor

Tasks (10)

Show 20 entries

Search:

Task Index ▲	ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Input Size / Records	Write Time	Shuffle Write Size / Records
0	23	0	SUCCESS	PROCESS_LOCAL	1	10.1.128.5	stdout stderr	2021-07-15 23:56:57	0.3 s		2.3 MiB / 15000	0.0 ms	75 B / 1
1	24	0	SUCCESS	PROCESS_LOCAL	0	10.1.128.4	stdout stderr	2021-07-15 23:56:57	0.2 s		2.3 MiB / 14999	0.0 ms	75 B / 1
2	25	0	SUCCESS	PROCESS_LOCAL	1	10.1.128.5	stdout stderr	2021-07-15 23:56:57	0.3 s		2.3 MiB / 15000	0.0 ms	75 B / 1

● DevCluster

Edit Permissions Clone Restart Terminate Delete

Configuration Notebooks (1) Libraries Event Log Spark UI Driver Logs Metrics Apps Spark Cluster UI - Master

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read
0	10.1.128.4:41935	Active	0	22.1 kB / 3.6 GiB	0.0 B	4	0	0	17	17	39 s (2 s)	37 MiB	750 B
driver	10.1.128.6:42891	Active	0	22.1 kB / 3.3 GiB	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B
1	10.1.128.5:38627	Active	0	22.1 kB / 3.6 GiB	0.0 B	4	0	0	17	17	28 s (2 s)	32.5 MiB	4.8 kB

```
# Reading customer csv files in a dataframe with specify
df_cust= spark.read.format("csv").option("header",True).
True).load("dbfs:/mnt/Gen2/Customer/csvFiles")
```

▼ (2) Spark Jobs

- ▶ Job 7 View (Stages: 1/1)
- ▶ Job 8 View (Stages: 1/1)

▶ df_cust: pyspark.sql.DataFrame = [C_CUSTKEY: integer

Cmd 6

```
#Creating the dataframe by specifying the schema using .schema() option. We don;t see any DAG getting created with schema is specified
df_cust_sch=
spark.read.format("csv").option("header",True).schema(cust_schema).load("/mnt/Gen2/Customer/csvFiles/")

▶ df_cust_sch: pyspark.sql.dataframe.DataFrame = [C_CUSTKEY: integer, C_NAME: string ... 6 more fields]
```

Command took 0.25 seconds -- by at /2021, 12:23:25 AM on DevCluster

Cmd 4

```
1 # Getting the number of partitions
2 print(f" Number of partitions with all files read = {df_cust.rdd.getNumPartitions()}")
```

Number of partitions with all files read = 10

Cmd 8

```
1 #Getting the number of partitions by specifying only one csv file.
2 print(f" Number of partitions with 1 file read = {df_cust_sch.rdd.getNumPartitions()}")
```

Number of partitions with 1 file read = 1

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

▼Details

```
-- Parsed Logical Plan --
GlobalLimit 21
+- LocalLimit 21
  +- Project [cast(C_MKTSEGMENT#2240 as string) AS C_MKTSEGMENT#2317, cast(AvgAcctBal#2284 as string) AS AvgAcctBa
    +- Filter (AvgAcctBal#2284 > cast(4500 as double))
    +- Filter (C_MKTSEGMENT#2240 = MACHINERY)
      +- Aggregate [C_MKTSEGMENT#2240], [C_MKTSEGMENT#2240, avg(C_ACCTBAL#2239) AS AvgAcctBal#2284]
        +- Relation[C_CUSTKEY#2234,C_NAME#2235,C_ADDRESS#2236,C_NATIONKEY#2237,C_PHONE#2238,C_ACCTBAL#2239,C]

-- Analyzed Logical Plan --
C_MKTSEGMENT: string, AvgAcctBal: string
GlobalLimit 21
+- LocalLimit 21
  +- Project [cast(C_MKTSEGMENT#2240 as string) AS C_MKTSEGMENT#2317, cast(AvgAcctBal#2284 as string) AS AvgAcctBa
```

```

== Physical Plan ==
AdaptiveSparkPlan isFinalPlan=false
+- == Current Plan ==
  HashAggregate(keys=[C_MKTSEGMENT#2240], functions=[finalmerge_count(merge count#2357L) AS count(1)#2340L]
  +- Exchange hashpartitioning(C_MKTSEGMENT#2240, 200), true, [id=#2016]
    +- HashAggregate(keys=[C_MKTSEGMENT#2240], functions=[partial_count(1) AS count#2357L])
      +- FileScan csv [C_MKTSEGMENT#2240] Batched: false, DataFilters: [], Format: CSV, Location: InMemory
iles, PartitionFilters: [], PushedFilters: [], ReadSchema: struct<C_MKTSEGMENT:string>
+- == Initial Plan ==
  HashAggregate(keys=[C_MKTSEGMENT#2240], functions=[finalmerge_count(merge count#2357L) AS count(1)#2340L]
  +- Exchange hashpartitioning(C_MKTSEGMENT#2240, 200), true, [id=#2016]
    +- HashAggregate(keys=[C_MKTSEGMENT#2240], functions=[partial_count(1) AS count#2357L])
      +- FileScan csv [C_MKTSEGMENT#2240] Batched: false, DataFilters: [], Format: CSV, Location: InMemory
iles, PartitionFilters: [], PushedFilters: [], ReadSchema: struct<C_MKTSEGMENT:string>

```

```

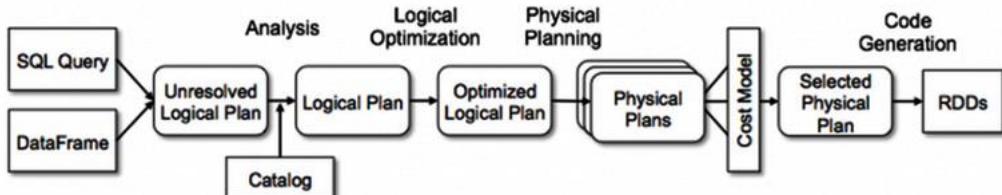
== Physical Plan ==
AdaptiveSparkPlan isFinalPlan=false
+- == Current Plan ==
  HashAggregate(keys=[C_MKTSEGMENT#2240], functions=[finalmerge_count(merge count#2362L) AS count(1)#2352L]
  +- Exchange hashpartitioning(C_MKTSEGMENT#2240, 200), true, [id=#2062]

```

```

DataFilters: [isnotnull(C_MKTSEGMENT#370), (C_MKTSEGMENT#370 = MACHINERY)], Format: CSV, Location:
InMemoryFileIndex[dbfs:/mnt/Gen2/Customer/csvFiles], PartitionFilters: [], PushedFilters:
[IsNotNull(C_MKTSEGMENT, EqualTo(C_MKTSEGMENT,MACHINERY)], ReadSchema:
struct<C_ACCTBAL:double,C_MKTSEGMENT:string>

```



↑ Upload + Add Directory ⏪ Refresh | ⏴ Rename └ Del

Authentication method: Access key ([Switch to Azure AD User Account](#))
Location: rawdata

Search blobs by prefix (case-sensitive)

Name	Modified	Access
<input type="checkbox"/> Customer		
<input type="checkbox"/> Orders		

Cmd 8

```
# Check the execution plan and you will find
df_ord_sch.join(df_cust_sch, df_ord_sch.O_C)
```

▼ (1) Spark Jobs

- Job 0 View (Stages: 4/4)
 - Stage 0: 8/8
 - Stage 1: 10/10
 - Stage 2: 200/200
 - Stage 3: 1/1

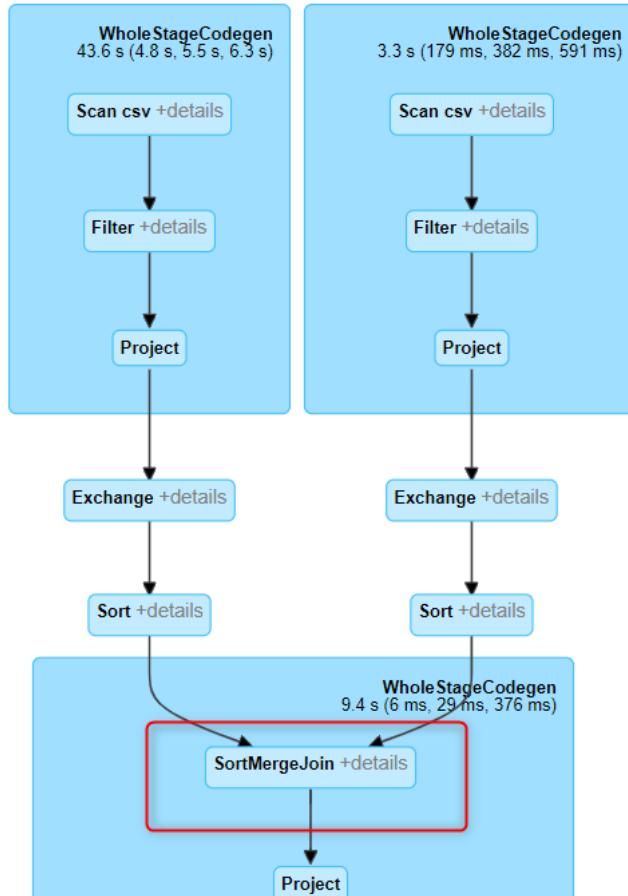
Out[8]: 1500000

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server

Completed Queries: 16

Completed Queries (16)

ID	Description	Submitted	Duration	Job IDs
15	# Check the execution plan and you will find so... +details	2021/07/15 19:09:27	12 s	[0]
14	show tables in `default` +details	2021/07/15 19:09:17	8 ms	



```
df = spark.read.format("csv").option("header","true").load("/mnt/Gen2/Customer/csvFiles/")
print("Number of partitions = " + str(df.rdd.getNumPartitions()))
df.write.mode("overwrite").option("path", "dbfs:/tmp")
```

▼ (2) Spark Jobs

- Job 21 View (Stages: 1/1)
 - Stage 29: 1/1
- Job 22 View (Stages: 1/1)
 - Stage 30: 10/10

df: pyspark.sql.dataframe.DataFrame = [C_CUSTKEY: string,

Number of partitions = 10

Command took 5.01 seconds -- by phanir@microsoft.com at 7/16/2021 11:49:47 AM

md 3

```
repartitionedDF = df.repartition(8)
print('Number of partitions: {}'.format(repartitionedDF.rdd.getNumPartitions()))
```

Number of partitions: 8

md 4

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server

Details for Job 22

Status: SUCCEEDED
Associated SQL Query: 613
Job Group: 8056729216216354018_4646673387408174941_077c7cc526734da2a09136f704d75fe0
Completed Stages: 1

- Event Timeline
- DAG Visualization

Completed Stages (1)

Page:	1	1 Pages. Jump to	1	Show	100	items in a page.	Go
Stage Id	Pool Name	Description	Submitted	Duration	Tasks:		
30	8056729216216354018	spark.conf.set("spark.sql.files.maxPartitionBy... saveAsTable at NativeMethodAccessorImpl.java:0 +details	2021/07/15 19:19:47	3 s	10/10		

```
1 # Check the duration of execution with default 200 partitions
2 spark.conf.set("spark.sql.shuffle.partitions", 200)
3 mktSegmentDF = df.groupBy("C_MKTSEGMENT").count().collect()
```

► (2) Spark Jobs

Command took 1.43 seconds -- by [REDACTED] on CookBookCluster

```
1 # Check the duration of execution with default 30 partitions
2 spark.conf.set("spark.sql.shuffle.partitions", 30)
3 mktSegmentDF = df.groupBy("C_MKTSEGMENT").count().collect()
```

► (2) Spark Jobs

Command took 0.61 seconds -- [REDACTED] on CookBookCluster

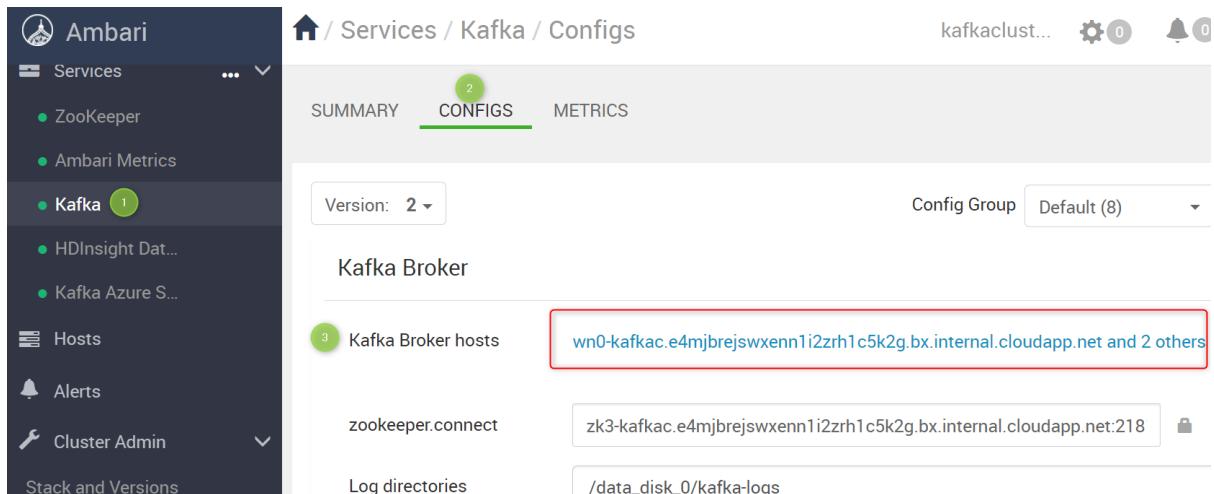
```
sortDF = df.sort('C_CUSTKEY')
display(sortDF)
```

► (1) Spark Jobs

► sortDF: pyspark.sql.dataframe.DataFrame = [C_CUSTKEY: string, C_NAME: string ... 6 more fields]

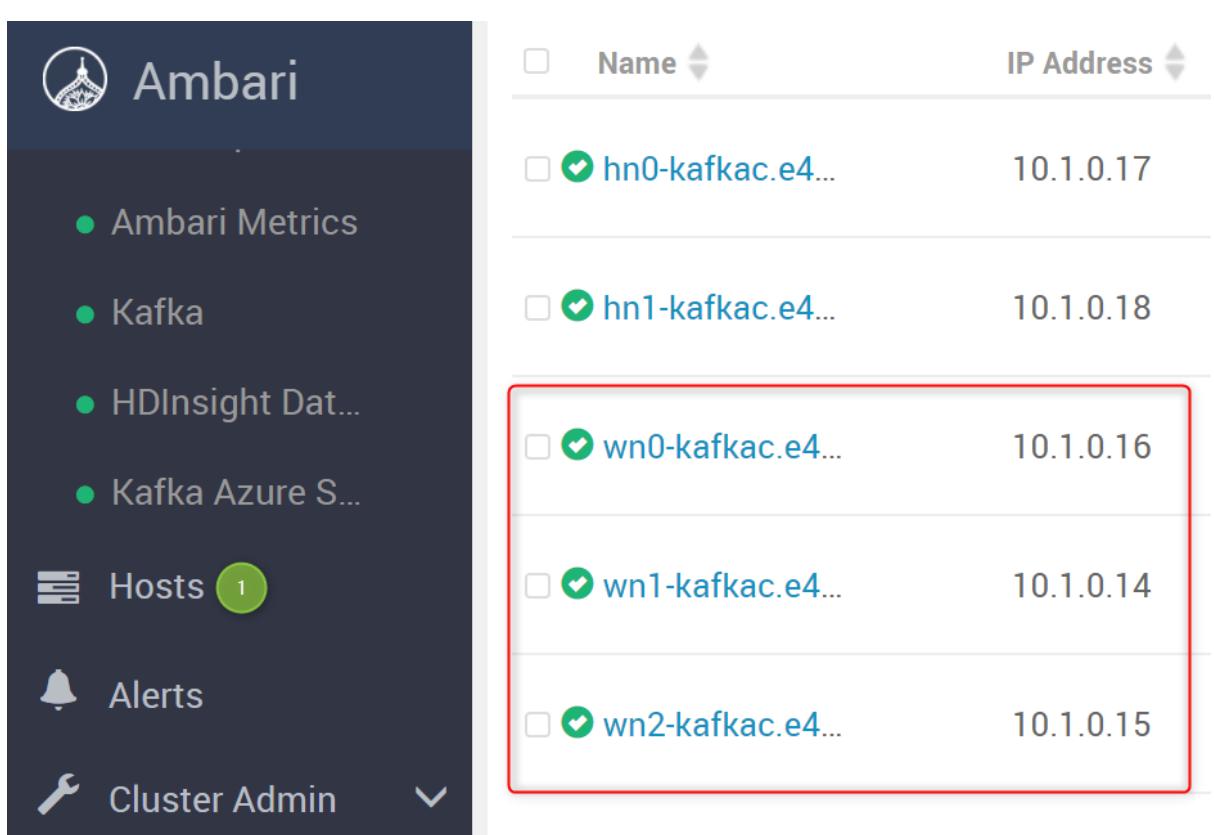
	C_CUSTKEY	C_NAME	C_ADDRESS	C_NATIONKEY	C_PHONE
1	1	Customer#000000001	IVhzIApeRb ot,c,E	0	10-989-741-2988
2	10	Customer#000000010	Lbrg3ElDieI0B10bB0Aymm	0	10-741-346-9870
3	100	Customer#000000100	mFowluhnHjp2GjCIYYavkW kUwOjlTCQ	0	10-749-445-4907
4	1000	Customer#000001000	Ne,bJD40ae8ZXo XFrOkEmACe	0	10-730-275-2976
5	10000	Customer#000010000	OI7fOBe9O7SEIX7xSxqR	0	10-923-910-8193
6	100000	Customer#000100000	Glw,RrxSCw6UDA5	0	10-488-516-6598

Chapter 4: Working with Streaming Data



The screenshot shows the Ambari interface for managing Kafka configurations. The left sidebar lists services like ZooKeeper, Ambari Metrics, Kafka (selected), HDInsight Data, and Kafka Azure S. The main area displays Kafka Broker configuration details. Under 'Kafka Broker hosts', it shows 'wn0-kafkac.e4mjbrejswxenn1i2zrh1c5k2g.bx.internal.cloudapp.net and 2 others'. Below that are 'zookeeper.connect' (set to 'zk3-kafkac.e4mjbrejswxenn1i2zrh1c5k2g.bx.internal.cloudapp.net:218') and 'Log directories' ('/data_disk_0/kafka-logs'). A red box highlights the host list.

Name	IP Address
hn0-kafkac.e4...	10.1.0.17
hn1-kafkac.e4...	10.1.0.18
wn0-kafkac.e4...	10.1.0.16
wn1-kafkac.e4...	10.1.0.14
wn2-kafkac.e4...	10.1.0.15



The screenshot shows the Ambari interface for managing hosts. The left sidebar lists services like Ambari Metrics, Kafka, HDInsight Data, and Kafka Azure S. The main area displays a list of hosts. A red box highlights the host list.

Name	IP Address
hn0-kafkac.e4...	10.1.0.17
hn1-kafkac.e4...	10.1.0.18
wn0-kafkac.e4...	10.1.0.16
wn1-kafkac.e4...	10.1.0.14
wn2-kafkac.e4...	10.1.0.15

Ambari Services

- ZooKeeper
- Ambari Metrics
- Kafka** (1)
- HDInsight Dat...
- Kafka Azure S...

Hosts

Alerts

log.roll.hours: 168

log.retention.hours: 168

listeners: PLAINTEXT://localhost:9092

Advanced kafka-broker

auto.create.topics.enable: true

```
producer = KafkaProducer(bootstrap_servers=['10.0.0.16:9092', '10.1.0.14:9092',
                                             '10.1.0.15:9092'],
                        value_serializer=lambda x:
                        dumps(x).encode('utf-8'))
```

Event Hubs Namespace

Search (Ctrl+ /)

+ Add

Shared access policies

Policy	Claims
RootManageSharedAccessKey	Manage, Send, Listen
readwrite	Manage, Send, Listen

Save Discard Delete ...

Manage

Send

Listen

Primary key: 7tn8

Secondary key: k

Connection string-primary key: Endpoint=sb://c...hub.servi...

TestCluster

Edit Permissions Clone Re

Configuration Notebooks (0) Libraries Event Log Spark UI Driver Logs

Uninstall Install New

	Name	Type	Status
<input type="checkbox"/>	com.microsoft.azure:azure-eventhubs-spark_2...	Maven	Installed

com.microsoft.azure:azure-eventhubs-spark_2.12-2.3.20

```

print( EventHubdf.isStreaming)
print( EventHubdf.printSchema())

```

```

True
root
|-- body: binary (nullable = true)
|-- partition: string (nullable = true)
|-- offset: string (nullable = true)
|-- sequenceNumber: long (nullable = true)
|-- enqueuedTime: timestamp (nullable = true)
|-- publisher: string (nullable = true)
|-- partitionKey: string (nullable = true)
|-- properties: map (nullable = true)
|   |-- key: string
|   |-- value: string (valueContainsNull = true)
|-- systemProperties: map (nullable = true)
|   |-- key: string
|   |-- value: string (valueContainsNull = true)

```

The screenshot shows the Azure portal interface for managing shared access policies in a Kafka-enabled Event Hubs namespace named "kafkaenabledeventhubns".

Left Panel (Settings):

- Shared access policies (highlighted with a red box and numbered 1)
- Scale
- Geo-Recovery
- Networking
- Encryption
- Properties
- Locks

Right Panel (Shared access policies screen):

Add Policy: A modal window is open to add a new policy.

Policy	Claims
RootManageSharedAccessKey	Manage, Send, Listen
sendreceivekafka	Manage, Send, Listen

Policy Details (sendreceivekafka):

- Manage (checked)
- Send (unchecked)
- Listen (checked)

Keys:

- Primary key:** Value is "4v", highlighted with a red box and numbered 2.
- Secondary key:** Placeholder field.

Connection strings:

- Connection string-primary key:** Value is "Endpoint=sb://kafkaenabledeventhubns.servicebus.wi...", highlighted with a red box and numbered 3.
- Connection string-secondary key:** Placeholder field.

Buttons at the top right: Save, Discard, Delete, ...

```
display(dbutils.fs.ls("/mnt/Gen2/LogFiles"))
```

► (2) Spark Jobs

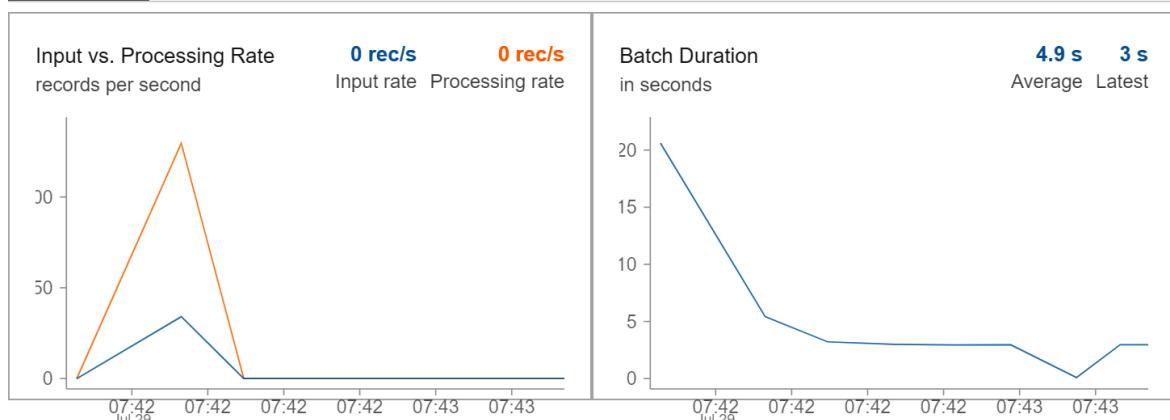
	path	name	size
1	dbfs:/mnt/Gen2/LogFiles/Log1.json	Log1.json	211
2	dbfs:/mnt/Gen2/LogFiles/Log2.json	Log2.json	209
3	dbfs:/mnt/Gen2/LogFiles/Log3.json	Log3.json	209
4	dbfs:/mnt/Gen2/LogFiles/Log4 - Copy.json	Log4 - Copy.json	209
5	dbfs:/mnt/Gen2/LogFiles/Log4.json	Log4.json	209

```
%fs head dbfs:/mnt/Gen2/LogFiles/Log3.json
```

```
{  
  "source": "Application 1",  
  "data": "Log data example 1",  
  "IP": "10.22.135.7",  
  "log_level": "DEBUG",  
  "log_type": "Error",  
  "log_app": "app",  
  "log_timestamp": "2020-08-25T17:45:13+07:00"  
}
```

▼ 0a223bb6-9302-435d-981b-55b7b4baf3da Last updated: 15 seconds ago

Dashboard Raw Data

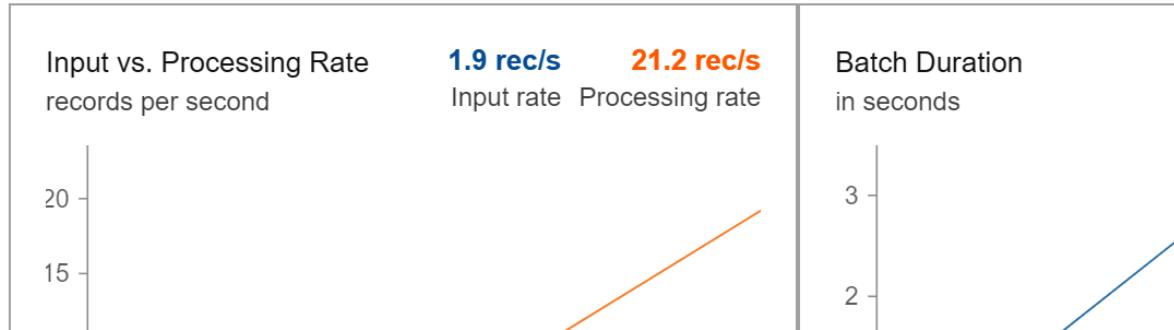


▶ (1) Spark Jobs

▼ 🐾 0a223bb6-9302-435d-981b-55b7b4baf3da

Last updated: 30 seconds ago

Dashboard Raw Data



```
%sql
SELECT * FROM VehiclechkpointKafkaEventHub_Delta_Agg_Tumbling ORDER BY Window desc
```

▶ (1) Spark Jobs

	window	id	count
1	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	19eeda6f-0489-4c68-937a-f91ba866643b	6
2	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	afa56211-c75c-4d76-9f17-2e550cd51fa8	5
3	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	66ce2ce6-3b23-4f92-a3b7-a7dfd3fe7caf	5
4	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	a529fb2a-5f7c-4aa1-976c-e4b683b919f2	5
5	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	c7bb7aa9-cc1c-453f-a86e-ca171c710e85	5
6	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	c91ff0d5-bffc-458f-9a15-25fac0fc6cc8	6

Showing all 10 rows.

```
%sql
SELECT * FROM VehiclechkpointKafkaEventHub_Delta_Agg_Tumbling ORDER BY Window desc
```

▶ (1) Spark Jobs

	window	id	count
1	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	a1a4b4ba-7f9c-4058-aa15-b79f3fd5dd7d	2
2	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	c7bb7aa9-cc1c-453f-a86e-ca171c710e85	1
3	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	67f5187f-d1c0-431b-844a-2f99f49d59e3	1
4	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	45e71b4f-f46e-4e58-959d-e2afa7830875	2
5	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	66ce2ce6-3b23-4f92-a3b7-a7dfd3fe7caf	1
6	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	a529fb2a-5f7c-4aa1-976c-e4b683b919f2	1

Showing all 30 rows.

```
%sql
SELECT * FROM VehiclechkpointKafkaEventHub_Delta_Agg_Tumbling WHERE id='c7bb7aa9-cc1c-453f-a86e-ca171c710e85'
```

▶ (2) Spark Jobs

	window	id
1	▶ {"start": "2021-07-30T01:54:00.000+0000", "end": "2021-07-30T01:55:00.000+0000"}	c7bb7aa9-cc1c-453f-a86e-ca171c710e85
2	▶ {"start": "2021-07-30T01:55:00.000+0000", "end": "2021-07-30T01:56:00.000+0000"}	c7bb7aa9-cc1c-453f-a86e-ca171c710e85
3	▶ {"start": "2021-07-30T01:53:00.000+0000", "end": "2021-07-30T01:54:00.000+0000"}	c7bb7aa9-cc1c-453f-a86e-ca171c710e85

```
%sql
select * from VehiclechkpointKafkaEventHub_Delta_Agg_Overlapping
where id ='c2c7cb35-2f97-4fab-ab23-62fe24eca7af'
order by window desc
```

▶ (1) Spark Jobs

	window	id	count
1	▶ {"start": "2021-07-30T03:08:00.000+0000", "end": "2021-07-30T03:10:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	1
2	▶ {"start": "2021-07-30T03:07:00.000+0000", "end": "2021-07-30T03:09:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	13
3	▶ {"start": "2021-07-30T03:06:00.000+0000", "end": "2021-07-30T03:08:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	34
4	▶ {"start": "2021-07-30T03:05:00.000+0000", "end": "2021-07-30T03:07:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	28
5	▶ {"start": "2021-07-30T03:04:00.000+0000", "end": "2021-07-30T03:06:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	6

```
%sql
select * from VehiclechkpointKafkaEventHub_Delta_Agg_Overlapping
where id ='c2c7cb35-2f97-4fab-ab23-62fe24eca7af'
order by window desc
```

▶ (1) Spark Jobs

	window	id	count
1	▶ {"start": "2021-07-30T03:08:00.000+0000", "end": "2021-07-30T03:10:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	2
2	▶ {"start": "2021-07-30T03:07:00.000+0000", "end": "2021-07-30T03:09:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	14
3	▶ {"start": "2021-07-30T03:06:00.000+0000", "end": "2021-07-30T03:08:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	34
4	▶ {"start": "2021-07-30T03:05:00.000+0000", "end": "2021-07-30T03:07:00.000+0000"}	c2c7cb35-2f97-4fab-ab23-62fe24eca7af	28

rawdata

Container

Search (Ctrl+ /)

Upload Change access level Refresh

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

_azuretmpfolder\$

Customer

CustomerDelta

Orders

Vehicle_Agg

Vehicle_Chkpoint1

[Upload](#) [Change access level](#) [Refresh](#) | [Delete](#) | [←](#)

Authentication method: Access key ([Switch to Azure AD User Account](#))

Location: [rawdata / Vehicle_Chkpoint1 / state / 0](#)

Search blobs by prefix (case-sensitive)

[+ Add filter](#)

Name	Modified
<input type="checkbox"/> [..]	
<input type="checkbox"/> 0	
<input type="checkbox"/> 1	
<input type="checkbox"/> 10	
<input type="checkbox"/> 100	

Authentication method: Access key ([Switch to Azure AD User Account](#))

Location: [rawdata / Vehicle_Chkpoint1 / offsets](#)

Search blobs by prefix (case-sensitive)

[+ Add filter](#)

Name	Modified
<input type="checkbox"/> [..]	
<input type="checkbox"/> __tmp_path_dir	7/30/2021, 9:29:52 AM
<input type="checkbox"/> 0	7/30/2021, 9:27:16 AM
<input type="checkbox"/> 1	7/30/2021, 9:28:07 AM
<input type="checkbox"/> 2	7/30/2021, 9:28:42 AM

«

Vehicle_Chkpoint1/offsets/1

Blob

Upload Change access level ...

to Azure AD User Account

Location: rawdata / Vehicle_Chkpoint1 / offsets

Search blobs by prefix (cas...)

Show deleted blobs

Add filter

Name	...
[...]	...
_tmp_path_dir	...
0	...
1	...

Save Discard Download Refresh Delete

Overview Versions Snapshots Edit Generate SAS

⚠ The file 'Vehicle_Chkpoint1/offsets/1' may not render correctly as it contains an unrec

```
1 v1
2 {"batchWatermarkMs":0,"batchTimestampMs":162761748722
3 ["kafkaenabledhub2":{"0":40}]]
```

Chapter 5: Integrating with Azure Key Vault, App Configuration, and Log Analytics

Home > Create a resource > Key Vault >

Create key vault

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *



Resource group *



[Create new](#)

Instance details

Key vault name * ⓘ



Region *



Pricing tier * ⓘ



Recovery options

Soft delete protection will automatically be enabled on this key vault. This feature allows you to recover or permanently delete a key vault and secrets for the duration of the retention period. This protection applies to the key vault and the secrets stored within the key vault.

To enforce a mandatory retention period and prevent the permanent deletion of key vaults or secrets prior to the retention

[Review + create](#)

[< Previous](#)

[Next : Access policy >](#)

Home > Create a resource > Key Vault >

Create key vault

Basics Access policy Networking Tags Review + create

Enable Access to:

- Azure Virtual Machines for deployment ⓘ
- Azure Resource Manager for template deployment ⓘ
- Azure Disk Encryption for volume encryption ⓘ

Permission model

- Vault access policy
- Azure role-based access control

[+ Add Access Policy](#)

Current Access Policies

Name	Email	Key Permissions	Secret Permissions	Certificate Permissions	Action
USER					
 Vinod Jaiswal	Vinod.Jaiswal [REDACTED]	9 selected	7 selected	15 selected	Delete

[Review + create](#) [< Previous](#) [Next : Networking >](#)

Create key vault

Basics Access policy Networking Tags Review + create

Network connectivity

You can connect to this key vault either publicly, via public IP addresses or service endpoints, or privately, using a private endpoint.

Connectivity method

- Public endpoint (all networks)
- Public endpoint (selected networks)
- Private endpoint

[Review + create](#)

[< Previous](#)

[Next : Tags >](#)

Home > cookbookkeyvaultdemo > cookbookkeyvaultdemo

cookbookkeyvaultdemo | Secrets

Key vault

2

Search (Ctrl+ /) <> + Generate/Import ⏪ Refresh ⏪ Restore Backup ⏪ Manage deleted secrets

Overview Activity log Access control (IAM) Tags Diagnose and solve problems Events

Settings

1 Keys Secrets Certificates Access policies Networking Security Properties Locks

There are no secrets available.

This screenshot shows the 'Secrets' blade in the Azure Key Vault interface. On the left, a sidebar lists various settings: Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Events, Settings, Keys, Secrets (which is highlighted with a green circle containing the number 1), Certificates, Access policies, Networking, Security, Properties, and Locks. The main area displays a table with three columns: Name, Type, and Status. A message at the top states, "There are no secrets available." There are two green circles with the numbers 1 and 2 placed over the 'Secrets' link and the top right corner of the main content area respectively.

Home > cookbookkeyvaultdemo > cookbookkeyvaultdemo >

Create a secret

Upload options Manual

Name * blobstoragesecret ✓

Value * ✓

Content type (optional) blob storage secret ✓

Set activation date

Set expiration date

Enabled Yes No

Create

This screenshot shows the 'Create a secret' form. It includes fields for Name (blobstoragesecret), Value (represented by a series of dots), Content type (blob storage secret), and Enabled status (Yes). Validation icons (green checkmarks) are present next to the Name, Value, and Content type fields. At the bottom is a large blue 'Create' button.

Home > cookbookkeyvaultdemo > cookbookkeyvaultdemo >

blobstoragesecret

Versions

+ New Version ⏪ Refresh 🗑 Delete ⏴ Download Backup

Version	Status	Activation date
CURRENT VERSION		
7a913485c27e4fb49171c2caf2ee0180	✓ Enabled	

Home >

KafkaStreamingWorkSpace

Azure Databricks Service

Search (Ctrl+ /) ⏪ Delete

- Overview
- Activity log
- Access control (IAM)
- Tags
- Settings
 - Virtual Network Peerings
- Encryption
- Properties
- Locks

Essentials

[View Cost](#) | [JSON View](#)

Status	Managed Resource Group
Active	databricks-rg-KafkaStreamingWorkSpace-[REDACTED]
Resource group	CookbookRG
Location	URL
East US	https://adb-[REDACTED].azuredatabricks.net
Subscription	Pricing Tier
[REDACTED]	premium
Subscription ID	Virtual Network
[REDACTED]	vnet-cookbook
Tags (change)	Private Subnet Name
Click here to add tags	private-subnet-adb

Microsoft Azure | Databricks

HomePage / Create Secret Scope

Create Secret Scope

A store for secrets that is identified by a name and backed by a specific store type. [Learn more](#)

Scope Name [?](#)
KeyVaultScope

Manage Principal [?](#)
Creator

Azure Key Vault [?](#)

DNS Name
<https://cookbookkeyvaultdemo.vault.azure.net/>

Resource ID
`/subscriptions/...`



cookbookkeyvaultdemo | Properties X

Key vault

Search (Ctrl+/) Save Discard Refresh

Events

Settings

Keys

Secrets

Certificates

Access policies

Networking

Security

Properties

Locks

Monitoring

Alerts

Name	cookbookkeyvaultdemo
Sku (Pricing tier)	Standard
Location	eastus
Vault URI	https://cookbookkeyvaultdemo.vault.azure.net/
Resource ID	<code>/subscriptions/...</code> resourceGroups/CookbookRG/providers...
Subscription ID	<code>...</code>
Subscription Name	<code>...</code>
Directory ID	<code>...</code>
Directory Name	<code>...</code>
Soft-delete	Soft delete has been enabled on this key vault

 **App Configuration** Add to Favorites
Microsoft 4.6 (18 ratings)
Azure benefit eligible [?]

Create

[Overview](#) [Plans](#) [Usage Information + Support](#) [Reviews](#)

Azure App Configuration allows developers to store, retrieve and manage access to application settings all in one place. It is easy to set up and simple to use from any application. It gives developers the ability to modify an application's behavior on demand without having to redeploy the application. App Configuration offers the following benefits:

Home > Create a resource > Marketplace > App Configuration >

Create App Configuration

Azure App Configuration provides a service to centrally manage application settings and feature flags. Modern programs, especially programs running in a cloud, generally have many components that are distributed in nature. Spreading configuration settings across these components can lead to hard-to-troubleshoot errors during an application deployment. Use App Configuration to store all the settings for your application and secure their accesses in one place. [Learn more](#)

Project Details

Subscription * [Redacted]

Resource group * CookbookRG [?]
[Create new](#)

Instance Details

Resource name * Enter resource name

Location * East US [?]

Pricing tier * Standard [?]
[View full pricing details](#)

Review + create

< Previous

Next: Tags >

Home > App Configuration > DevAppconfigurationRes

»  **DevAppconfigurationRes | Configuration explorer** ...
App Configuration

Search (Ctrl+/) [<] [>] 2 Create Refresh

Authentication method: Access keys (Switch to Azure AD) [?]

Date : Select date Keys : Select key Labels : Select label

Loaded 1 key-values with 1 unique keys. | Show values | Expand all | Collapse all

Key	Value	Label
StorageKey	(Hidden value)	(No label)

Operations

[?] Configuration explorer 1

» **DevAppconfigurationRes | Access keys** ...

Access keys 1

Search (Ctrl+ /) <> Show values

Primary key
Regenerate
Id (credential)
..... Copy

Secret
..... Copy

Connection string
..... Copy 2

Scale up
Identity
Encryption
Private endpoint connections
Properties
Locks

Create Log Analytics workspace ...

is collected and stored.

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * ⓘ

..... ▼

Resource group * ⓘ

CookbookRG ▼

[Create new](#)

Instance details

Name * ⓘ

DevLogAnalyticsWorkspacedemo ✓

Region * ⓘ

Central US ▼

Review + Create

[« Previous](#)

[Next : Pricing tier >](#)

Home >

DevLogAnalyticsWorkspace

Log Analytics workspace

Search (Ctrl+ /) Delete

View Cost | JSON View

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Settings

Locks

Agents management

Agents configuration

Custom logs

Computer Groups

Essentials

Resource group (change)
cookbookrg

Status
Active

Location
East US

Subscription (change)

Subscription ID

Tags (change)
Click here to add tags

Workspace Name
DevLogAnalyticsWorkspace

Workspace ID

Pricing tier
Pay-as-you-go

Access control mode
Use resource or workspace permissions

Operational issues
OK

Get started with Log Analytics

This screenshot shows the 'Essentials' section of the Log Analytics workspace settings. It displays the workspace name 'DevLogAnalyticsWorkspace', its location 'East US', and its current status as 'Active'. The 'Operational issues' section shows a green checkmark next to 'OK', indicating no operational problems at the time the screenshot was taken.

Home > DevLogAnalyticsWorkspace

DevLogAnalyticsWorkspace | Logs

Log Analytics workspace

Search (Ctrl+ /) Feedback Queries Query explorer

New Query 1 Run Time range : Last 24 hours Save Share New alert rule Export

Type your query here or click one of the queries to start

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Settings

Locks

Agents management

Agents configuration

Custom logs

Computer Groups

Linked storage accounts

Queries History

This screenshot shows the main 'Logs' page for the workspace. It features a 'New Query 1' editor with a 'Run' button and a time range selector set to 'Last 24 hours'. Below the editor is a text input field with placeholder text: 'Type your query here or click one of the queries to start'. The left sidebar contains navigation links for various workspace components, and the right side shows a 'Queries History' section which is currently empty.

Queries

Always show Queries ⓘ | [Community Git repo ↗](#) | [Documentation ↗](#) [X](#)

Query packs: Select query packs

Category 

 Search

 Add filter

★ Favorites

All Queries

Applications

Audit

Azure Monitor

Azure Resources

Containers

Databases

Desktop Analytics

IT & Management ...

Top 3 browser exceptions

What were the highest reported exceptions today?

[Run](#)

Example query

Slowest pages

What are the 3 slowest pages, and how slow are they?

[Run](#)

Example query

Failed requests – top 10

What are the 3 slowest pages, and how slow are they?

[Run](#)

Example query

Failing dependencies

Which 5 dependencies failed the most today?

[Run](#)

Example query

Diagnostic setting

...

 Save  Discard  Delete  Feedback

A diagnostic setting specifies a list of categories of platform logs and/or metrics that you want to collect from a resource, and one or more destinations that you would stream them to. Normal usage charges for the destination will occur. [Learn more about the different log categories and contents of those logs](#)

Diagnostic setting name *

Collect Databricks Metrics 

Category details

log

Audit

metric

AllMetrics

Destination details

Send to Log Analytics workspace

Archive to a storage account

Stream to an event hub

Send to partner solution

Diagnostic setting

...

Save Discard Delete Feedback

A diagnostic setting specifies a list of categories of platform logs and/or metrics that you want to collect from a resource, and one or more destinations that you would stream them to. Normal usage charges for the destination will occur. [Learn more about the different log categories and contents of those logs](#)

Diagnostic setting name *

Collect Databricks Metrics



Category details

log

Audit

metric

AllMetrics

Destination details

Send to Log Analytics workspace

Subscription

[REDACTED]

Log Analytics workspace

DevLogAnalyticsWorkspace (eastus)

Archive to a storage account

Stream to an event hub

Send to partner solution

Home > Log Analytics workspaces >

» DevLogAnalyticsWorkspace

Search (Ctrl+ /)

Delete

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Settings

Locks

Agents management

Agents configuration

Custom logs

Computer Groups

Select one or more data sources to connect to the workspace

Azure virtual machines (VMs)

Windows and Linux Agents management

Azure Activity logs

Storage account log

System Center Operations Manager

Add monitoring solutions that provide insights for applications and services in your environment

View solutions

Create alerts to proactively detect any issue that arise in your workspace

Learn more

Maximize your Log Analytics experience

Search and analyze logs

Use Log Analytics rich query language to analyze logs
[View logs](#)

Manage alert rules

Notify or take action in response to important information in your data

[Set alerts](#)

Manage usage and costs

Understand your usage of Log Analytics and estimate your costs for each month

[Manage costs](#)

DevLogAnalyticsWorkspace | Logs

Log Analytics workspace

Search (Ctrl+ /)

- [Agents configuration](#)
- [Custom logs](#)
- [Computer Groups](#)
- [Linked storage accounts](#)
- [Network Isolation](#)

General

- [Workspace summary](#)
- [Workbooks](#)
- [Logs](#)
- [Solutions](#)
- [Usage and estimated costs](#)
- [Properties](#)
- [Service Map](#)

Workspace Data Sources

- [Virtual machines](#)

New Query 1

DevLogAnalyticsW... Select scope

Tables Queries Functions ...

Time range: Last 24 hours | Save Share +

Type your query here or click one of the queries to start

Search

Filter Group by: Solution

Collapse all

Favorites

You can add favorites by clicking on the ★ icon

- LogManagement**
- ▶ [DatabricksAccounts](#)
- ▶ [DatabricksClusters](#)
- ▶ [DatabricksDBFS](#)
- ▶ [DatabricksJobs](#)
- ▶ [DatabricksNotebook](#)
- ▶ [DatabricksSecrets](#)
- ▶ [DatabricksSSH](#)
- ▶ [DatabricksWorkspace](#)

Queries History

No queries history

ActionName	RequestId	Response	RequestParams	Type
attachNotebook	68156378-2965-451d-b944-259d12e8...	{"statusCode":2...	{"notebookId":"452834995439556","clusterId":"0302-165541-d...	DatabricksNo
attachNotebook	987de542-72e4-4105-999c-ae718f6f4...	{"statusCode":2...	{"notebookId":"452834995439556","clusterId":"0302-165541-d...	DatabricksNo
attachNotebook	b2577f28-3d95-45d9-acd5-5c5bb5e...	{"statusCode":2...	{"notebookId":"452834995439556","clusterId":"0302-165541-d...	DatabricksNo

Chapter 6: Exploring Delta Lake in Azure Databricks

Location: rawdata

Search blobs by prefix (case-sensitive)

	Name	Modified
<input type="checkbox"/>	 Customer	
<input type="checkbox"/>		
<input type="checkbox"/>		
<input type="checkbox"/>	 Orders	

Authentication method: Access key (Switch to Azure)

Location: rawdata / Orders

Search blobs by prefix (case-sensitive)

	Name
<input type="checkbox"/>	 [..]
<input type="checkbox"/>	 newParquetFiles
<input type="checkbox"/>	 parquetFiles

```
query_source1 = append_kafkadata_stream(topic='eventhubsource1')
```

Cancel ••

▶ (1) Spark Jobs ━━━━━━

▶ ⚙ edd4228e-e858-4f6e-a1c1-5f5acd0f62f4 Last updated: 15 seconds ago

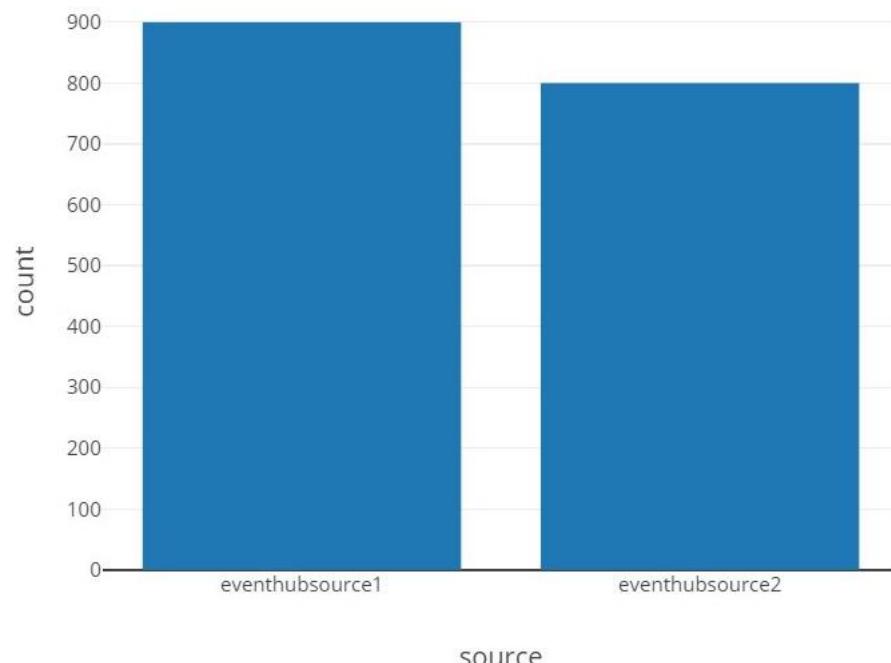
Cmd 9

```
query_source2 = append_kafkadata_stream(topic='eventhubsource2')
```

Cancel ••

▶ (1) Spark Jobs ━━━━━━

▶ ⚙ 76343cb9-5960-48d2-9c1a-d1f49ff27749 Last updated: 15 seconds ago



Authentication method: Access key (Switch)

Location: rawdata / Customer

Search blobs by prefix (case-sensitive)

	Name	Modified
<input type="checkbox"/>	[..]	
<input type="checkbox"/>		
<input type="checkbox"/>		
<input type="checkbox"/>		
<input type="checkbox"/>	parquetFiles	
<input type="checkbox"/>	parquetFiles_Daily	

Authentication method: Access key (Switch to Azure AD User Account)

Location: rawdata / Customer / delta

Search blobs by prefix (case-sensitive)

Name
<input type="checkbox"/> [..]
<input type="checkbox"/> _delta_log
<input type="checkbox"/> part-00000-ced1e1a9-799e-49e8-a628-df1cff240361-c000.snappy.parquet

Location: rawdata / Customer / delta / _delta_log

Search blobs by prefix (case-sensitive)

Name

<input type="checkbox"/>	[..]
<input type="checkbox"/>	_tmp_path_dir
<input type="checkbox"/>	00000000000000000000000000000000.crc
<input type="checkbox"/>	00000000000000000000000000000000.json

Customer/delta/_delta_log/00000000000000000000000000000000.json ...

Blob

Save Discard Download Refresh | Delete

Overview Versions Edit **Edit** Generate SAS

```
1 {"commitInfo":{"timestamp":1628417004020,"userId":"1231225035485455","userName":"...","operation":"WRITE"},"
2 {"protocol":{"minReaderVersion":1,"minWriterVersion":2}}
3 {"metaData":{"id":"67bfb2f9-75e0-41d9-a4b4-107750160f90","format":{"provider":"parquet","options":{}}, "schemaString":"{...typ
4 {"add":{"path":"part-00000-ced1ea9-799e-49e8-a628-df1cff240361-c000.snappy.parquet","partitionValues":{}, "size":8618802,"mo
5 }
```

Customer/delta/_delta_log/00000000000000000000000000000000.crc ...

Blob

Save Discard Download Refresh | Delete

Overview Versions Edit **Edit** Generate SAS

The file 'Customer/delta/_delta_log/00000000000000000000000000000000.crc' may not render correctly as it contains an unrecognized extension.

```
1 {"tableSizeBytes":8618802,"numFiles":1,"numMetadata":1,"numProtocol":1,"numTransactions":0}
2 |
```

Customer/delta/_delta_log/00000000000000000000000000000001.json ...

Blob

Save Discard Download Refresh | Delete

Overview Versions Edit **Edit** Generate SAS

```
1 |,"operation":"WRITE","operationParameters":{"mode":"Append","partitionBy":[]}, "notebook": {"notebook
2 |Values":{},"size":3691463,"modificationTime":1628417370000,"dataChange":true,"stats":{"\\"numRecords\\":45000
3 |}
```

Customer/delta/_delta_log/0000000000000000000002.json ...

Blob

Save Discard Download Refresh Delete

Overview Versions Edit Generate SAS

```
1 {"commitInfo":{"timestamp":1628417520625,"userId":"1231225035485455","userName":...,"operation":"DELETE"}  
2 {"remove":{"path":"part-00000-ced1e1a9-799e-49e8-a628-df1cff240361-c000.snappy.parquet","deletionTimestamp":1628417520623,...  
3 {"add":{"path":"part-00000-527e3cd6-fbce-474e-be7c-00821ede3c5c-c000.snappy.parquet","partitionValues":{},"size":8620958,"m  
4
```

Customer/delta/_delta_log/0000000000000000000003.json ...

Blob

Save Discard Download Refresh Delete

Overview Versions Edit Generate SAS

```
1 com" "operation":"UPDATE","operationParameters":{"predicate":"(C_CUSTKEY#4112 = 101275)"}, "notebook": {"notebookId":...  
2 "imestamp":1628417846199,"dataChange":true,"extendedFileMetadata":true,"partitionValues":{}, "size":8620958}  
3 .ues":{}, "size":8620802, "modificationTime":1628417847000, "dataChange":true, "stats": {"\\"numRecords\\":104999, \\"minValu  
4
```

Vehicle_Delta/_delta_log/00000000000000000000000000000000.json ...

Blob

Save Discard Download Refresh Delete

Overview Versions Edit Generate SAS

```
1 ", "operation":"STREAMING UPDATE","operationParameters":{"outputMode":"Append", "queryId": "76...":2774...  
2
```

Overview Versions Edit Generate SAS

```
1 titionMetrics":{"numRemovedFiles":0,"numOutputRows":75,"numOutputBytes":3440,"numAddedFiles":1}}}  
2  
3 d-4e91-ba08-1f5cd3dd\","timestamp":\\"2021-04-16T16:59:57.818Z\","rpm":99,"speed":100,"kms":995.  
4
```

```
%sql  
DELETE FROM CustomerDelta WHERE C_MKTSEGMENT='FURNITURE'
```

► (3) Spark Jobs
Error in SQL statement: ConcurrentAppendException: Files were added to the root of the table by a concurrent update. Please try the operation again.

```
%sql  
DELETE FROM CustomerDeltaPartition WHERE C_MKTSEGMENT='FURNITURE'
```

► (4) Spark Jobs

OK

```
%sql
DELETE FROM CustomerDeltaPartition WHERE C_MKTSEGMENT='MACHINERY' AND C_CUSTKEY=1377
```

▶ (4) Spark Jobs

⊕ Error in SQL statement: ConcurrentAppendException: Files were added to partition [C_MKTSEGMENT=MACHINERY] by a concurrent update. Please try the operation again.

Conflicting commit: {"timestamp":1628440339088,"userId":"55","userName":"", "operationType": "WRITE", "operationParameters": {"mode": "Append", "partitionBy": ["C_MKTSEGMENT"]}, "notebook": {"notebookId": "2"}, "clusterId": "", "readVersion": 33, "isolationLevel": "Serializable", "isBlindAppend": true, "operationMetrics": {"numFiles": "4", "numOutputBytes": "2461672", "numOutputRows": "29949"}}

Refer to <https://docs.microsoft.com/azure/databricks/delta/concurrency-control> for more details.

```
%sql
DELETE FROM CustomerDeltaPartition WHERE C_MKTSEGMENT='MACHINERY' AND C_CUSTKEY=1377
```

▶ (8) Spark Jobs

OK

```
display(ordersDF.filter("O_ORDERSTATUS = '0'").groupBy("O_OrderYear", "O_ORDERPRIORITY").agg(count("*").alias("TotalOrders")).orderBy("O_OrderYear", "O_ORDERPRIORITY", ascending=False).limit(20))
```

▶ (2) Spark Jobs

	O_OrderYear	O_ORDERPRIORITY	TotalOrders
1	1996	1-URGENT	46014
2	1996	4-NOT SPECIFIED	45776
3	1997	4-NOT SPECIFIED	45754
4	1997	2-HIGH	45689
5	1997	1-URGENT	45662
6	1996	3-MEDIUM	45657
7	1996	5-LOW	45639

Showing all 20 rows.



Command took 1.70 seconds --

```
display(spark.sql("DROP TABLE IF EXISTS Orders"))

display(spark.sql("CREATE TABLE Orders USING DELTA LOCATION '/mnt/Gen2Source/Orders/OrdersDelta'"))

display(spark.sql("OPTIMIZE Orders ZORDER BY (O_ORDERPRIORITY)"))
```

▶ (17) Spark Jobs

OK
OK

	path	metrics
1	null	▶ {"numFilesAdded": 7, "numFilesRemoved": 28, "filesAdded": {"min": 2178106, "max": 3726077, "avg": 3452511, "totalFiles": 7, "totalSize": 24167579}, "filesRemoved": {"min": 443929, "max": 1024790, "avg": 887174, "totalFiles": 28, "totalSize": 24840888}, "partitionsOptimized": 7, "zOrderStats": {"strategyName": "minCubeSize(107374182400)", "inputCubeFiles": {"num": 0, "size": 0}, "inputOtherFiles": {"num": 28, "size": 24840888}, "inputNumCubes": 0, "mergedFiles": {"num": 28, "size": 24840888}, "numOutputCubes": 7, "mergedNumCubes": null}, "numBatches": 1}

▶ [Table] ordersDeltaDF: pyspark.sql.dataframe.DataFrame = [O_ORDERKEY: integer, O_C

	O_OrderYear	O_ORDERPRIORITY	TotalOrders
1	1996	1-URGENT	46014
2	1996	4-NOT SPECIFIED	45776
3	1997	4-NOT SPECIFIED	45754
4	1997	2-HIGH	45689
5	1997	1-URGENT	45662
6	1996	3-MEDIUM	45657
7	1996	5-L OW	45639

Showing all 20 rows.



Command took 0.72 seconds --

```
display(spark.sql("DROP TABLE IF EXISTS Orders"))

display(spark.sql("CREATE TABLE Orders USING DELTA LOCATION '/mnt/Gen2Source/Orders/OrdersDelta'"))

display(spark.sql("OPTIMIZE Orders ZORDER BY (O_ORDERPRIORITY)"))
```

▶ (20) Spark Jobs

OK
OK

	path	metrics
1	null	[{"numFilesAdded": 7, "numFilesRemoved": 28, "filesAdded": {"min": 2178106, "max": 3726077, "avg": 3452511, "totalFiles": 7, "totalSize": 24167579}, "filesRemoved": {"min": 443929, "max": 1024790, "avg": 887174, "totalFiles": 28, "totalSize": 24840888}, "partitionsOptimized": 7, "zOrderStats": {"strategyName": "minCubeSize(107374182400)", "inputCubeFiles": {"num": 0, "size": 0}, "inputOtherFiles": {"num": 28, "size": 24840888}, "inputNumCubes": 0, "mergedFiles": {"num": 28, "size": 24840888}, "numOutputCubes": 7, "mergedNumCubes": null}, "numBatches": 1}]

```
df_cust.write.format("delta").mode("append").saveAsTable("customer")
```

▶ (1) Spark Jobs

com.databricks.sql.transaction.tahoe.schema.InvariantViolationException: NOT NULL constraint violated for column: C_NAME.

```
%sql
-- Properties column of DESCRIBE output has the details about constraints
DESCRIBE DETAIL customer;
```

	partitionColumns	numFiles	sizeInBytes	properties
1	00	0	0	{"delta.constraints.mktsegment": "C_MKTSEGMENT = 'BUILDING'"}

Showing all 1 rows.



```
display(spark.sql("DESCRIBE HISTORY customerdelta"))
```

▶ (2) Spark Jobs

	version	timestamp	userId	userName	operation	operationParameters
15	5	2021-08-08T11:11:01.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
16	4	2021-08-08T11:10:59.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
17	3	2021-08-08T11:10:57.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
18	2	2021-08-08T11:10:55.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
19	1	2021-08-08T11:10:52.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
20	0	2021-08-08T11:10:50.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}

default.customerdelta

Refresh

Details History

Filter

Refresh

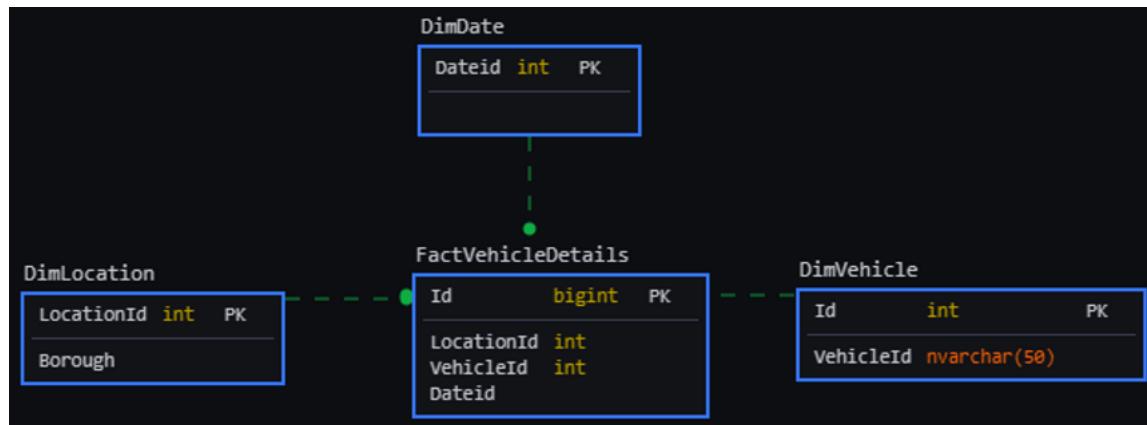
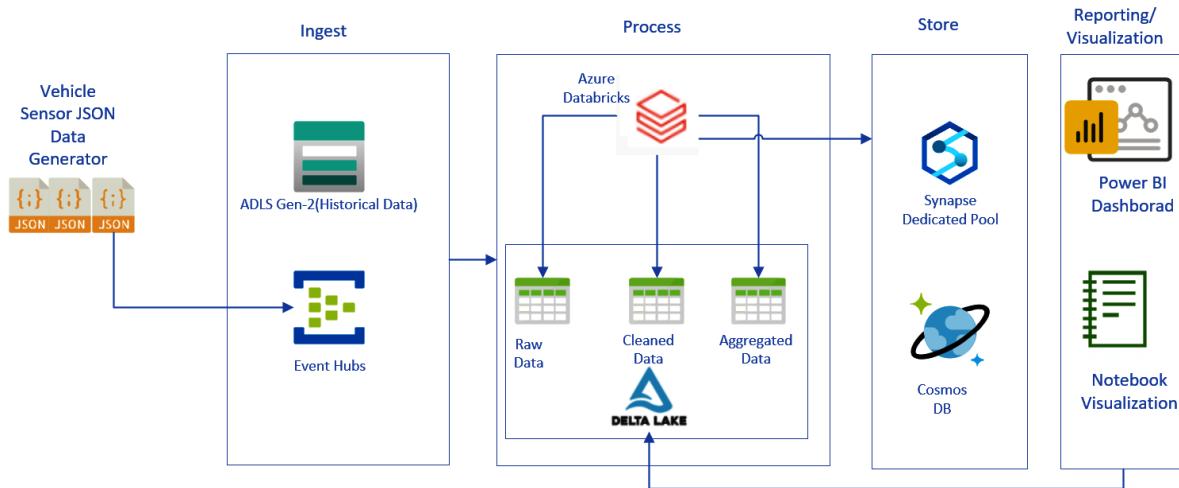
	version	timestamp	userId	userName	operation	operationParameters
15	5	2021-08-08T11:11:01.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
16	4	2021-08-08T11:10:59.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
17	3	2021-08-08T11:10:57.000+0000	1231225035485455		WRITE	▶ {"mode": "Append", "partition": 0}
18	2	2021-08-08T11:10:55.000+0000	1231225035485455	pham@microsoft.com	WRITE	▶ {"mode": "Append", "partition": 0}

Location: rawdata / CustomerDelta

Search blobs by prefix (case-sensitive)

Name	Modified	Access tier
□ [..]		
□ _delta_log		
□ part-00000-0625ed... 8/8/2021, 4:41:10 PM	Hot (Inferred)	
□ part-00000-0885cb... 8/8/2021, 4:40:55 PM	Hot (Inferred)	

Chapter 7: Implementing Near-Real-Time Analytics and Building a Modern Data Warehouse



cookbookcosmosdb | Keys ...

Azure Cosmos DB account

Search (Ctrl+/)

Private Endpoint Connections

CORS

Dedicated Gateway

Keys

Advisor Recommendations

Add Azure Cognitive Search

Add Azure Function

Advanced security (preview)

Locks

Containers

Read-write Keys Read-only Keys

URI: <https://cookbookcosmosdb.documents.azure.com:443/>

PRIMARY KEY

SECONDARY KEY

PRIMARY CONNECTION STRING

AccountEndpoint=https://cookbookcosmosdb.documents.azure.com:443/;AccountKey= ..

SECONDARY CONNECTION STRING

AccountEndpoint=https://cookbookcosmosdb.documents.azure.com:443/;AccountKey= ..

cookbookblobstorage1 | Access keys

Storage account

Search (Ctrl+ /) < Show keys Set rotation reminder Refresh

Networking Azure CDN Access keys Shared access signature Encryption Security

Data management Geo-replication Data protection Object replication Blob inventory

Access keys authenticate your applications' requests to this storage account. Keep your keys in a Key Vault, and replace them often with new keys. The two keys allow you to replace one while still

Remember to update the keys with any Azure resources and apps that use this storage account.

Storage account name cookbookblobstorage1

key1
Last rotated: 3/15/2021 (125 days ago)
Rotate key

Key
.....

Connection string
.....

VehicleInformationDB ([REDACTED]) | Conn...

SQL database

Search (Ctrl+ /) < ADO.NET JDBC ODBC PHP Go

Power Automate (preview)

Settings Compute + storage Connection strings Maintenance Properties Locks

JDBC (SQL authentication)
`jdbc:sqlserver://[REDACTED].svr.database.windows.net:1433;database=VehicleInformationDB;user=sqladmin in @vehicleinformationsvr;password=[REDACTED];encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.database.windows.net;loginTimeout=30;`

Download JDBC driver for SQL server

sqldpool (synapsedemoworkspace11/sqldpool) | Connection strings

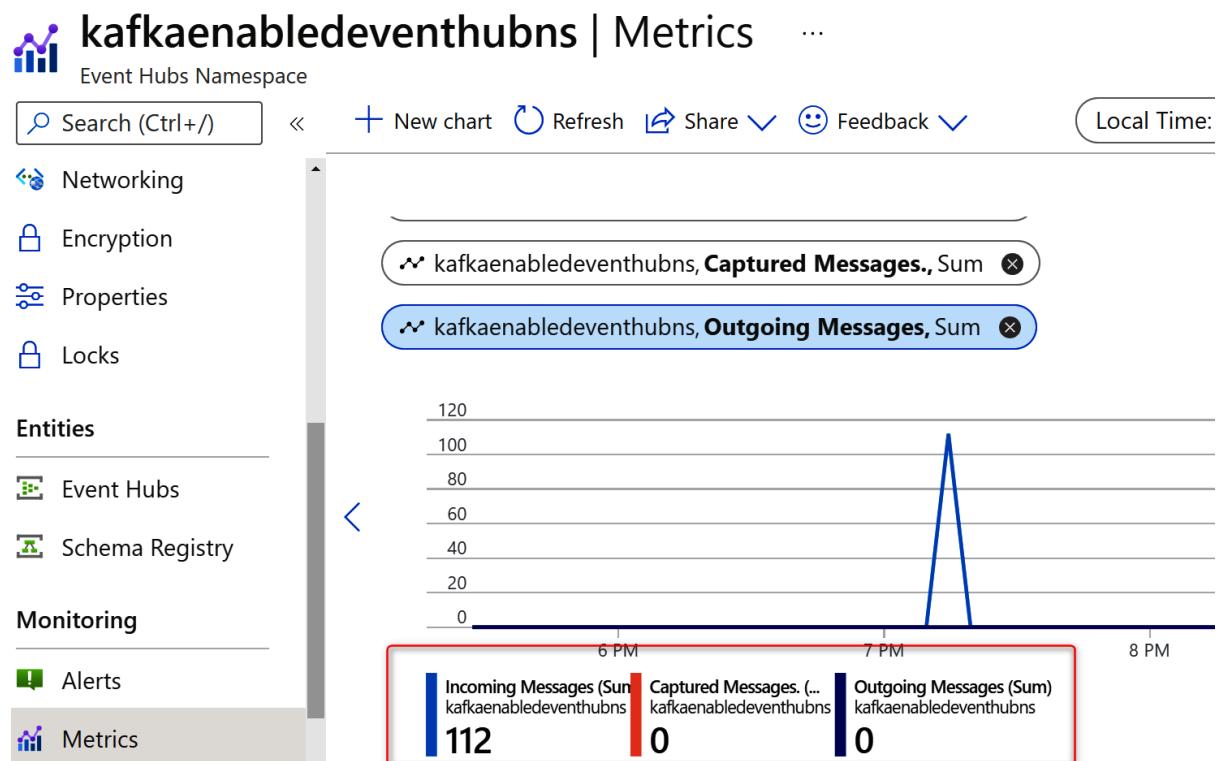
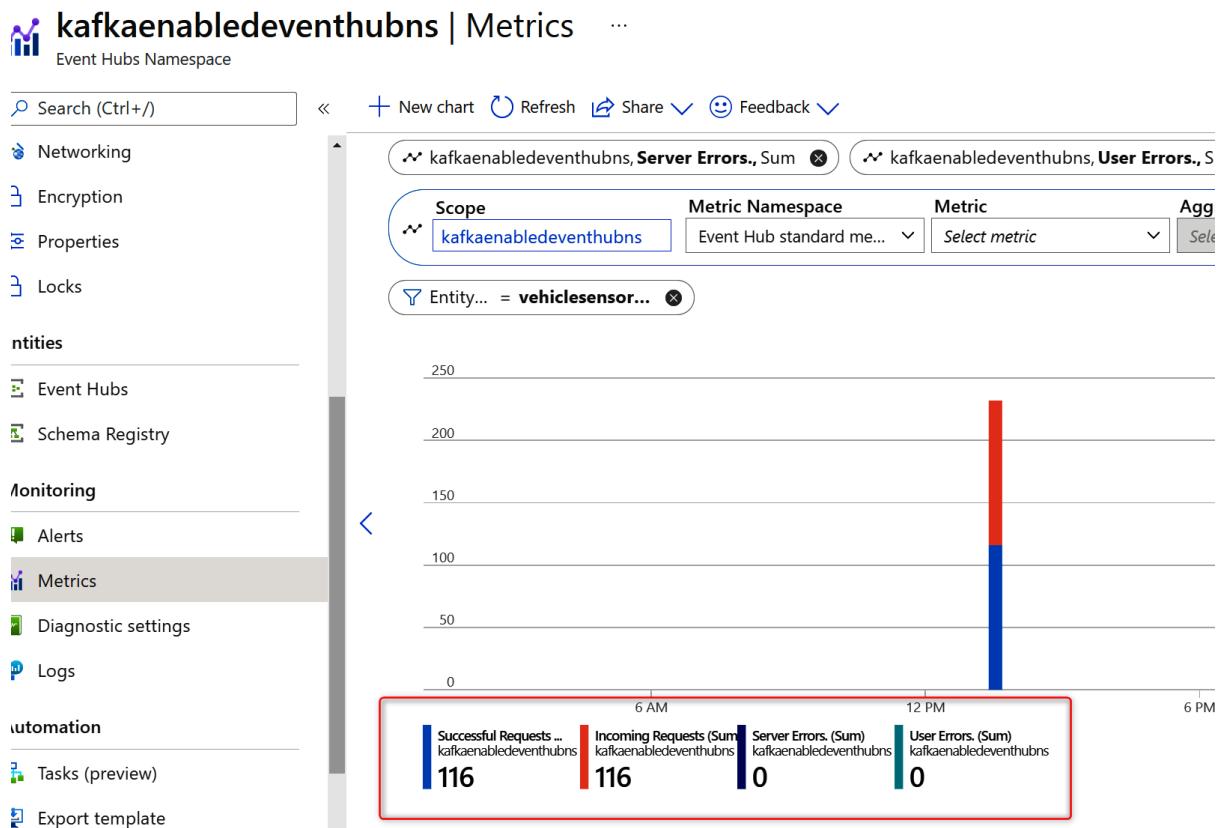
Dedicated SQL pool

Search (Ctrl+ /) < ADO.NET JDBC ODBC PHP

Workload management Maintenance schedule Geo-backup policy Connection strings Properties Locks

JDBC (SQL authentication)
`jdbc:sqlserver://synapsedemoworkspace11.sql.azuresynapse.net:1433;database=sqldpool;user=sqladminuser@synapsedemoworkspace11;password=[REDACTED];encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.sql.azuresynapse.net;loginTimeout=30;`

JDBC (Active Directory password authentication)





cookbookstoragegen2 | Containers

Storage account

Search (Ctrl+/)

«

+ Container



Change access le

Overview

Activity log

Tags

Diagnose and solve problems

Access Control (IAM)

Search containers by prefix

Name

rawdata

sensordata

sensordata ...

Container

Search (Ctrl+/)

«

Upload Add Directory Refresh | Rename

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Editor (preview)

Authentication method: Access key ([Switch to Azure AD User Account](#))
Location: sensordata / historicalvehicledata / Hist / data

Search blobs by prefix (case-sensitive)

Name

Modified

[...]

_spark_metadata

part-00000-1a8fb63d-6437-42bf-aff8-eca... 5/7/2021, 12:00:00 AM

part-00000-1d0a550a-030c-471c-bd8d-fe... 5/7/2021, 12:00:00 AM

part-00000-259de091-2238-4103-bf9b-38... 5/7/2021, 9:20:00 AM

part-00000-28e6902c-6eb5-4cf9-902d-50... 5/7/2021, 12:00:00 AM

 cookbookadresource

Data factory (V2)

Search (Ctrl+ /) Delete

Overview Location: East US

Activity log Subscription (change)

Access control (IAM) Subscription ID

Tags

Diagnose and solve problems

Getting started

Open Azure Data Factory Studio

Start authoring and monitoring your data pipelines and data flows.

Open 

Read documentation

Learn how to be productive quickly. Explore concepts, tutorials, and samples.

Learn more 

Networking

Managed identities

Properties

Locks

vehicleinformationdb | Data Explorer

Azure Cosmos DB account

Search (Ctrl+ /) New Container  

Quick start

Notifications

 Data Explorer

Settings

Features

Replicate data globally

Default consistency

Backup & Restore

Firewall and virtual networks

SQL API 

DATA

vehicle

Scale

vehicleinformationdb

Items

Settings

Stored Procedures

User Defined Function...

Triggers

sqldwpool (synapsesdemoworkspace11/sqldwpool) | Connection strings

Dedicated SQL pool

Search (Ctrl+ /) <>

Geo-backup policy

Connection strings (highlighted with a red box)

Properties

Locks

Security

ADO.NET JDBC ODBC PHP

JDBC (SQL authentication)

```
jdbc:sqlserver://synapsesdemoworkspace11.sql.azuresynapse.net:1433;database=sqldwpool;user=sqladminuse
=(your_password_here);encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.sql.azuresynapse.
```

vehicleinformationdb | Data Explorer

Azure Cosmos DB account

SQL API

New Item Update Discard Delete Upload Item

Items

DATA

SELECT * FROM c Edit Filter

vehicle

- Scale
- ▼ vehicleinformation
- Items
- Settings
- Stored Procedures
- User Defined Functions
- Triggers

NOTEBOOKS

- Gallery
- My Notebooks
- 2. Visualization.ipynb

	id	/id
b51480fc-af07-491...	b51480fc-af07-491d-b9...	
288e3ee8-ea29-48...	288e3ee8-ea29-489b-9c...	
90cce91-416e-45...	90cce91-416e-453f-b5...	
4b1e829f-5f32-4ae...	4b1e829f-5f32-4ae9-b4...	
2832b3ec-222c-40...	2832b3ec-222c-4049-9a...	
c79935cc-0b88-44...	c79935cc-0b88-44ae-97...	
115e6d02-4a43-41...	115e6d02-4a43-4131-b...	
5cb4ae8e-19de-40...	5cb4ae8e-19de-40e1-96...	
aa7e09d0-92cb-4c...	aa7e09d0-92cb-4cc0-99...	
04ac43cf-ed3c-4fb...	04ac43cf-ed3c-4fb7-a1...	

Load more

1 {
2 "Category": "Coupe",
3 "Hour": 3,
4 "eventtime": 1621655106170722,
5 "Make": "Chevrolet",
6 "source": "vehiclesensoreventhub",
7 "rpm": 46,
8 "speed": 94,
9 "lfi": 0,
10 "long": -73.949997,
11 "kms": 7794,
12 "Month": 5,
13 "Borough": "BROOKLYN",
14 "Year": 2021,
15 "Model": "300 CE",
16 "id": "b51480fc-af07-491d-b914-fec46d7d7d47",
17 "Day": 22,
18 "lat": 40.650002,
19 "Location": "185 Erasmus Street",
20 "_rid": "5Z3ALXW6ngBAAAAAAA==",
21 "_self": "dbs/-5Z3AA==/colls/-5Z3ALXW6ng=/docs/",
22 "_etag": "\"00008f04-0000-0100-0000-609589000000",
23 "_attachments": "attachments/",
24 "_ts": 1620412672
25 }

SQL SERVER

... SELECT top 5 * Untitled-1 CTR

CONNECTIONS

synapsesdemoworkspace11

Tables

- dbo.Customer
- dbo.databricks_streaming_checkpoint_0e08b4d...
- dbo.databricks_streaming_checkpoint_1121578...
- dbo.databricks_streaming_checkpoint_312df0...
- dbo.databricks_streaming_checkpoint_474e0e3...
- dbo.Dim_Location
- dbo.Dim_VehicleMaster
- dbo.vehicleinformation
- dbo.vehicletripinformation
- Views
- Programmability
- Security

RESULTS

LocationId	Borough	Location	Latitude	Longitude
1	BRONX	1475 Thieriot A...	40.837048	-73.865433
33	STATEN ISLAND	630 Richmond ...	40.579021	-74.151535
39	MANHATTAN	339 East 12th S...	40.73061	-73.984016
48	BROOKLYN	185 Erasmus St...	40.650002	-73.949997
108	QUEENS	Thrilla, New York	40.742054	-73.769417

Id	VehicleId	Make	Model	Category	ModelYear
9	c79935cc-0b88...	Chrysler	300	Sedan	2017
22	5cb4ae8e-19de...	Saab	3-Sep	Hatchback, Co...	1999
29	4b1e829f-5f32...	BMW	1 Series	Coupe, Convert...	2013
31	288e3ee8-ea29...	Toyota	4Runner	SUV	2013
34	90cce91-416e...	Toyota	2500 Crew Cab	Pickup	2021

id	eventtime	rpm	speed	kms	lfi
b51480fc-af07-...	2020-10-19 03:...	74	86	1594	1
90cce91-416e...	2021-05-07 03:...	80	87	4741	1
90cce91-416e...	2021-05-07 03:...	80	87	4741	1

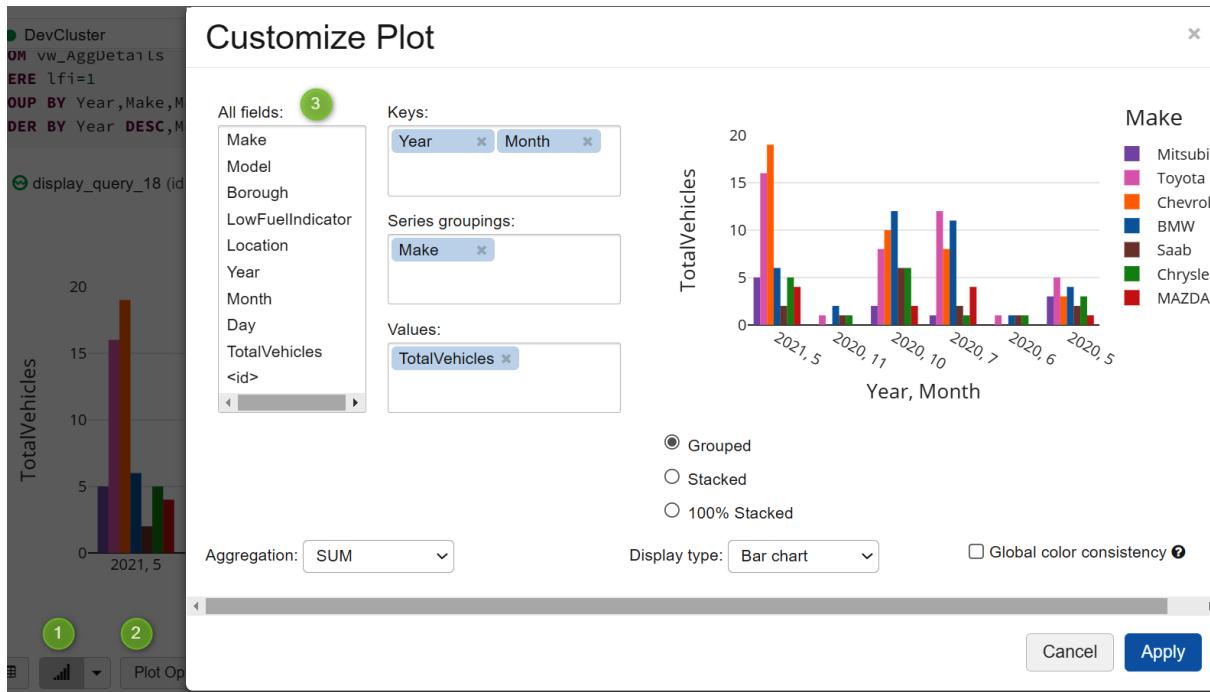
```

GROUP BY Year,Make,Model,coalesce(Borough,"UnKnown"),Location,
ORDER BY Year DESC,Month DESC,TotalVehicles DESC

```

► display_query_18 (id: fde636c1-8cfcd-405a-a44a-7bebf9d64f6a) Last update

	Make	Model	Borough	LowFuelIn
1	Mitsubishi	3000GT	QUEENS	1
2	Toyota	2500 Crew Cab	QUEENS	1
3	Toyota	4Runner	STATEN ISLAND	1
4	Toyota	2500 Crew Cab	BRONX	1
5	Chevrolet	300 CE	MANHATTAN	1
6	BMW	M Series	STATEN ISLAND	1
7	Mitsubishi	3000GT	BROOKLYN	1
8	Bar	Quantile	MANHATTAN	1
9	Scatter	Histogram	BRONX	1
10	Map	Box plot	BRONX	1
11	Line	Q-Q plot	STATEN ISLAND	1
12	Area	Pivot	MANHATTAN	1
Pie		Legacy charts ►		
Showing 1 - 12 of 12 rows				
 				



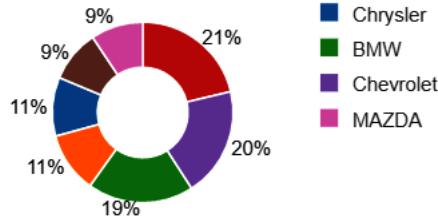
```
SELECT Make,Borough,Year,Month,COUNT(*) as TotalVehicles
FROM vw_AggDetails
GROUP BY Make,Borough,Year,Month
```

► (1) Spark Jobs

► ⚡ display_query_19 (id: 228f3487-e444-4c83-8e17-cd953db6137c) *Last updated*

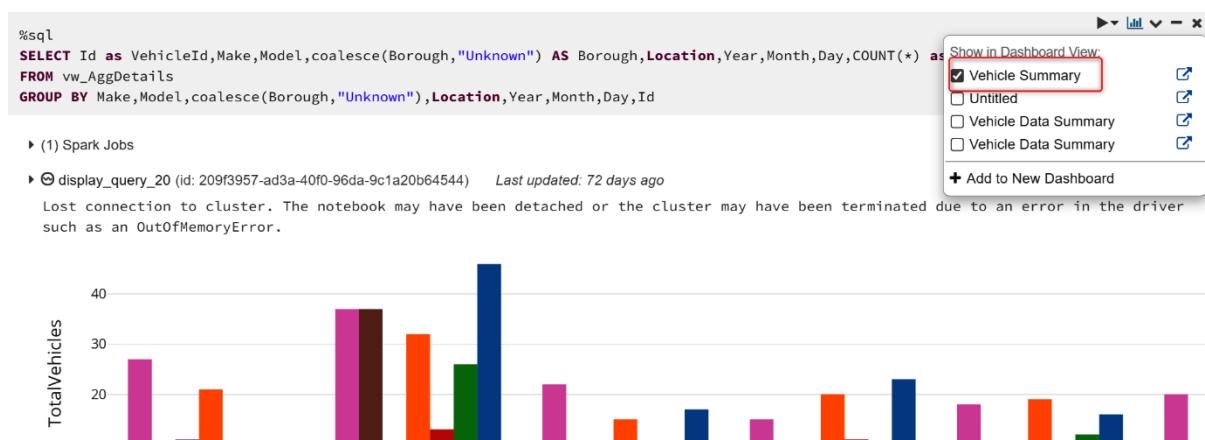
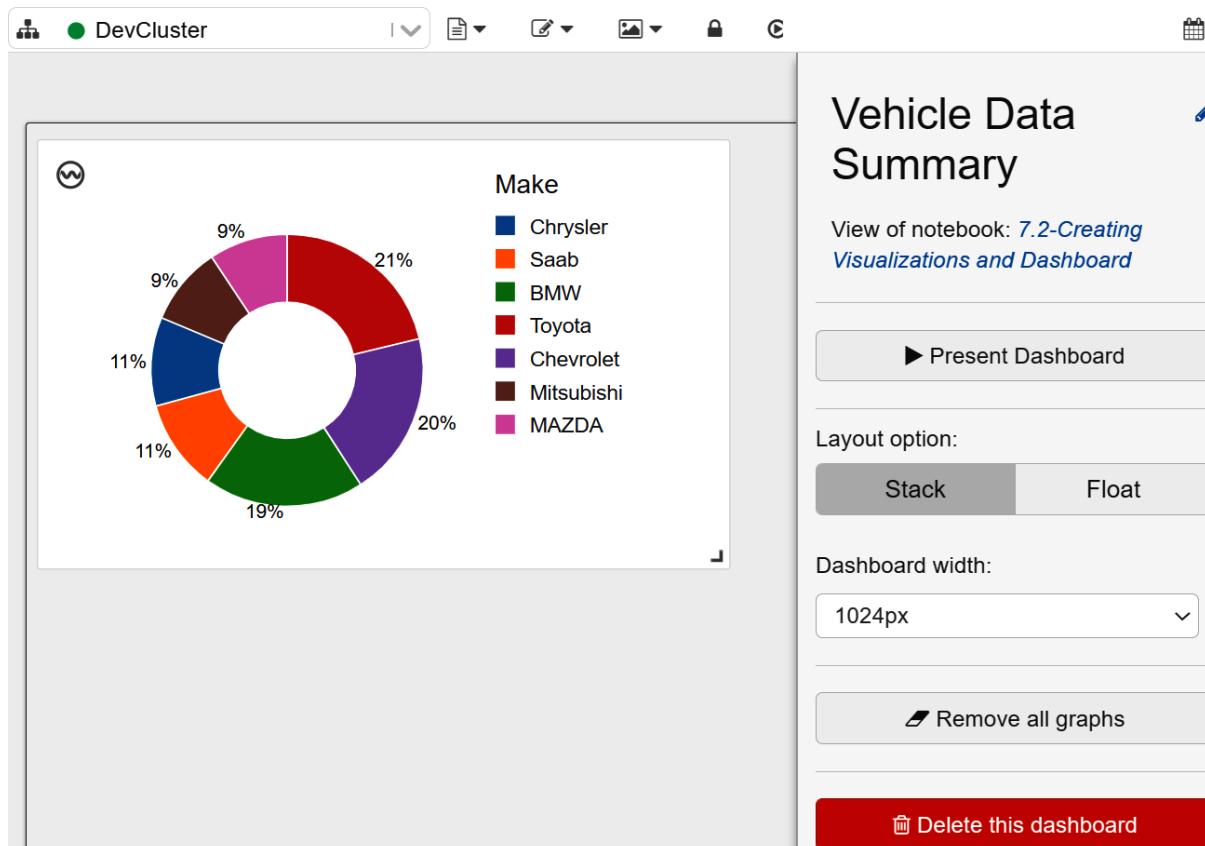
Lost connection to cluster. The notebook may have been detached or the driver has an error.

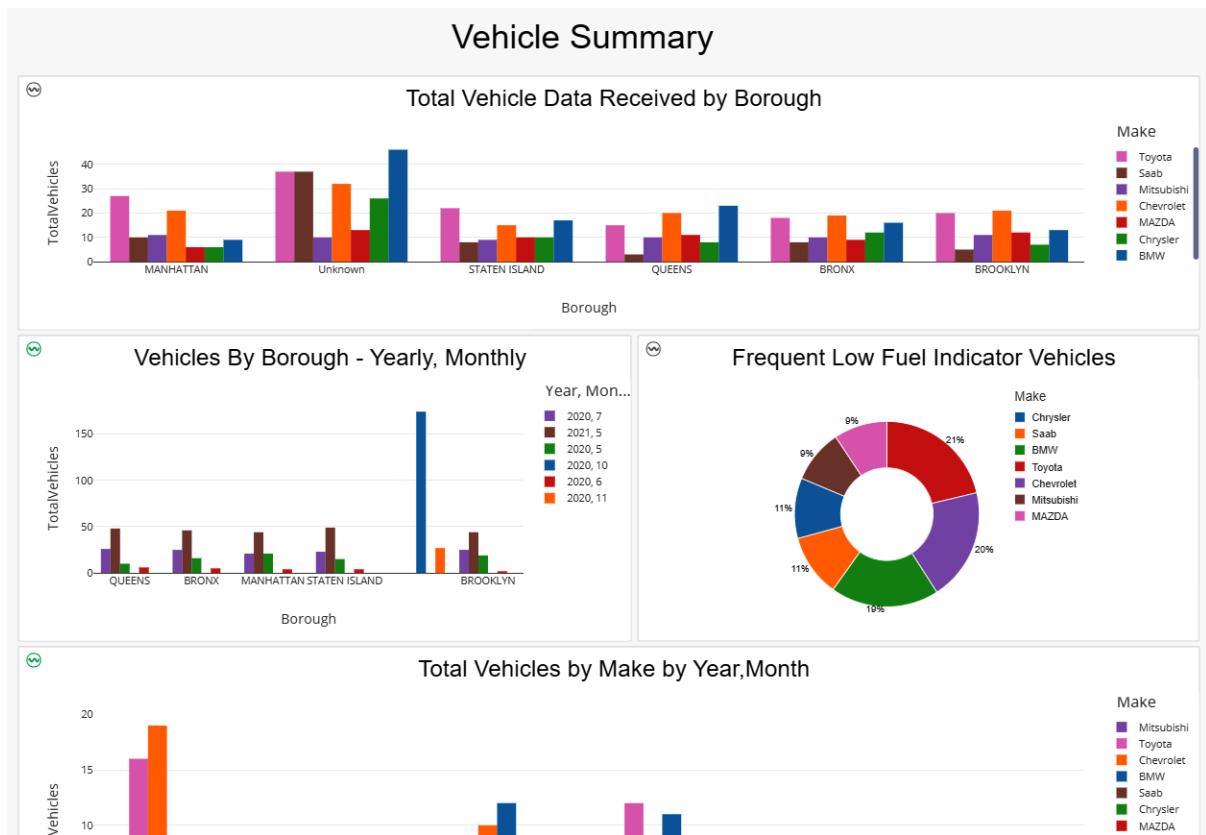
Make



7.2-Creating Visualizations and Dashboard (Python)

KafkaStreamingWork...





DevCluster


[Edit](#)
[Permissions](#)
[Clone](#)

[Configuration](#)
[Notebooks \(1\)](#)
[Libraries](#)
[Event Log](#)
[Spark UI](#)
[Driver Log](#)

Azure Data Lake Storage Credential Passthrough [?](#)

Enable credential passthrough for user-level data access

[Spark](#)
[Tags](#)
[SSH](#)
[Logging](#)
[Init Scripts](#)
[JDBC/ODBC](#)
[Permissions](#)

Server Hostname

adb- .1.azuredatabricks.net

Port

443

Protocol

HTTPS

HTTP Path

sql/protocolv1/o/ /0302-

Get Data

The screenshot shows the 'Get Data' interface with a search bar at the top. On the left, a sidebar lists categories: All, File, Database, Power Platform, Azure (which is selected and highlighted in yellow), Online Services, and Other. The main pane displays a list of Azure services, with 'Azure Databricks' highlighted by a red rectangle. At the bottom, there are buttons for 'Certified Connectors' and 'Template Apps', and 'Connect' and 'Cancel' buttons.

Azure Databricks

Server Hostname ⓘ

adb.

HTTP Path ⓘ

sql/protocolv1/o/

▲ Advanced Options (optional)

Database (optional) ⓘ

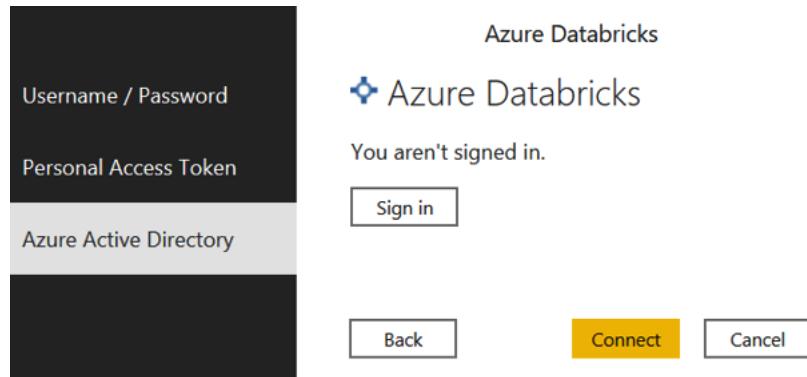
Example: abc

Batch Size (rows) (optional) ⓘ

Example: 10000

Data Connectivity mode ⓘ

Import DirectQuery



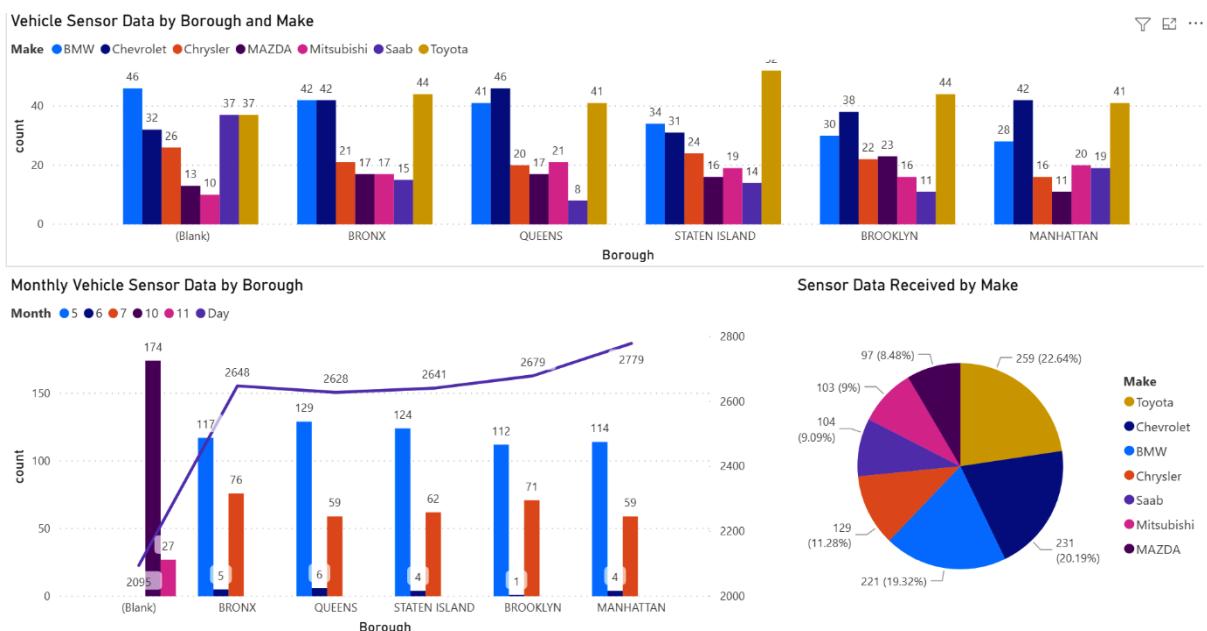
ADB

1.1.azuredatabricks.net:443

- SPARK [6]
- default
- deltadb
- office_hours
- tpchdb
- vehicle
- vehiclesensor [5]**
 - vehicledelta_bronze
 - vehicledelta_historical
 - vehicledelta_historical_test
 - vehicledelta_silver
 - vehicledeltaaggregated**

["start":2020-10-31 03:00:00,"end":2020-10-31 04:00:00]	Chrysler	
{"start":2021-05-23 22:00:00,"end":2021-05-23 23:00:00}	Toyota	STATEN ISLAND
{"start":2021-05-24 01:00:00,"end":2021-05-24 02:00:00}	Toyota	STATEN ISLAND
{"start":2021-05-08 03:00:00,"end":2021-05-08 04:00:00}	BMW	STATEN ISLAND
{"start":2021-05-09 03:00:00,"end":2021-05-09 04:00:00}	Chevrolet	BRONX
{"start":2020-06-01 19:00:00,"end":2020-06-01 20:00:00}	MAZDA	QUEENS
{"start":2021-07-18 13:00:00,"end":2021-07-18 14:00:00}	Toyota	BRONX
{"start":2021-05-24 04:00:00,"end":2021-05-24 05:00:00}	BMW	MANHATTAN
{"start":2020-05-22 19:00:00,"end":2020-05-22 20:00:00}	Chrysler	STATEN ISLAND
{"start":2021-05-23 22:00:00,"end":2021-05-23 23:00:00}	BMW	MANHATTAN
{"start":2021-07-18 23:00:00,"end":2021-07-19 00:00:00}	Mitsubishi	QUEENS
{"start":2021-05-11 03:00:00,"end":2021-05-11 04:00:00}	Saab	STATEN ISLAND
{"start":2021-05-07 03:00:00,"end":2021-05-07 04:00:00}	Chevrolet	MANHATTAN
{"start":2021-05-23 18:00:00,"end":2021-05-23 19:00:00}	BMW	MANHATTAN
{"start":2021-07-19 04:00:00,"end":2021-07-19 05:00:00}	BMW	QUEENS
{"start":2020-05-28 19:00:00,"end":2020-05-28 20:00:00}	Mitsubishi	BRONX
{"start":2021-07-18 15:00:00,"end":2021-07-18 16:00:00}	Saab	QUEENS
{"start":2020-05-27 19:00:00,"end":2020-05-27 20:00:00}	BMW	QUEENS
{"start":2021-05-08 03:00:00,"end":2021-05-08 04:00:00}	Chevrolet	BROOKLYN
{"start":2020-10-29 03:00:00,"end":2020-10-29 04:00:00}	Chevrolet	
{"start":2021-07-18 18:00:00,"end":2021-07-18 19:00:00}	Chrysler	BROOKLYN
{"start":2020-10-20 03:00:00,"end":2020-10-20 04:00:00}	Chevrolet	
{"start":2021-05-23 17:00:00,"end":2021-05-23 18:00:00}	BMW	BRONX

Load Transform Data Cancel



CookbookRG Resource group

Search (Ctrl+ /) Create Edit columns Delete resource group Refresh Export to CSV Open queue

Cost Management

- Cost analysis
- Cost alerts (preview)
- Budgets
- Advisor recommendations

Monitoring

- Insights (preview)
- Alerts
- Metrics
- Diagnostic settings
- Logs
- Advisor recommendations

Essentials

Subscription (change) Deployment
Subscription ID 1 Failed, 22 Succeeded
Tags (change) Location
Click here to add tags East US

Filter for any field... Type == all X Location == all X Add filter No grouping

Showing 1 to 28 of 28 records. Show hidden types Type ↑
 Name ↑ vehicleinformationdb Azure Cosmos DB account

Data Factory Validate all Publish all 1

Factory Resources

Filter resources by name +

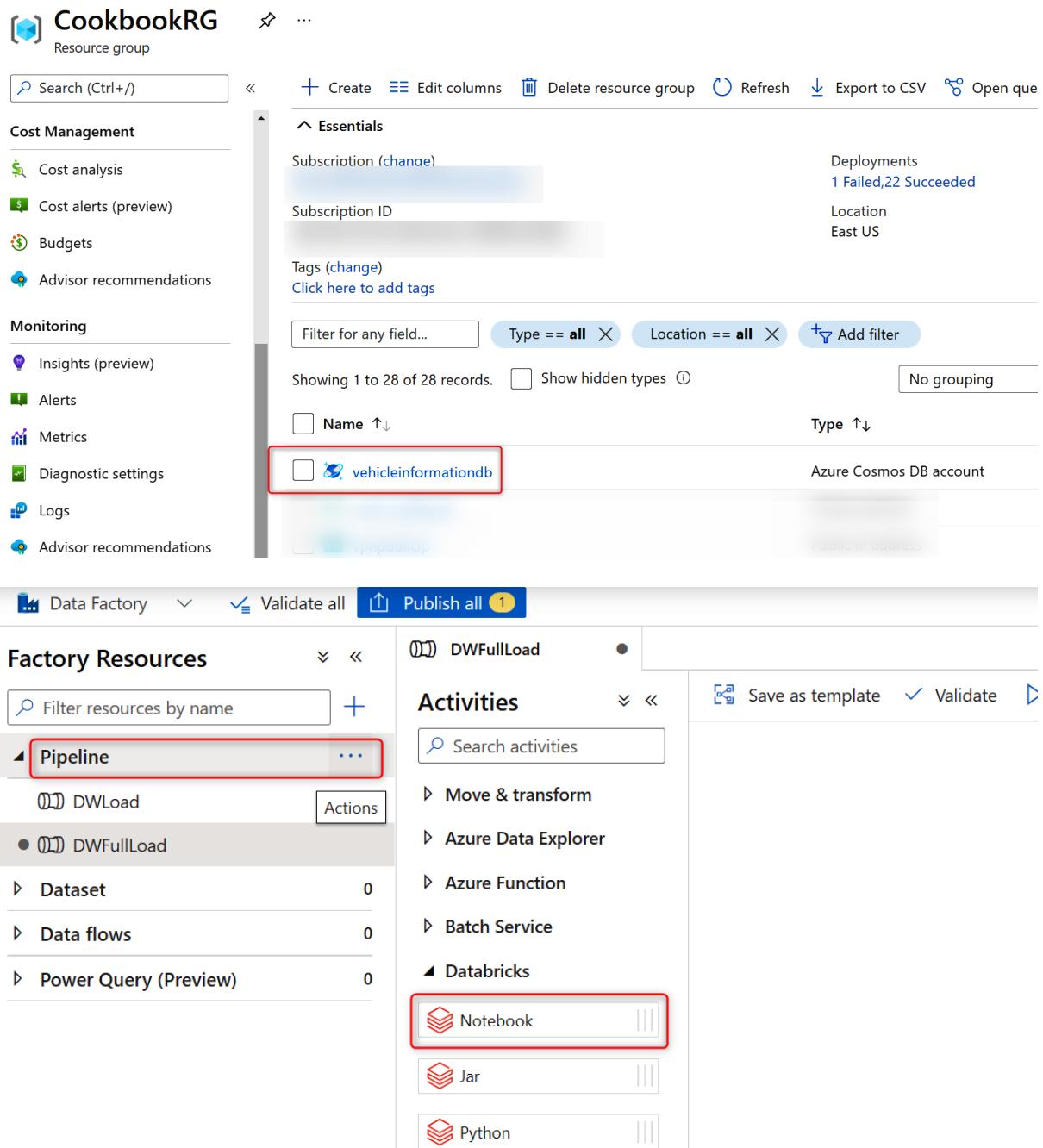
Pipeline ...

- DWLoad
- DWFullLoad**
- Dataset 0
- Data flows 0
- Power Query (Preview) 0

Activities

Save as template Validate

- Move & transform
- Azure Data Explorer
- Azure Function
- Batch Service
- Notebook**
- Jar
- Python



DWLoad

Save as template Validate Debug Add trigger

Connect via integration runtime * [?](#)
AutoResolveIntegrationRuntime

Account selection method *
Enter manually

Databrick Workspace URL * [?](#)
https://adb-... .azuredatabricks.net

Authentication type *
Access Token

Access token Azure Key Vault
Access token * [?](#)
.....

Select cluster
 New job cluster Existing interactive cluster Existing instance

Existing cluster ID * [?](#)

General Azure Databricks Settings User properties

Factory Resources

Filter resources by name +

▲ Pipeline 1

- DWLoad
- Dataset 0
- Data flows 0
- Power Query (Preview) 0

Activities

Search activities

- Move & transform
- Azure Data Explorer
- Azure Function
- Batch Service
- Databricks
- Data Lake Analytics
- General

DWLoad

Save as template Validate Debug Add trigger

```

graph LR
    A[Notebook: LoadingFromSilverTables] --> B[Stored procedure: Load in Facts and Dims]
  
```

Chapter 8: Azure Databricks SQL Analytics

The screenshot shows the Microsoft Azure Databricks portal interface. On the left, a dark sidebar menu includes options like Data Science & Engineering, Machine Learning, and SQL, with SQL highlighted by a red box. Other items like Recents and Search are also visible. The main area features a large "Azure Databricks" logo and a "Explore the Quickstart Tutorial" section with a brief description of how to set up a cluster, run queries, and display results. A "New Notebook" button is located at the bottom of the sidebar.

The screenshot shows the "Admin Console" section of the Azure Databricks portal. It displays the user's sign-in information ("Signed in as [redacted].com") and a "User Settings" menu. The "Admin Console" option is highlighted by a red box. Other menu items include Manage Account and Log Out. Below this is a "Workspaces" section with a checked item for "ADBCookBook".

Add User

Add a user to Databricks. You can add any user who belongs to the Azure Active Directory tenant of your Azure Databricks workspace.

Email

CancelOK

Admin Console

[Users](#) [Groups](#) [Global Init Scripts](#) [Workspace Settings](#)[+ Add User](#)

Username	Name	Admin	Workspace access	Databricks SQL access
user@xyz.com		<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Admin Console

[Users](#) [Groups](#) [Global Init Scripts](#) [Workspace Settings](#)[+ Add User](#)

Username	Name	Admin	Workspace access	Databricks SQL access	Allow cluster creation
user@xyz.com		<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Microsoft Azure | Databricks

Portal

The sidebar menu includes:

- SQL
- Create
- Dashboards
- Queries
- Alerts
- Data
- SQL Endpoints (with a red notification dot)

The main screen shows a table for managing SQL Endpoints:

Actions	Active / Max	Size	State
Start	0 / 1	2X-Small	Stopped

Buttons include:

- Create SQL Endpoint
- Start

New SQL Endpoint

Preview

X

Name:

DemoSQLEndpoint

Cluster Size i :

2X-Small

4 DBU ▾

Auto Stop:

After

60

minutes of inactivity.

Multi-cluster Load Balancing i :
[Preview](#)

Cluster Count: Min

1

Max

1

Photon i :

Tags i :

Key

Value

Cancel

Create

SQL Endpoints

[Create SQL Endpoint](#)

Search SQL endpoints...

Name	State	Size	Active / Max	Actions
DemoSQLEndpoint	<input checked="" type="checkbox"/> Running	2X-Small	1 / 1	Stop ⋮

{ } { }

☰

⚡

LIMIT 1000

[Save](#)

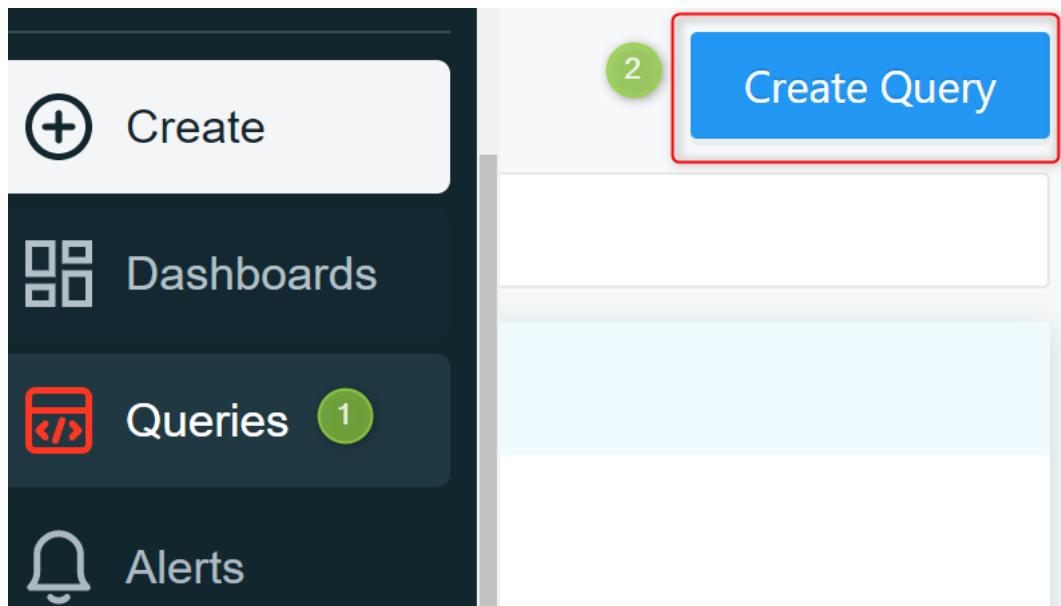
[Run](#)

```
1 SHOW GRANT `user@xyz.com` ON TABLE default.customer;
```

Table

[+ Add Visualization](#)

Principal	ActionType	ObjectType	ObjectKey
user@xyz.com	READ_METADATA	TABLE	`default`.`customer`
user@xyz.com	SELECT	TABLE	`default`.`customer`
user@xyz.com	USAGE	DATABASE	default



```
1 select COUNT(*) as TotalCustomer,C_MKTSEGMENT
2 group by C_MKTSEGMENT
3 order BY COUNT(*) desc
4
5 |
```

Table : Count(*),C_MKTSEGMENT

TotalCustomer C_MKTSEGMENT

- + Add to Dashboard
- Download as CSV File
- Download as TSV File
- Download as Excel File

: Edit Visualization 4 rows

Customer Query

tpch



{ } { }

≡



DemoSQLendpoint ✓

default

Filter tables & columns...

- customer
- customer_maxfiles
- customerdelta
- customerdeltapartition
- customersegments
- customersinglefile

```
1 select COUNT(*) as
2 group by C_MKTSEGME
3 order BY COUNT(*) c
4
5 default.customerde
```

Table : Counter

TotalCustomer C_MK

- Dashboards
- Queries
- Alerts
- Data
- SQL Endpoints
- Query History
- Help
- Settings

Query History

Me

com

Last 14 days



All SQL endpoints

Status

Query	SQL Endpoi...	Started At	Duration	User
Listing columns 'catalog : null, sch...	DemoSQLe...	2021-07-21 15:06	982 ms	
Listing tables 'catalog : null, sche...	DemoSQLe...	2021-07-21 15:06	81 ms	
show tables in default -- user_id: {}	DemoSQLe...	2021-07-21 15:06	214 ms	
describe database extended default -...	DemoSQLe...	2021-07-21 15:06	278 ms	
show grant on DATABASE default -- us...	DemoSQLe...	2021-07-21 15:06	290 ms	
show databases -- user_id: {}	DemoSQLe...	2021-07-21 15:06	232 ms	

✓ show databases -- user_id...	De
✓ select count(*) from cust...	De
✓ select count(*) from orde...	De
✓ select * from orders LIMI...	De
✓ Listing columns 'catalog ...	De
✓ Listing tables 'catalog : ...	De
✓ SHOW DATABASES -- user_id...	De
✗ select * from orders LIMI...	De
✗ select * from orders LIMI...	De
✗ select * from orders LIMI...	De
✗ select * from orders LIMI...	De
✗ select * from democustome...	De

Overview Execution Details

```

1 select
2   count(*)
3   from
4     customer
5   LIMIT
6   1000

```

ID	e2	104-4980ce
Status	✓ Finished	
Start time	2021-07-21 15:04:01.945	
End time	2021-07-21 15:04:06.627	
Duration	3.56 s	
User		
SQL Endpoint	DemoSQLEndpoint	

Customer Query tpch



dd description



⌚ Refresh Schedule

Never

```

1 select COUNT(*) as Total
2 group by C_MKTSEGMENT
3 order BY COUNT(*) desc
4
5 default.customerdefault

```

Refresh Schedule

Refresh every

Never ▾

Minutes

Schedule

Cancel

OK

Never

1 minute

Query History

Query	SQL Endpoint	Details
⑤ select * from customer limit 100000	DemoSQLe...	Open ↗
⑤ select * from customer LIMIT 1000	DemoSQLe...	
⑤ Listing columns 'catalog : null, sch...	DemoSQLe...	
⑤ Listing tables 'catalog : null, sche...	DemoSQLe...	
⑤ select * from customer LIMIT 1000	DemoSQLe...	
⑤ Listing columns 'catalog : null, sch...	DemoSQLe...	
⑤ Listing tables 'catalog : null, sche...	DemoSQLe...	
⑤ select * from customer LIMIT 1000	DemoSQLe...	
⑤ Listing columns 'catalog : null, sch...	DemoSQLe...	
⑤ Listing tables 'catalog : null, sche...	DemoSQLe...	
⑤ select * from customer LIMIT 1000	DemoSQLe...	
⑤ Listing columns 'catalog : null, sch...	DemoSQLe...	
⑤ Listing tables 'catalog : null, sche...	DemoSQLe...	
⑤ Listing columns 'catalog : null, sch...	DemoSQLe...	
⑤ Listing tables 'catalog : null, sche...	DemoSQLe...	
⑤ Listing columns 'catalog : null, sch...	DemoSQLe...	
⑤ Listing tables 'catalog : null, sche...	DemoSQLe...	
⑤ show tables in default -- user id: {}	DemoSQLe...	

Execution Details

Duration 5.07 s 100%
 ● Loading metadata & optimizing 383 ms 8%
 ● Execution 2.13 s 42%
 ● Result fetching 2.56 s 51%

Rows returned 100,000

IO
 Rows read 135,000
 Bytes read 10.58 MB
 Bytes read from cache 0 %
 Bytes written 0 bytes

Files & Partitions
 Files read 3
 Partitions read 0

ID 3ccf8e6b-364c-421c-a9ee-cea...
 Status Finished
 Start time 2021-07-21 15:23:15.478
 End time 2021-07-21 15:23:21.663
 Duration 5.07 s
 User [REDACTED]
 SQL Endpoint DemoSQLendpoint
 Details [Open ↗](#)

Details for Query 52

 Download screen as png

Submitted Time: 2021/07/21 09:53:17

Duration: 2 s

Succeeded Jobs: 12 13

Expand all the details in the query plan visualization

Scan parquet default.customer +details

Stages: 20.0 21.0

cache writes size (uncompressed) total (min, med, max)	20.5 MiB (6.8 MiB, 6.8 MiB, 6.8 MiB)
time spent waiting fetching data from cloud storage total (min, med, max)	592 ms (163 ms, 213 ms, 215 ms)
time spent in the cache locality manager in milliseconds total (min, med, max)	5 ms (0 ms, 0 ms, 5 ms)
number of files read	3
filesystem read data size total (min, med, max)	10.6 MiB (3.5 MiB, 3.5 MiB, 3.5 MiB)
cache async file status fetch waiting time total (min, med, max)	0 ms (0 ms, 0 ms, 0 ms)
scan time total (min, med, max)	763 ms (212 ms, 263 ms, 288 ms)
filesystem read data size (sampled) total (min, med, max)	21.2 MiB (7.0 MiB, 7.0 MiB, 7.1 MiB)
filesystem read time (sampled) total (min, med, max)	824 ms (268 ms, 276 ms, 280 ms)
metadata time	0 ms
size of files read	10.6 MiB
cache hits size total (min, med, max)	0.0 B (0.0 B, 0.0 B, 0.0 B)

```
1  SELECT
2    COUNT(*),
3    SUM(count) AS TotalCountofVehicles
4  FROM
5    vehiclesensor.vehicledeltaaggregated
6  WHERE Borough={{ BoroughName }} and Make ={{ MakeValue }}
7
8
```

BoroughName



MakeValue



Query Parameter Examples

1

```
1 SELECT
2   COUNT(*),
3   SUM(count
4 FROM
5   vehiclese
6 WHERE Boro
7
8
9 -- SELECT
10 -- COUNT(
11 -- case
```

BoroughName
STATEN ISLAN

Add Parameter

X

* Keyword:

Choose a keyword for this parameter

* Title:

2

Type: Text

▼

Cancel

Add Parameter

```
1 SELECT
2   COUNT(*) AS TotalVehicles,
3   case
4     when lfi = 0 then "Low"
5     else "High"
6   end as LowFuelIndicator
7 FROM
8   vehiclesensor.vehicledelta_silver
9 WHERE
10  eventtime BETWEEN '{{ Date Range.start }}' AND '{{ Date Range.end }}'
11  group BY LowFuelIndicator
```

Date Range



2021-05-23

→ 2021-05-26



Table

⋮

TotalVehicles	LowFuelIndicator
83	High
211	Low

```
8     vehiclesensor.vehicledelta_silver
9     WHERE
10    eventtime BETWEEN '{{ Date Range.start }}' AND
11    group BY LowFuelIndicator
```

Date Range

2021-05-23 → 2021-05-26

This week Jul 18 - Jul 24

This month July

This year 2021

Last week Jul 11 - Jul 17

Last month June

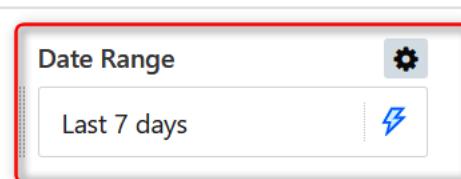
Last year 2020

Last 7 days Jul 14 - Today

Last 14 days Jul 7 - Today

Last 30 days Jun 21 - Today

Last 60 days May 22 - Today



Table

TotalVehicles	LowFuelIndicator
49	High
158	Low

Boroughs

X

* Title: Boroughs

Type: Query Based Dropdown List

Query: Boroughs

X

Select query to load dropdown values from

Allow multiple values

Quotation: Double Quotation Mark

V

Placed in query as: "value1","value2","value3"

Cancel

OK

```
12    Borough in ({{ Boroughs }})  
13    group BY LowFuelIndicator,Borough
```

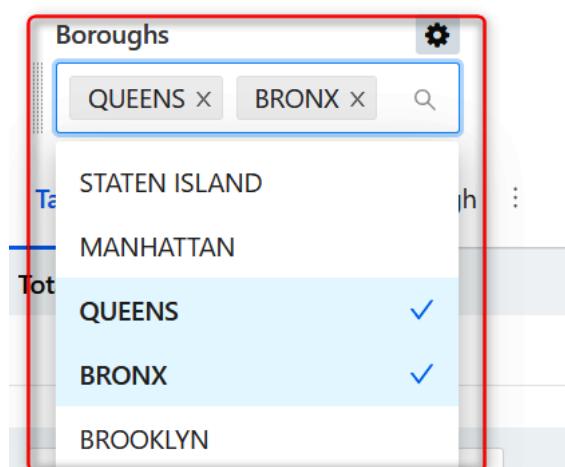


Table :

Borough

BRONX x

TotalVehicles	LowFuelIndicator	Borough
58	High	BRONX
140	Low	BRONX

Clear

Select All

Values

BRONX ✓

BROOKLYN

MANHATTAN

QUEENS

STATEN ISLAND

BRONX x

TotalVehicles LowFuelIndicator

58 High

140 Low

Boroughs



QUEENS x

BRONX x

MANHATTAN x



Table :

+ Add Visualization

TotalVehicles LowFuelIndicator Borough

140 Low BRONX

58 High BRONX

General X Axis Y Axis Series Colors Data Labels

Chart Type

Bar

Horizontal Chart

X Column

Borough

Y Columns

TotalVehicles

Group By

LowFuelIndicator

Visualization Name

LFI vehicles by Borough

General X Axis Y Axis Series Colors **Data Labels**

Show Data Labels

Number Values Format

0,0[,]00000

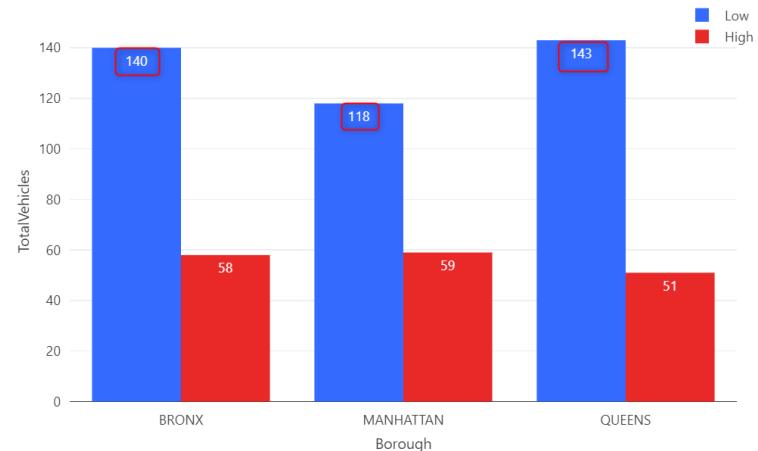
Percent Values Format

0[,]00%

Date/Time Values Format

YYYY-MM-DD HH:mm

Data Labels



Visualization Type

Funnel

Steps	Value	% Max	% Previous
QUEENS	143	100%	100%
BRONX	140	97.90%	97.90%
MANHATTAN	118	82.52%	84.29%
MANHATTAN	59	41.26%	50%
BRONX	58	40.56%	98.31%
QUEENS	51	35.66%	87.93%

General Appearance

Step Column

Borough

Step Column Title

Steps

Value Column

TotalVehicles

Boroughs 

QUEENS X BRONX X MANHATTAN X 

Table : LFI vehicles by Borough : **Funnel** :

Steps	Value
-------	-------

Add Visualization Widget 

Customer Query 

Choose Visualization

Table 

TABLE
Table 

COUNTER
Counter

CHART
Chart

Cancel **Add to Dashboard**

Query Parameter Examples

X

Choose Visualization

LFI vehicles by Borough

V

Title

LFI vehicles by Borough - Query Parameter Examples

...

Description

(empty text area)

Parameters

Title	Keyword	Default Value	Value Source
Boroughs 	<code>{{ Boroughs }}</code>	QUEENS, BRONX	Dashboard 

Cancel

Add to Dashboard

★ Vehicle and Customer Stat

Share

Schedule 

Refresh 

:

Boroughs

QUEENS X BRONX X MANHATTAN X 

Counter - Vehicle Agg

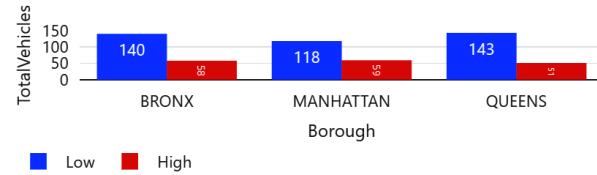
652

Total Count

 2 months ago

LFI vehicles by Borough - Query Parameter Examples



 just now

Table - Customer Query

TotalCustomer	C_MKTSEGMENT
30,189	HOUSEHOLD
30,142	BUILDING
29,968	FURNITURE
29,949	MACHINERY

Line Chart - Query Parameter Examples



Vehicle and Customer Stat

Share **Schedule** Refresh

Boroughs: QUEENS X BRONX X MANHATTAN X

Refresh: Select interval

Never

Every 1 minute

Every 5 minutes

Every 10 minutes

Every 30 minutes

Every 1 hour

Every 12 hours

SQL Endpoint : Counter - Vehicle Agg

652 Total Count

QUEENS 143

Overview **Connection Details** Monitoring

Server Hostname: adb-[REDACTED] 43.3.azuredatabricks.net

Port: 443

Protocol: https

HTTP Path: /sql/1.0/endpoints/

JDBC URL: jdbc:spark://adb-[REDACTED].azuredatabricks.net:443/default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/endpoints/84;

OAuth URL: https://adb-[REDACTED].azuredatabricks.net/oidc

Use these details to connect to this endpoint from BI tools

Create a [personal access token](#)

Azure Databricks

Server Hostname ⓘ

adb-.3.azuredatabricks.net

HTTP Path ⓘ

/sql/1.0/endpoints/{

► Advanced Options (optional)

Database (optional) ⓘ

Example: abc

Batch Size (rows) (optional) ⓘ

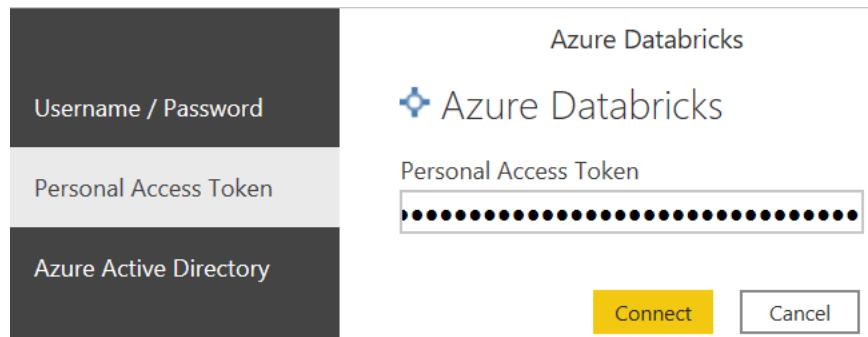
Example: 10000

Data Connectivity mode ⓘ

Import

DirectQuery

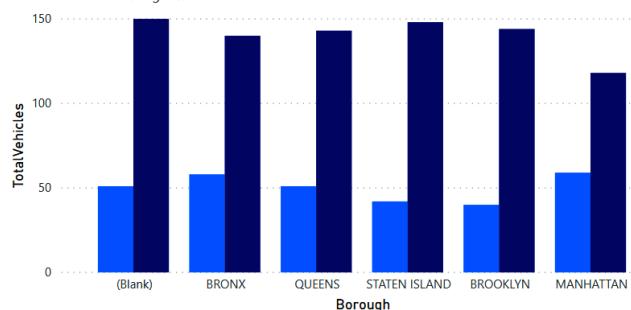
OK Cancel



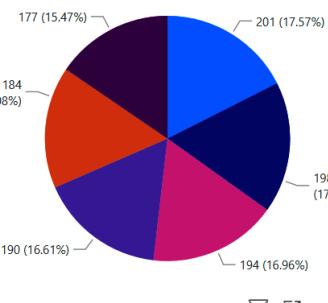
TotalVehicles	LowFuelIndicator	Borough
148	Low	STATEN ISLAND
51	High	QUEENS
118	Low	MANHATTAN
42	High	STATEN ISLAND
58	High	BRONX
40	High	BROOKLYN
140	Low	BRONX
59	High	MANHATTAN
143	Low	QUEENS
144	Low	BROOKLYN
51	High	null
150	Low	null

TotalVehicles by Borough and LowFuelIndicator

LowFuelIndicator ● High ● Low

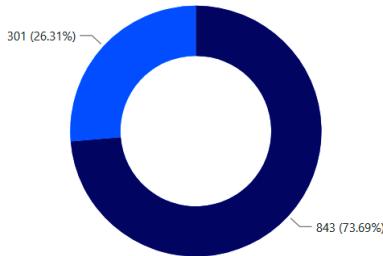


TotalVehicles by Borough



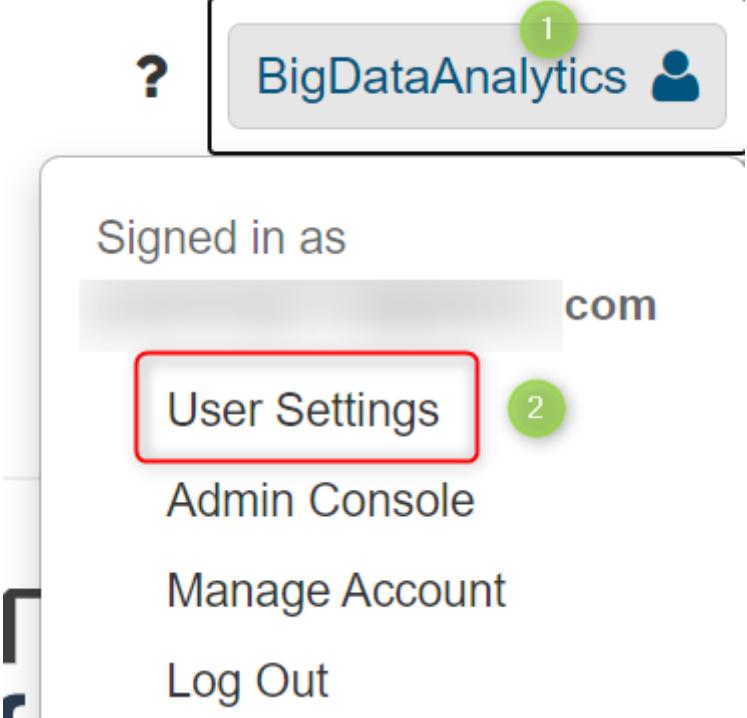
Borough
● (Blank)
● BRONX
● QUEENS
● STATEN ISLAND
● BROOKLYN
● MANHATTAN

TotalVehicles by LowFuelIndicator



LowFuelIndicator
● Low
● High

Chapter 9: DevOps Integrations and Implementing CI/CD for Azure Databricks



The screenshot shows the Azure Databricks User Settings page. At the top, there is a navigation bar with a question mark icon, the text "BigDataAnalytics", and a user profile icon. A green circle with the number "1" is positioned above the user profile icon. Below the navigation bar, the text "Signed in as" is followed by a blurred email address ending in ".com". A green circle with the number "2" is positioned next to the blurred email address. A red box highlights the "User Settings" button. Below the "User Settings" button are several other options: "Admin Console", "Manage Account", and "Log Out".

Git provider

Azure DevOps Services

Token or app password (optional) ?

Token or app password with repo read/write perm iss

Git provider username or email (optional) ?

Save

Azure DevOps

CookbookProj / Repos / Files / CookbookRepo

CookbookProj

Overview Boards Repos Files Commits Pushes

CookbookRepo

M README.md

main / Type to find a file

Files

Contents History

Name ↑

M README.md

Introduction

BigDataAnalytics

Jul 24 2021, 8:12 AM IST. Exit

Git: Not linked 2

1

Jul 24 2021, 8:12 AM IST

admin@demomoveorg.onmicrosoft.com

All changes saved Save now

.windows.net/".format(storageAccount)

lication(ADLSGen2App) was have created as

The screenshot shows the Azure DevOps interface for a project named 'CookbookProj'. In the center, a repository named 'CookbookRepo' is displayed, containing a single file 'README.md'. The left sidebar provides navigation links for Overview, Boards, Repos, Files, Commits, and Pushes. The bottom of the screen features a status bar with the date and time ('Jul 24 2021, 8:12 AM IST'), a message about Git being unlinked, and a notification for saved changes. A red box highlights the 'Git: Not linked' message, and a green circle with the number '1' highlights a notification for saved changes.

The screenshot shows the Azure DevOps interface. At the top, the URL is https://dev.azure.com/DemoCookbookOrg/CookbookProj/_git/CookbookRepo. The page title is 'CookbookProj' and the repository name is 'CookbookRepo'. The left sidebar includes links for Overview, Boards, Repos (with a green notification badge '1'), and Files. The main content area shows a file named 'README.md'.

Git Preferences

Status Link Unlink

Link https://dev.azure.com/DemoCookbookOrg/CookbookProj/_

Branch

Path in Git Repo [notebooks/6.1-Reading Writing to Delta Tables.py](#)

Tip: You can also import/export multiple notebooks or an entire folder through [Workspace API](#) to your computer and check-in to your favorite version control system.

Save Notebook Revision

The file you linked to does not exist on Git yet. Would you like to make a commit and save the current version?

First Commit

Close

Save



BigDataAnalytics



Git: Synced



Jul 24 2021, 8:17 AM IST



admin@demomoveorg.onmicrosoft.com
Commit fd8cc6c982

First Commit

All changes saved [Save now](#)

Azure DevOps interface showing the 'Pull requests' page for the 'CookbookProj' repository. The sidebar shows 'CookbookProj' selected under 'Repos'. The main area displays a single pull request with the following details:

- Pull requests** (button)
- Mine**, **Active**, **Completed**, **Abandoned** (filter buttons)
- New pull request** (button)
- Customize view** (button)
- You updated 3p Dev 6m ago** (notification)
- Create a pull request** (button)

: / Repos / Files / CookbookRepo

Search

CookbookRepo

notebooks

PY 6.1-Reading Writing to Delt...

M README.md

Dev / notebooks

notebooks

Contents History

Name ↑	Last change	Commits
PY 6.1-Reading Wri...	11m ago	fd8cc6

Adding new comment

```
storageAccount="cookbookadlsgen2storage"
mountpoint = "/mnt/Gen2Source"
storageEndPoint ="abfss://rawdata@{}.dfs.core.windows.net"
print ('Mount Point =' +mountpoint)
```

Git: Synced



Jul 24 2021, 8:33 AM IST



admin@demomoveorg.onmicrosoft.com

[Commit d050d4b2e3](#)

Adding Comment

All changes saved

[Save now](#)

New pull request

Dev ▾ into main ▾ ↗

Overview Files 1 Commits 1

1 changed file

The screenshot shows a GitHub pull request interface. At the top, it says "New pull request". Below that, the merge branch is set to "Dev" and the target branch is "main". There are tabs for "Overview", "Files" (which is selected and highlighted with a blue border), and "Commits". It indicates there is 1 changed file. On the left, there's a sidebar for the repository "CookbookRepo" which contains a "notebooks" folder. Inside "notebooks", there's a file named "PY 6.1-Reading Writing to Delta Tabl...". The main area shows the content of this file. Line 2, which contains "+ # Adding new comment", is highlighted with a red box.

```
1 1 # Databricks notebook source
2 2 + # Adding new comment
3 3 storageAccount="cookbookadlsgen2storage"
4 4 mountpoint = "/mnt/Gen2Source"
5 5 storageEndPoint ="abfss://rawdata@{}.dfs.
```

New personal access token

Personal access tokens function like ordinary OAuth access tokens. They can be HTTPS, or can be used to [authenticate to the API](#) over Basic Authentication.

Note

for ADB

What's this token for?

Select scopes

Scopes define the access for personal tokens. [Read more about OAuth scopes.](#)

<input checked="" type="checkbox"/> repo	Full control of private repositories
<input type="checkbox"/> repo:status	Access commit status
<input type="checkbox"/> repo_deployment	Access deployment status
<input type="checkbox"/> public_repo	Access public repositories
<input type="checkbox"/> repo:invite	Access repository invitations
<input type="checkbox"/> security_events	Read and write security events
<input type="checkbox"/> workflow	Update GitHub Action workflows

Generating tokens

To generate a GitHub personal access token, follow the [GitHub](#) permission.

To generate a Bitbucket Cloud app password, follow the [Bitbucket](#) must have "read" and "write" permission under repository.

To generate a GitLab personal access token, follow the [GitLab](#) "write_repository" permission.

To generate a Azure DevOps personal access token, follow the have "Full access" scope.

Git provider

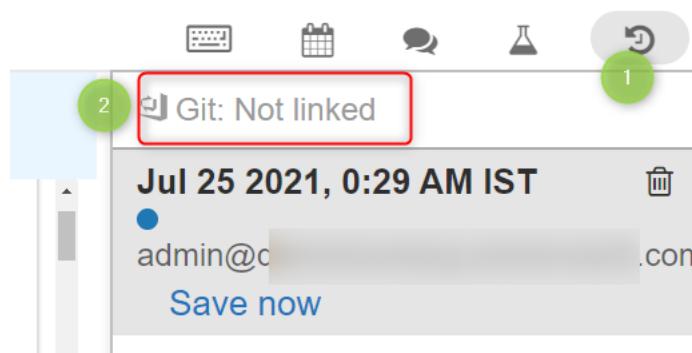
GitHub

Token or app password

.....

Git provider username or email [?](#)

Save



Status Link Unlink

Link <https://github.com/> /ADBRepo

Branch

Path in Git Repo [notebooks/6.1-Reading Writing to Delta Tables.py](#)

Tip: You can also import/export multiple notebooks or an entire folder through [Workspace API](#) to your computer and check-in to your favorite version control system.

The screenshot shows a GitHub repository named 'ADBRepo'. The 'Code' tab is selected. A dropdown menu shows 'main'. The commit history lists a single commit from 'rusum21' titled 'First Commit' for the file '6.1-Reading Writing to Delta Tables.py'. The commit message includes the text '#Adding comment for GitHub' and 'display(dbutils.fs.ls("/mnt/Gen2Source/Customer/parquetFiles"))'.

/ ADBRepo Private

<> Code Issues Pull requests Actions Projects Security

main ADBRepo / notebooks /

rusum21 First Commit

6.1-Reading Writing to Delta Tables.py First Commit

```
#Adding comment for GitHub
display(dbutils.fs.ls("/mnt/Gen2Source/Customer/parquetFiles"))
```

Git Preferences

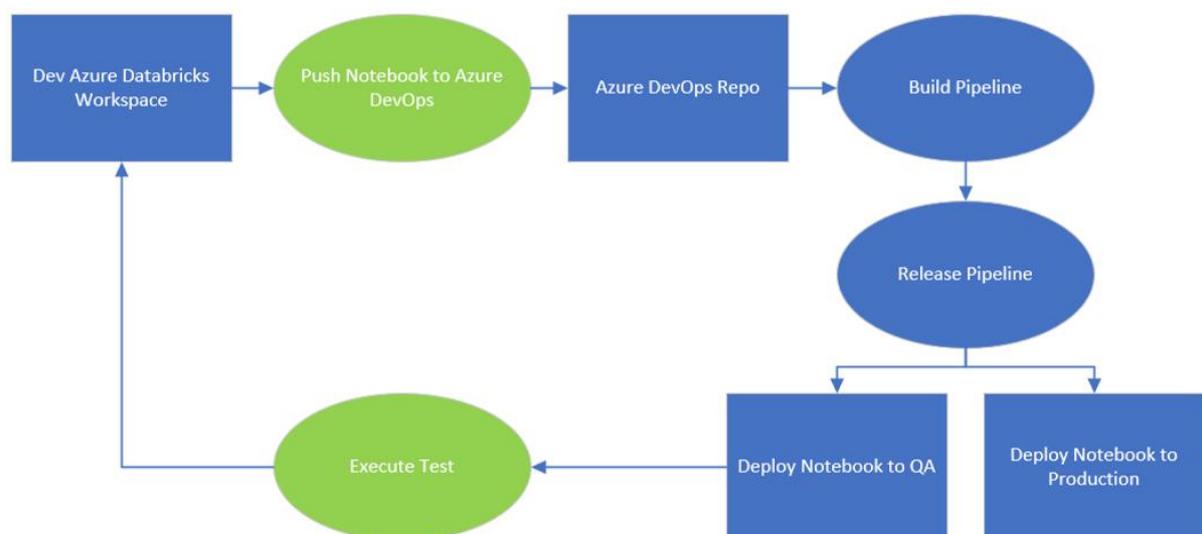
Status Link Unlink

Link <https://github.com/> /ADBRepo

Branch dev-sprint1

Path in Git Repo notebooks/6.1-Reading Writing to Delta Tables.py

Tip: You can also import/export multiple notebooks or an entire folder through [Workspace API](#) to your computer and check-in to your favorite version control system.



CookbookProj +

- Overview
- Boards
- Repos
- Pipelines
 - Pipelines 1
 - Environments
 - Releases
 - Library
 - Task groups
- Project settings <<

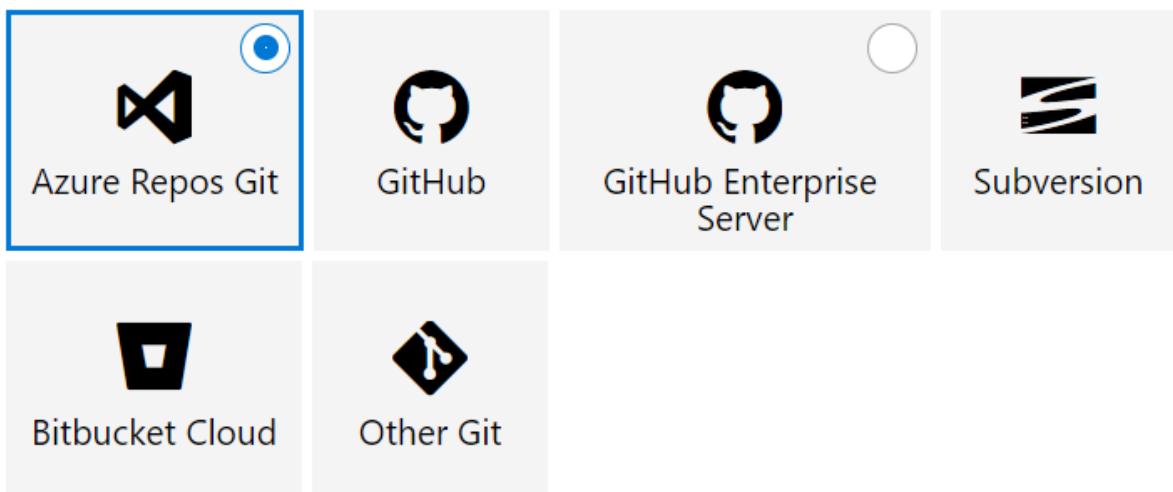
New pipeline

Where is your code?

- Azure Repos Git YAML
Free private Git repositories, pull requests, and code search
- Bitbucket Cloud YAML
Hosted by Atlassian
- Github YAML
Home to the world's largest community of developers
- Github Enterprise Server YAML
The self-hosted version of GitHub Enterprise
- Other Git
Any generic Git repository
- Subversion
Centralized version control by Apache

Use the classic editor 2 to create a pipeline without YAML.

Select a source



Team project

▼

Repository

▼

Default branch for manual and scheduled builds

▼

Select a template

Or start with an [Empty job](#)

Configuration as code



YAML

Looking for a better experience to configure your pipelines using YAML files? Try the new YAML pipeline creation experience. [Learn more](#)

Featured



.NET Desktop

Build and test a .NET or Windows classic desktop solution.



Android

The screenshot shows the Azure Pipelines interface for creating a pipeline named 'gent job 1'. The pipeline consists of a single step: 'Get sources' (CookbookRepo, main branch). A '+' button is highlighted with a red box (step 1). To the right, the 'Add tasks' menu is open, showing the 'publish build artifacts' task. This task is also highlighted with a red box (step 2). A callout box points to the 'Add' button at the bottom of the task list, which is also highlighted with a red box (step 3).

Get sources
CookbookRepo main

gent job 1
Run on agent

...

Add tasks

Refresh

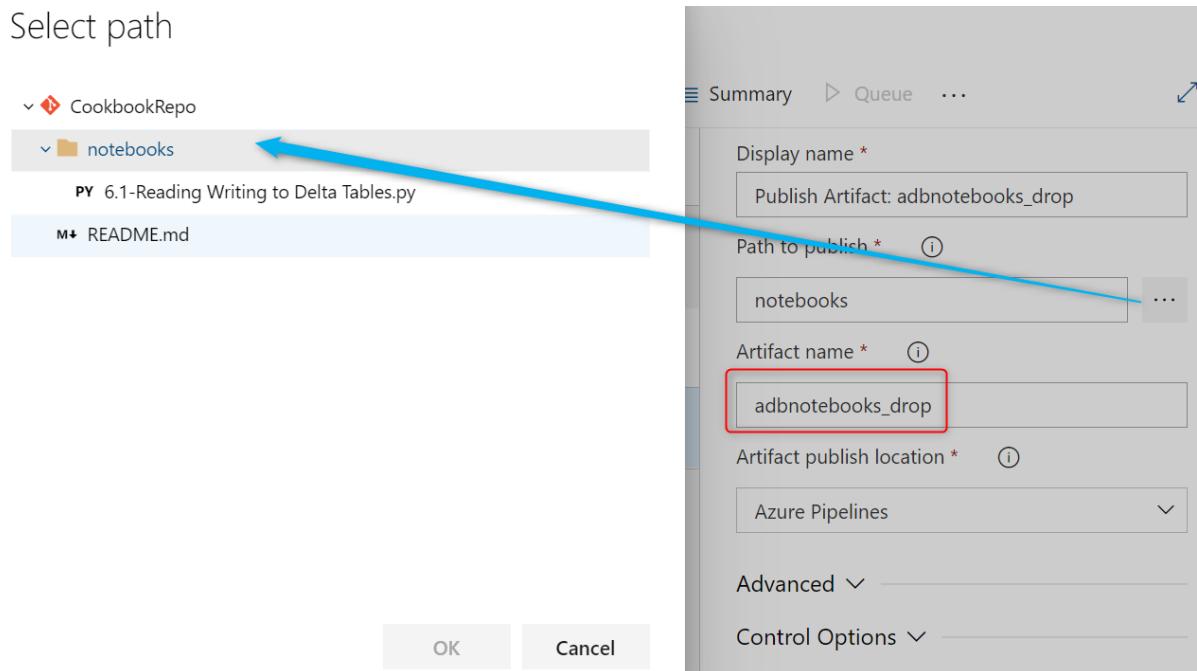
publish build artifacts

...

Publish build artifacts
Publish build artifacts to Azure Pipelines or a Windows file share

Add

Select path



Summary

Manually run by admin [View 8 changes](#)

Repository and version	Time started and elapsed	Related	Tests and coverage
❖ CookbookRepo ↳ main ↳ 018e95b	📅 Today at 6:58 PM ⌚ 16s	⌚ 0 work items ⌚ 1 published; 1 consumed	Get started

Jobs

Name	Status	Duration
Agent job 1	Success	⌚ 7s

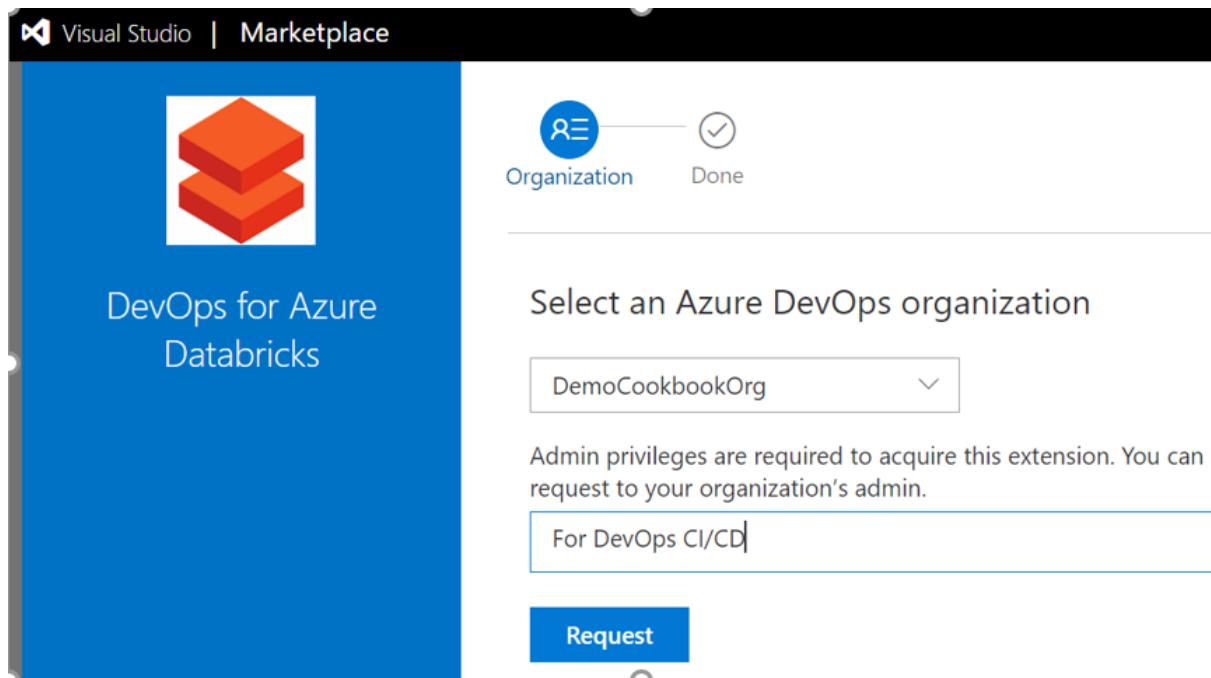
The screenshot shows the Azure DevOps Pipelines interface for the project "CookbookProj". The left sidebar has a red box around the "Releases" item, which is currently selected. The main area shows the "All pipelines > Pipeline" view. A red box highlights the "Empty job" template under the heading "Select a template". Below it, there's a "Featured" section with three cards:

- Azure App Service deployment**: Deploy your application to Azure App Service. Web App on Windows, Linux, containers, WebJobs.
- Deploy a Java app to Azure App Service**: Deploy a Java application to an Azure Web App.
- Deploy a Node.js app to Azure App**: Deploy a Node.js application to an Azure Web App.

The screenshot shows the "Pipeline" configuration screen for the "CookbookProj" project. The top navigation bar includes tabs for Pipeline, Tasks, Variables, Retention, Options, and History. The "Tasks" tab is selected. The interface is divided into two main sections: "Artifacts" and "Stages".

- Artifacts**: A section for adding artifacts, with a button "+ Add" and a placeholder "+ Add an artifact".
- Stages**: A section for adding stages, currently showing one stage named "Stage 1" with a red box around it. It indicates "1 job, 0 task".

Below the stages, there are tabs for Tasks, Variables, Retention, Options, and History. On the right side, there's a "Marketplace" section featuring the "DevOps for Azure Databricks" extension by Microsoft DevLabs, which has 3,251 installs. Buttons for "Get it free" and "Learn more" are available.



The screenshot shows the Azure DevOps pipeline editor for the 'Deploy Notebook' pipeline. The top navigation bar includes 'All pipelines > Deploy Notebook', 'Pipeline', 'Tasks', 'Variables', 'Retention', and 'Options'. The pipeline interface has two main sections: 'Artifacts' and 'Stages'. In the 'Artifacts' section (1), there's a button to 'Add an artifact' and a note that 'Schedule not set'. In the 'Stages' section (2), there's a 'Build' artifact type selected, indicated by a green circle with the number '2'. To the right, there are sections for 'Project', 'Source (build pipeline)', 'Default version', 'Source alias', and a note about artifact versions. A '5 more artifact types' link is also present. The 'Build' artifact details show 'CookbookProj' as the project, 'CookbookProj-CI' as the source pipeline, 'Latest' as the default version, and '_CookbookProj-CI' as the source alias. A callout note (3) explains that artifacts will be published from the latest successful build of the source pipeline. At the bottom right is a large blue 'Add' button (4).

All pipelines > ⚙ Deploy Notebooks

Save ...

Pipeline Tasks Variables Retention Options History

Stage 1	Deployment process	...
Agent job	Run on agent	+

Add tasks | Refresh

databricks



Configure Databricks
Configure Databricks CLI

Add

Pipeline Tasks Variables Retention Options History

Stage 1	Deployment process	...
Agent job	Run on agent	+
Configure Databricks CLI	Configure Databricks	✓

Task version 0.*

Display name *

Configure Databricks CLI

Workspace URL *

https://adb-.azuredatabricks.net

Access Token *

dab-3

All pipelines > ⚙ Deploy Notebooks

Save Create release ...

Pipeline Tasks Variables Retention Options History

Stage 1	Deployment process	...
Agent job	Run on agent	?
Configure Databricks CLI	Configure Databricks	+

Add tasks | Refresh

databricks



Deploy Databricks Notebooks

Recursively deploys Notebooks from given folder to a Databricks Workspace

3

Add

Select a file or folder

The screenshot shows a file browser interface with the following structure:

- Linked artifacts
- _CookbookProj-Cl (Build)
- adbnotebooks_drop** (highlighted with a red box)
- 6.1-Reading Writing to Delta Tables.py

A blue arrow points from the highlighted 'adbnotebooks_drop' item to the right pane, which displays the task configuration for 'Deploy Notebooks to Workspace'.

Task version: 0.*

Display name: Deploy Notebooks to Workspace

Notebooks folder: \$(System.DefaultWorkingDirectory)/_CookbookProj-Cl/adbnotebooks_drop

Workspace folder: /Shared

↑ Deploy Notebooks > Release-1 > Stage 1 ✓ Succeeded

← Pipeline Tasks Variables Logs Tests | Deploy Cancel Refresh ...

Deployment process Succeeded

Agent job Succeeded

Agent job

Started: Azure Pipelines · Agent: Hosted Agent

Initialize job · succeeded

Download artifact - _CookbookProj-Cl - ... · succeeded

Configure Databricks CLI · succeeded

Deploy Notebooks to Workspace · succeeded

Finalize Job · succeeded

All pipelines > Deploy Notebooks

Pipeline Variables Variables Retention Options History

Pipeline variables

Variable groups 3

Add variable groups to your pipeline. [Learn more about variable groups](#)

Link variable group Manage variable groups

Predefined variables

1 2 3 4

1: Releases (highlighted in green)

2: Variable groups (highlighted in green)

3: Predefined variables

4: Link variable group and Manage variable groups

Library > UAT Variable Group*

Variable group | Save | Clone | Security | Pipeline permissions | Approvals

Variable group name

UAT Variable Group

Description

Link secrets from an Azure key vault as variables ⓘ

Variables

Name ↑	Value
Token	*****
URL	https://adb-... .azuredatabrick
+ Add	

Library

Variable groups Secure files

Search variable groups

Name ↓

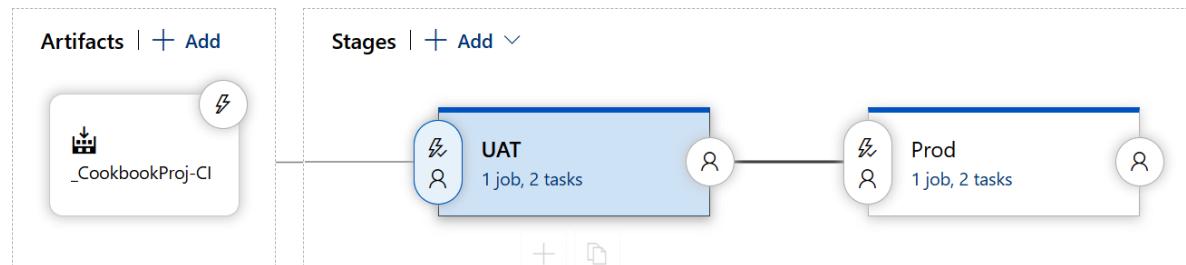
fx Prod Variable Group

fx UAT Variable Group

All pipelines > Deploy Notebooks

Save Create release ..

Pipeline Tasks Variables Retention Options History



All pipelines > Deploy Notebooks

Pipeline Tasks Variables Retention Option

Pipeline variables Variable groups Predefined variables

Prod Variable Group (2)

UAT Variable Group (2)

Variable group scope

Release

Stages

UAT

Link

This screenshot shows the 'Variables' tab in an Azure Pipeline named 'Deploy Notebooks'. The 'Variable groups' section is selected. On the right, a modal window titled 'Prod Variable Group (2)' displays the 'UAT Variable Group (2)' entry, which is highlighted with a red box. Below this, the 'Variable group scope' section shows 'Stages' selected, with 'UAT' also highlighted with a red box. A large blue 'Link' button is at the bottom of the modal.

Pipeline Tasks Variables Retention Options History

Pipeline variables Variable groups Predefined variables

Name	Value
UAT Variable Group (2)	Scopes: UAT
Token	*****
URL	https://adb-
Prod Variable Group (2)	... Scopes: Prod
Token	*****
URL	https://adb-

Link variable group | **Manage variable groups**

This screenshot shows the 'Variables' tab in the same pipeline. The 'Variable groups' section is selected. It lists two groups: 'UAT Variable Group (2)' and 'Prod Variable Group (2)'. Each group has a table row with 'Name' and 'Value' columns. The 'Value' column for 'UAT Variable Group (2)' contains 'Scopes: UAT' and the 'Value' column for 'Prod Variable Group (2)' contains '... Scopes: Prod'. At the bottom, there are buttons for 'Link variable group' and 'Manage variable groups'.

All pipelines > Deploy Notebooks

Save

Pipeline Tasks Variables Retention Options History

UAT

Deployment process

...

Display name *

Configure Databricks CLI

Agent job

Run on agent

+



Configure Databricks CLI

Configure Databricks



Deploy Notebooks to Work...

Deploy Databricks Notebooks

Workspace URL *



\$(URL)

Access Token *



\$(Token)

Control Options

Output Variables

Deploy Notebooks > Release-3

Pipeline

Variables

History

+ Deploy

Cancel

Refresh

Edit

...

Release

Manually triggered

by admin

7/24/2021, 11:23 PM

Artifacts



_CookbookProj-Cl

Stages

UAT

Succeeded

on 7/24/2021, 11:24 PM

Prod

Succeeded

on 7/24/2021, 11:25 PM

Library > Prod Variable Group*

Variable group

 Save

 Clone

 Security

 Pipeline permissions

Prod Variable Group

Description

Link secrets from an Azure key vault as variables 

Azure subscription * | [Manage](#) 



 Scoped to subscription 'Visual Studio Enterprise'

Key vault name * [Manage](#) 

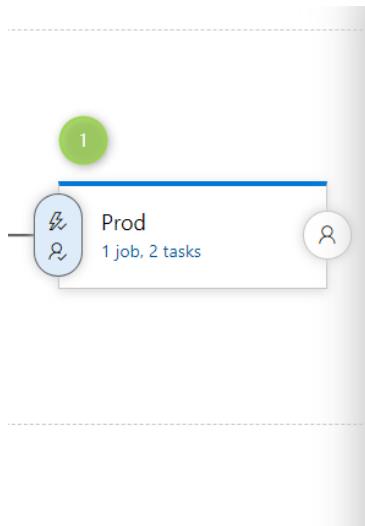
CookbookKeyVault



Variables

Last refreshed: 6 hours ago

Delete	Secret name	Content type	Status	Expiration date
	ADBTOKEN		Enabled	Never



Pre-deployment conditions

Prod

Triggers

Define the trigger that will start deployment to this stage

2

Enabled

 Pre-deployment approvals 

Select the users who can approve or reject deployments to this stage

Approvers 

 admin  Search users and groups for approvers

Timeout 

30

Days 

↑ Deploy Notebooks > Release-4

Pipeline Variables History + Deploy Cancel Approve multiple Refresh Edit ...

Prod
Pending approval
On admin for 12 minutes

Pre-deployment approval pending
On admin since 7/24/2021, 11:39 PM
Approve

Automatic trigger
Deployment triggered on 7/24/2021, 11:39 PM

C CookbookProj +

Boards
Repos
Pipelines
Pipelines 1
Environments
Releases
Library
Task groups

... > CookbookProj-CI

Tasks Variables Triggers Options History Save & queue Discard Summary ...

Continuous integration
CookbookRepo Enabled

Scheduled
No builds scheduled

Build completion
Build when another build completes

CookbookRepo
Enable continuous integration (checked)
Batch changes while a build is in progress (unchecked)

Branch filters
Type: Include Branch specification: main

All pipelines > Deploy Notebooks

Pipeline Tasks Variables Retention Options History

Artifacts + Add Continuous deployment trigger
1 _CookbookProj-CI

Stages + Add
2 Enabled
Creates a release every time a new build is available.

Build branch filters
3 Type: Include Build branch: main

Continuous deployment trigger
Build: _CookbookProj-CI

Pull request trigger



#6 Merged PR 7: dd

on CookbookProj-Cl

i This run has been retained forever by main (Branch).

Summary Releases

Manually run by A admin

Repository and version

❖ CookbookRepo

↳ main ↳ 018e95b

Time started and elapsed

📅 Today at 12:02 AM

⌚ 14s

Related

∅ 0 work items

☒ 1 published; 1 consumed

❖ CookbookRepo

↳ ARMTemplates

↳ AzureDatabricks

📄 azuredeploy.json

📄 azuredeploy.parameters.json

➢ notebooks

M README.md

↳ main ↴

Files

Contents

Name ↑

📁 ARMTemp

📁 notebook

M README.

All pipelines > New release pipeline

Pipeline Tasks Variables Retention Options

Artifacts | + Add

Add an artifact

Schedule not set

Add an artifact

Source type

Build ✓ Azure Re...

GitHub

TFVC

5 more artifact types ▾

Project * CookbookProj

Source (repository) * CookbookRepo

Default branch * main

Default version *

Override template parameters

Name	Value
disablePublicIp	false
workspaceName	"DevBigDataWS"
pricingTier	"premium"
location	"east us"

OK Cancel

Archie... Save Create release ...

Template location * Linked artifact

Template * \$(System.DefaultWorkingDirectory)/_CookbookRepo/ARMTemplates/AzureDatabricks/azuredeploy.json

Template parameters * \$(System.DefaultWorkingDirectory)/_CookbookRepo/ARMTemplates/AzureDatabricks/azuredeploy.parameters.json

Override template parameters - disablePublicIp false - workspaceName "DevBigDataWS" - pricingTier "premium" - location "east us"

Deployment mode * Incremental

Advanced ^

↑ DeployAzureDatabricksService > Release-2 > Stage 1 ✓ Succeeded

← Pipeline Tasks Variables Logs Tests | Deploy Cancel Refresh Download all logs

Deployment process Succeeded

Agent job Succeeded

Agent job

Pool: Azure Pipelines · Agent: Hosted Agent

- Initialize job · succeeded
- Download Artifacts · succeeded
- ARM Template deployment: Deploy ADB · succeeded
- Finalize Job · succeeded

Environment	Resource Group	Workspace Name
Dev	CookbookRG	BigDataAnalytics
UAT	UATCookbookRG	UATBigDataAnalytics
Production	ProdCookbookRG	ProdBigDataAnalytics

Chapter 10: Understanding Security and Monitoring in Azure Databricks

Home > cookbookadlsgen2storage1

cookbookadlsgen2storage1 | Access Control

Storage account

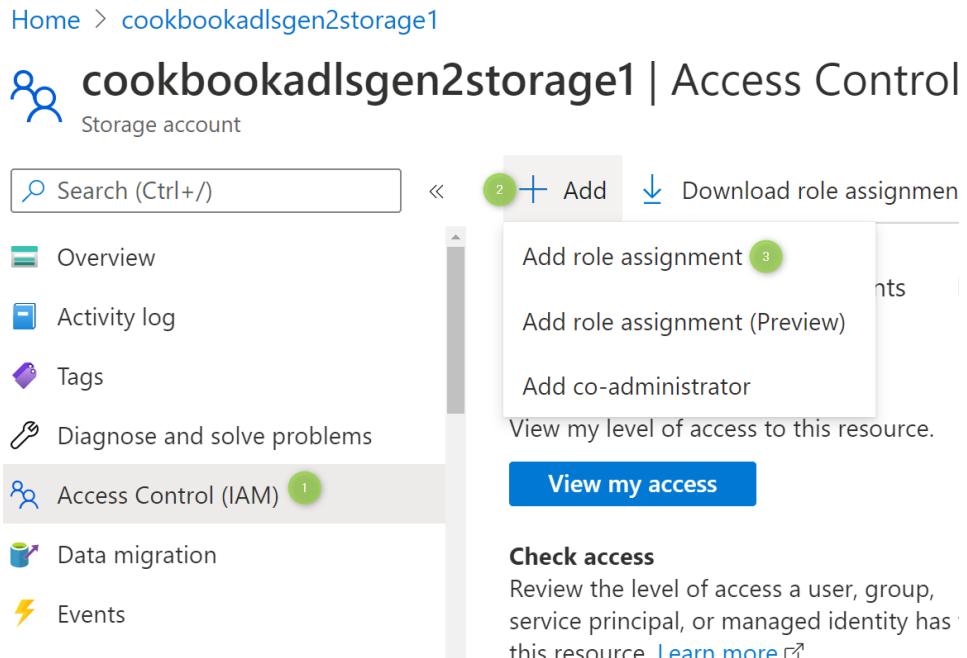
Search (Ctrl+ /) << Add Download role assignments

- Overview
- Activity log
- Tags
- Diagnose and solve problems
- Access Control (IAM) 1
- Data migration
- Events

Add role assignment 3
Add role assignment (Preview)
Add co-administrator
View my level of access to this resource.

View my access

Check access
Review the level of access a user, group, service principal, or managed identity has to this resource. [Learn more](#)



e1 | Access Control (IA)

Download role assignments

Access Role assignments Roles

level of access to this resource.

my access

the level of access a user, group, principal, or managed identity has to source. [Learn more](#)

group, or service principal

by name or email address

Add role assignment

Role [i](#)

Select a role

Select a role

Storage Account Backup Contributor Role [i](#)

Storage Account Contributor [i](#)

Storage Account Key Operator Service Role [i](#)

Storage Blob Data Contributor [i](#)

Storage Blob Data Owner [i](#)

Storage Blob Data Reader [i](#)

Storage Blob Delegator [i](#)

Storage File Data SMB Share Contributor [i](#)

Storage File Data SMB Share Elevated Contributor [i](#)

Add role assignment

X

Role ⓘ

Storage Blob Data Reader ⓘ



Assign access to ⓘ

User, group, or service principal



Select ⓘ

sqluser

No users, groups, or service principals found.

Selected members:



SQLUser

sqluser@demomoveorg.onmicrosoft.c...

[Remove](#)

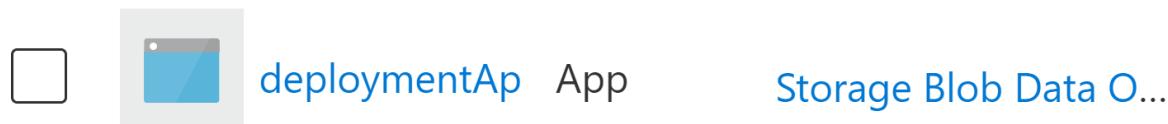
[Save](#)

[Discard](#)

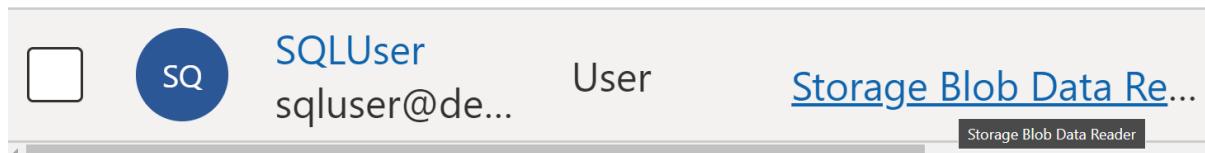
Storage Blob Data Reader				
<input type="checkbox"/>	 AU	Appuser appuser@demomoveo...	User	Storage Blob Data Re... This resource
<input type="checkbox"/>	 SQ	SQLUser sqluser@demomoveo...	User	Storage Blob Data Re... This resource

Key Vault Administrator				
<input type="checkbox"/>	 AU	admin user1 admin@demomoveo...	User	Key Vault Administrat... Subscription (Inherited)
Storage Blob Data Contributor				
<input type="checkbox"/>	 App			Storage Blob Data Co... This resource
Storage Blob Data Owner				
<input type="checkbox"/>	 App			Storage Blob Data O... This resource
Storage Blob Data Reader				
<input type="checkbox"/>	 SQ	SQLUser sqluser@demomoveo...	User	Storage Blob Data Re... This resource

Storage Blob Data Owner



Storage Blob Data Reader



Home > cookbookadlsgen2storage1 > sourcedata

sourcedata | Access Control (IAM) ...

Container

A screenshot of the Azure portal showing the "Access Control (IAM)" settings for the "sourcedata" container. The left sidebar shows options like "Search (Ctrl+/", "Overview", "Diagnose and solve problems", "Access Control (IAM)", "Settings", "Shared access tokens", "Manage ACL", "Access policy", "Properties", and "Metadata". The main area shows a table of role assignments:

	User	Role	Inheritance
<input type="checkbox"/>	admin@dem...	User	Key Vault Administrat... Subscription (Inherite...
<input type="checkbox"/>	App	Storage Blob Data Contributor	Storage Blob Data Co... Parent resource (Inher...
<input type="checkbox"/>	App	Storage Blob Data Reader	Storage Blob Data Re... This resource
<input type="checkbox"/>	SQLUser	Storage Blob Data Reader	Storage Blob Data Re... Parent resource (Inher...

A screenshot of the Azure portal showing the "Manage Access" dialog for the "sourcedata" container. The left sidebar shows the "EXPLORER" view with "cookbookadlsgen2storage1 (ADLS Gen2)" selected. The main area shows a table of entities with one row highlighted:

Name	Access Tier	Access Tier Last Modified	Last Modified	Blob Type	Content Type
Customer					

The "Add Entity" section contains a search bar with "sqluser" and a "Search" button. The result list shows "SQLUser" with the email "sqluser@demomoveorg.onmicrosoft.com".

Manage Access

Managing permissions for: sourcedata/Customer.

[Learn more about access control lists \(ACLs\).](#)

Users, groups, and service principals:

\$superuser (Owner)



\$superuser (Owning Group)



Other

Mask

SQLUser

sqluser@demomoveorg.onmicrosoft.com



Add

Permissions for: SQLUser

Read Write Execute

Access

Default *

Read and write permissions will only work for an entity if the entity also has execute permissions on all parent directories, including the container (root directory).

* The default ACL determines permissions for new children of this directory. Changing the default ACL does not affect children that already exist. [Learn more about default ACLs.](#)

OK

Cancel

The screenshot shows the 'Manage Access' dialog box from the Azure Storage Explorer. At the top, there are standard file operations: Upload, Download, Open, New Folder, Select All, Copy, Paste, Rename, Move, and Manage ACLs. Below the title 'Manage Access' is a message: 'Managing permissions for: sourcedata/Customer/parquetFiles.' A link 'Learn more about access control lists (ACLs)' is provided. The main area is titled 'Users, groups, and service principals:' and lists four items: '\$superuser (Owner)', '\$superuser (Owning Group)', 'Other', and 'Mask'. Each item has a pencil icon for editing.

The screenshot shows the Azure Storage Explorer interface with a folder named 'Customer' selected. The context menu for this folder is open, displaying the following options: Open, Download, Copy, Paste, Change Access Tier..., Get Shared Access Signature..., Manage Access Control Lists..., Propagate Access Control Lists..., Acquire Lease, and Break Lease... . The 'Propagate Access Control Lists...' option is highlighted with a red box.

← → ↘ ↙ sourcedata > Customer

Name

Manage Access

Manage Access

Managing permissions for: sourcedata/Customer/parquetFiles.

Learn more about access control lists (ACLs).

Users, groups, and service principals:

\$superuser (Owner)

\$superuser (Owning Group)

Other

Mask

SQLUser
sqluser@demomoveorg.onmicrosoft.com

Add

Activities

Clear completed

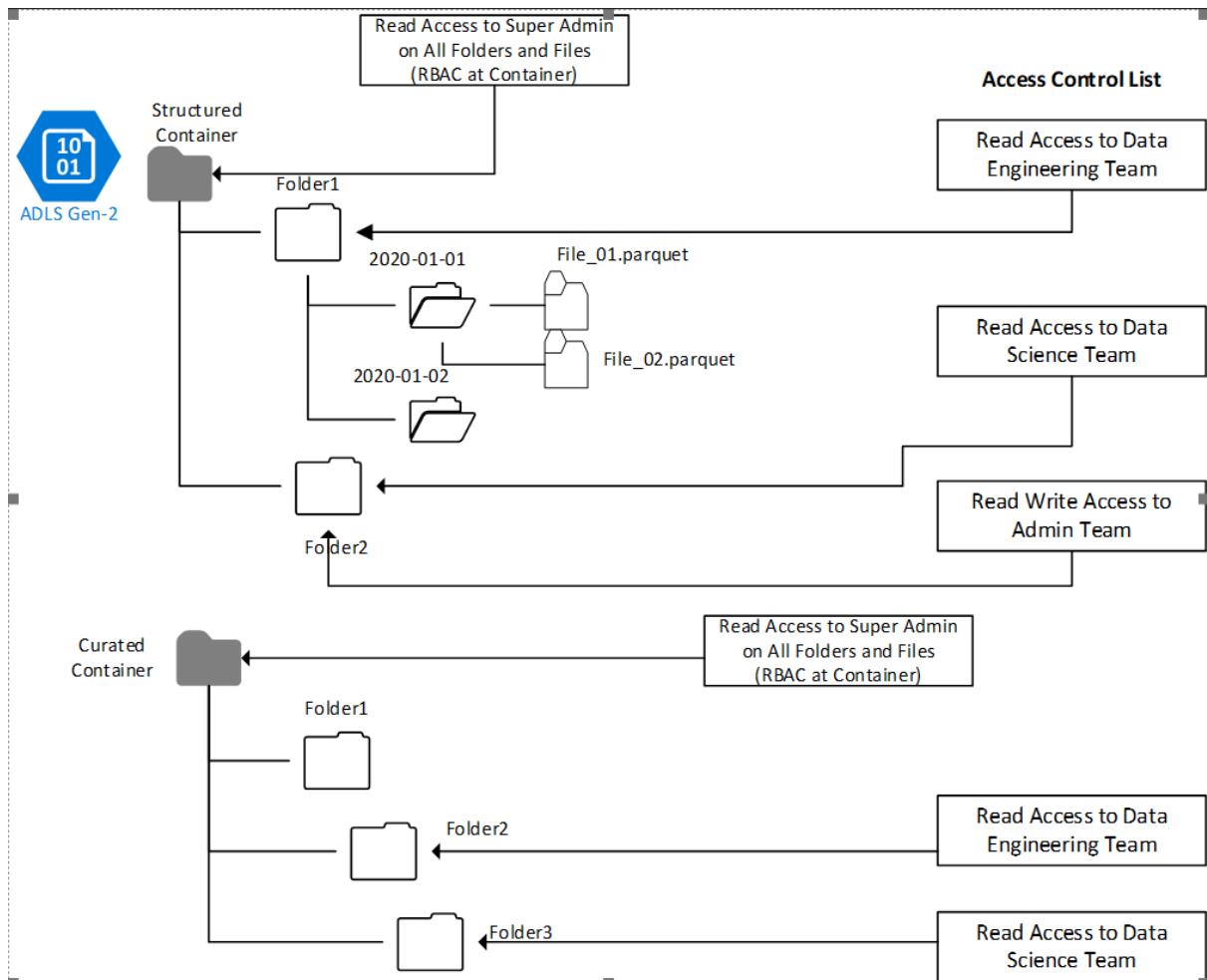
Showing 1 to 2 of 2 cached results

Permissions for: \$superuser

Access

Read Write Execute

The screenshot shows the 'Manage Access' blade for a specific folder path: 'sourcedata > Customer'. It displays the current permissions for the folder 'parquetFiles'. There are two entries in the 'Users, groups, and service principals:' list: '\$superuser (Owner)' and '\$superuser (Owning Group)'. Below this is a 'Mask' section containing 'SQLUser' with the email 'sqluser@demomoveorg.onmicrosoft.com'. An 'Add' button is available to grant permissions to other users or groups. At the bottom, it shows the 'Permissions for: \$superuser' section, where the 'Access' checkbox is checked, and the 'Read', 'Write', and 'Execute' checkboxes are all selected (indicated by checked boxes). The 'Activities' and 'Clear completed' buttons are visible at the bottom left, and the status 'Showing 1 to 2 of 2 cached results' is displayed at the bottom right.



New Cluster

Cancel

Create Cluster

1 Workers: 14 GB Memory, 4 Cores

1 Driver: 14 GB Memory, 4 Cores

Standard

Databricks Runtime Version [?](#)

[Learn more](#)

Runtime: 8.0 (Scala 2.12, Spark 3.1.1)

Note Databricks Runtime 8.x uses Delta Lake as the default table format. [Learn more](#)

Autopilot Options

Enable autoscaling [?](#)

Terminate after minutes of inactivity [?](#)

Worker Type [?](#)

Workers

Standard_DS3_v2

14 GB Memory, 4 Cores, 0.75 DBU

1

[?](#)

New Configure separate pools for workers and drivers for flexibility. [Learn more](#)

Driver Type

Same as worker

14 GB Memory, 4 Cores, 0.75 DBU

▼ Advanced Options

Azure Data Lake Storage Credential Passthrough [?](#)

Enable credential passthrough for user-level data access

Single User Access [?](#)

com

Only one user is allowed!

New Cluster

Cancel

Create Cluster

1 Workers:56 GB Memory, 8 Cores, 2 DBU

1 Driver:56 GB Memory, 8 Cores, 2 DBU [?](#)

Cluster Mode [?](#)

High Concurrency

Databricks Runtime Version [?](#)

[Learn more](#)

Runtime: 8.0 (Scala 2.12, Spark 3.1.1)

Note Databricks Runtime 8.x uses Delta Lake as the default table format. [Learn more](#)

Autopilot Options

Enable autoscaling [?](#)

Terminate after minutes of inactivity [?](#)

Worker Type [?](#)

Workers

Standard_DS13_v2

56 GB Memory, 8 Cores, 2 DBU

[?](#) Spot instances [?](#)

New Configure separate pools for workers and drivers for flexibility. [Learn more](#)

Driver Type

Same as worker

56 GB Memory, 8 Cores, 2 DBU

▼ Advanced Options

Azure Data Lake Storage Credential Passthrough [?](#)

Enable credential passthrough for user-level data access and only allow Python and SQL commands

[+ Create](#) [Edit columns](#) [Delete resource group](#) [Refresh](#) [Export to CSV](#) [Open qu](#)

^ Essentials

Subscription ([change](#))

Deployments
1 Failed, 4 Succeeded

Subscription ID

Location
East US

Tags ([change](#))

[Click here to add tags](#)

Filter for any field...

Type == all [X](#)

Location == all [X](#)

[+ ↴ Add filter](#)

Showing 1 to 5 of 5 records. Show hidden types [\(i\)](#)

No grouping

Name ↑↓

Type ↑↓

 cookbookadlsgen2storage1

Storage account

 PassthroughWorkspace

Azure Databricks Service



cookbookadlsgen2storage1 | Containers

[↗](#) [...](#)

Storage account

Search (Ctrl+ /)

[+ Container](#) [!\[\]\(94c1d1a0891a053538583e9f6cb61576_img.jpg\) Change access level](#) [!\[\]\(8d7b4178fb3a8d943c307ac739637331_img.jpg\) Restore](#)

 Containers

Search containers by prefix

 File shares

Name

rawdata

 Queues

sourcedata

 Tables

Security + networking

storage

Appuser assignments - cookbookadlsgen2storage1

Assignments for the selected user, group, service principal, or managed identity at this scope or inherited to

+ Add

Search by assignment name or description

Role assignments (1) ⓘ

Role	Description	Scope	Group assignments
Reader	View all resources, but does not al...	This resource	--

Deny assignments (0) ⓘ

Classic administrators (0) ⓘ

My access

View my le...

Check acc...

View m...

Check acc...

Review the...

Admin Console / Groups / Edit group

users

Members Entitlements Parent Groups



User or Group Name ▲

Type ▼

appuser@demomoveorg.onmicrosoft.com

User

Permission Settings for: HighConcurrencyCluster

x

NAME

PERMISSION

Can Manage

admins

Can Manage

inherited

appuser@demomoveorg.onmicrosoft.com

Can Attach To

+ Add

Cancel

Save

sourcedata | Access Control (IAM) ...

Container

Search (Ctrl+ /) « Add Download role assignments Edit columns Refresh Remove

Overview Diagnose and solve problems Access Control (IAM)

Storage Blob Data Contributor

<input type="checkbox"/>		App	Storage Blob Data Co...	Parent resource
--------------------------	--	-----	-------------------------	-----------------

Storage Blob Data Owner

<input type="checkbox"/>		App	Storage Blob Data Ow...	Parent resource
--------------------------	--	-----	-------------------------	-----------------

Storage Blob Data Reader

<input type="checkbox"/>		Appuser appuser@de...	User	Storage Blob Data Re...	This resource
--------------------------	--	--------------------------	------	-------------------------	---------------

```

1 #Reading the files from ADLS Gen-2 using AAD authentication of the user.
2 display(dbutils.fs.ls("abfss://sourcedata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles"))

```

▶ (3) Spark Jobs

	path	name
1	abfss://sourcedata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00000-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-56-1-c000.snappy.parquet	part-00000-tid-78773
2	abfss://sourcedata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00001-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-57-1-c000.snappy.parquet	part-00001-tid-78773
3	abfss://sourcedata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00002-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-58-1-c000.snappy.parquet	part-00002-tid-78773
4	abfss://sourcedata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00003-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-59-1-c000.snappy.parquet	part-00003-tid-78773

Operation failed: "This request is not authorized to perform this operation using this permission.", 403, GET, https://cookbookadlsgen2storage1.dfs.core.windows.net/rawdata?upn=false&resource=filesystem&maxResults=500&timeOut=90&recursive=false, AuthorizationPermissionMismatch, "This request is not authorized to perform this operation using this permission. RequestId:3bc110f1-001f-008c-32a1-6c2309000000 Time:2021-06-29T04:45:08.911058Z"

Command took 0.68 seconds -- by appuser@demomoveorg.onmicrosoft.com at 6/29/2021, user

Caused by: org.apache.spark.SparkException: Job aborted due to stage failure: Task 0 in stage 15.0 failed 4 times, most recent failure Lost task 0.3 in stage 15.0 (TID 26) (10.139.64.5 executor 0): Operation failed: "This request is not authorized to perform this operation using this permission.", 403, PUT, https://cookbookadlsgen2storage1.dfs.core.windows.net/sourcedata/Customer/delta/part-00000-54a93d8-1f28-48b2-8a33-a5b9282b3863-c000.snappy.parquet?resource=file&timeout=90, AuthorizationPermissionMismatch, "This request is not authorized to perform this operation using this permission. RequestId:23f003f7-601f-0033-462b-8914ac000000 Time:2021-08-04T12:22:50.657494Z"

Key Vault Administrator					
<input type="checkbox"/>	AU	admin user1 mov...	User	Key Vault Administrator	Subscription (Inherited)
Reader					
<input type="checkbox"/>	AP	Appuser appuser@demomo...	User	Reader	Resource group (Inherited)
Storage Blob Data Contributor					
<input type="checkbox"/>		App		Storage Blob Data Contribut...	This resource
Storage Blob Data Owner					
<input type="checkbox"/>		App		Storage Blob Data Owner	This resource
Storage Blob Data Reader					
<input type="checkbox"/>	SQ	SQLUser sqluser@demomov...	User	Storage Blob Data Reader	This resource

Users, groups, and service principals:

<input type="checkbox"/> rawdata	\$superuser (Owner)	<input type="button" value="Edit"/>
<input type="checkbox"/> Name	\$superuser (Owning Group)	<input type="button" value="Edit"/>
<input type="checkbox"/> parquetFiles	Other	
<input type="checkbox"/> parquetFiles_Daily	Mask	
	Appuser appuser@demomoveorg.onmicrosoft.com	<input type="button" value="Delete"/>

Add

Permissions for: Appuser

<input checked="" type="checkbox"/> Access	<input type="checkbox"/> Read	<input type="checkbox"/> Write	<input checked="" type="checkbox"/> Execute
<input type="checkbox"/> Default *	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Manage Access

Managing permissions for: rawdata/Customer.
Learn more about access control lists (ACLs).

Users, groups, and service principals:

\$superuser (Owner)	
\$superuser (Owning Group)	
Other	
Mask	
Appuser appuser@demomoveorg.onmicrosoft.com	

Add

Permissions for: Appuser

Read	Write	Execute
<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

cookbookadlsgen2storage1 (ADLS Gen2)

- blob Containers
 - rawdata
 - sourcedata
- File Shares
- Queues
- Tables

Users, groups, and service principals:

\$superuser (Owner)	
\$superuser (Owning Group)	
Other	
Mask	
Appuser appuser@demomoveorg.onmicrosoft.com	

Add

Permissions for: Appuser

Read	Write	Execute
<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

```
1 #Reading the files from ADLS Gen-2 using AAD authentication of the user with ACLs granted to folders
2 display(dbutils.fs.ls("abfss://rawdata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles"))
```

▶ (2) Spark Jobs

	path	name
1	abfss://rawdata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00000-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-56-1-c000.snappy.parquet	part-00000-tid-7877307C
2	abfss://rawdata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00001-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-57-1-c000.snappy.parquet	part-00001-tid-7877307C
3	abfss://rawdata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00002-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-58-1-c000.snappy.parquet	part-00002-tid-7877307C
4	abfss://rawdata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00003-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-59-1-c000.snappy.parquet	part-00003-tid-7877307C
5	abfss://rawdata@cookbookadlsgen2storage1.dfs.core.windows.net/Customer/parquetFiles/part-00004-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-60-1-c000.snappy.parquet	part-00004-tid-7877307C

operation failed: "This request is not authorized to perform this operation using this permission.", 403, GET, https://cookbookadlsgen2storage1.dfs.core.windows.net/rawdata?upn=false&resource=filesystem&maxResults=500&directory=Customer/parquetFiles_Daily&timeout=90&recursive=false, AuthorizationPermissionMismatch, "This request is not authorized to perform this operation using this permission. RequestId:6afba26f-801f-0066-29a8-6c0427000000 Time:2021-06-29T05:35:43.8238029Z"

Users, groups, and service principals:

\$superuser (Owner)	
\$superuser (Owning Group)	
Other	
Mask	
Appuser appuser@demomoveorg.onmicrosoft.com	

Add

Permissions for: Appuser

Read	Write	Execute
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Caused by: Operation failed: "This request is not authorized to perform this operation using this permission.", 403, GET, https://cookbookadlsgen2storage1.dfs.core.windows.net/rawdata/Customer/parquetFiles/part-00000-tid-7877307066025976411-3a62bd06-bc9c-467f-ae7b-d5db8ca40833-56-1-c000.snappy.parquet?timeout=90, AuthorizationPermissionMismatch, "This request is not authorized to perform this operation using this permission. RequestId:d5024788-701f-0096-55aa-6c42d6000000 Time:2021-06-29T05:51:03.3595431Z"

← → ↘ ↙ rawdata > Customer

Name	Access Tier	Access Tier Last Modi
parquetFiles		
parquetF		

parquetFiles

parquetF

-
-
-
-
-
-
-
-
-

+ Subnet + Gateway subnet Refresh | Manage users

Search subnets

Name ↑↓	IPv4 ↑↓	IPv6 ↑↓
default	10.0.0.0/24	-
private-subnet	10.0.2.0/24	-
public-subnet	10.0.1.0/24	-

Create an Azure Databricks workspace

...

Manage all your resources.

Subscription * ⓘ



Resource group * ⓘ

CookbookStorageRG



[Create new](#)

Instance Details

Workspace name *

vnet-adbworkspace



Region *

East US



Pricing Tier * ⓘ

Standard (Apache Spark, Secure with Azure AD) ↗

Deploy Azure Databricks workspace with Secure Cluster Connectivity (No Public IP) ⓘ

Yes No

Deploy Azure Databricks workspace in your own Virtual Network (VNet)

Yes No

Virtual Network * ⓘ

adb-vnet1



Two new subnets will be created in your Virtual Network

Implicit delegation of both subnets will be done to Azure Databricks on your behalf

Public Subnet Name *

public-subnet

Public Subnet CIDR Range * ⓘ

10.0.1.0/24



Private Subnet Name *

private-subnet

Private Subnet CIDR Range * ⓘ

10.0.2.0/24



[Review + create](#)

[< Previous](#)

[Next : Advanced >](#)

vnet-adbworkspace Azure Databricks Service

Search (Ctrl+ /) Delete

Overview

- Activity log
- Access control (IAM)
- Tags

Settings

- Virtual Network Peerings
- Encryption
- Properties

Essentials

Status	Managed Resource Group databricks-rg-vnet-adbworkspace-hn7cv
Active	URL https://adb-19.azure.databricks.net
Resource group	CookbookStorageRG
Location	Pricing Tier standard
Subscription	Virtual Network adb-vnet1
Subscription ID	Private Subnet Name private-subnet

cookbookadlsgen2storage1 | Networking

Storage account

Search (Ctrl+/)

Networking

- Containers
- File shares
- Queues
- Tables

Security + networking

- Networking (1)
- Access keys
- Shared access signature

Allow access from Selected networks

Configure network security for your storage accounts. [Learn more](#)

Virtual networks

Add existing virtual network Add new virtual network

Virtual Network	Subnet	Address range
No network selected.		

Firewall

Add IP ranges to allow access from the internet or your on-prem

Add your client IP address ('49.37.148.245')

Address range

Subscription *

Virtual networks * adb-vnet1

Subnets * public-subnet (Service endpoint required)

The following networks don't have service endpoints enabled for 'Microsoft.Storage'. Enabling access will take up to 15 minutes to complete. After starting this operation, it is safe to leave and return later if you do not wish to wait.

Virtual network Service endpoint status

Enable

All networks Selected networks

[Configure network security for your storage accounts. Learn more](#)

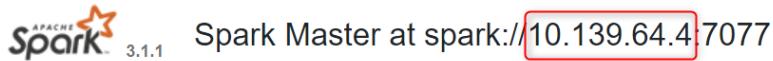
Virtual networks

[Add existing virtual network](#) [Add new virtual network](#)

Virtual Network	Subnet	Address range	Endpoint Status
adb-vnet1	1 public-subnet	10.0.1.0/24	✓ Enabled

⚡ HighConcurrencyCluster ⚡

Configuration Notebooks (2) Libraries Event Log Spark UI Driver Logs Metrics Apps **Spark Cluster UI - Master**



Spark Master at spark://10.139.64.4:7077

URL: spark://10.139.64.4:7077
Alive Workers: 2
Cores in use: 8 Total, 8 Used
Memory in use: 17.8 GiB Total, 14.2 GiB Used
Resources in use:
Applications: 1 Running, 0 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

▼Workers (2)

Worker Id	Address	State	Cores
worker-20210805010536-10.139.64.6-35581	10.139.64.6:35581	ALIVE	4 (4 Used)
worker-20210805010540-10.139.64.5-41393	10.139.64.5:41393	ALIVE	4 (4 Used)

⚡ HighConcurrencyCluster ⚡

Configuration Notebooks (2) Libraries Event Log **Spark UI** Driver Logs Metrics Apps

Live Metrics

Ganglia UI ↗

⚡ HighConcurrencyCluster ⚡

Configuration Notebooks (2) Libraries Event Log Spark UI Driver Logs Metrics

Live Metrics

Ganglia UI ↗

⚡ HighConcurrencyCluster

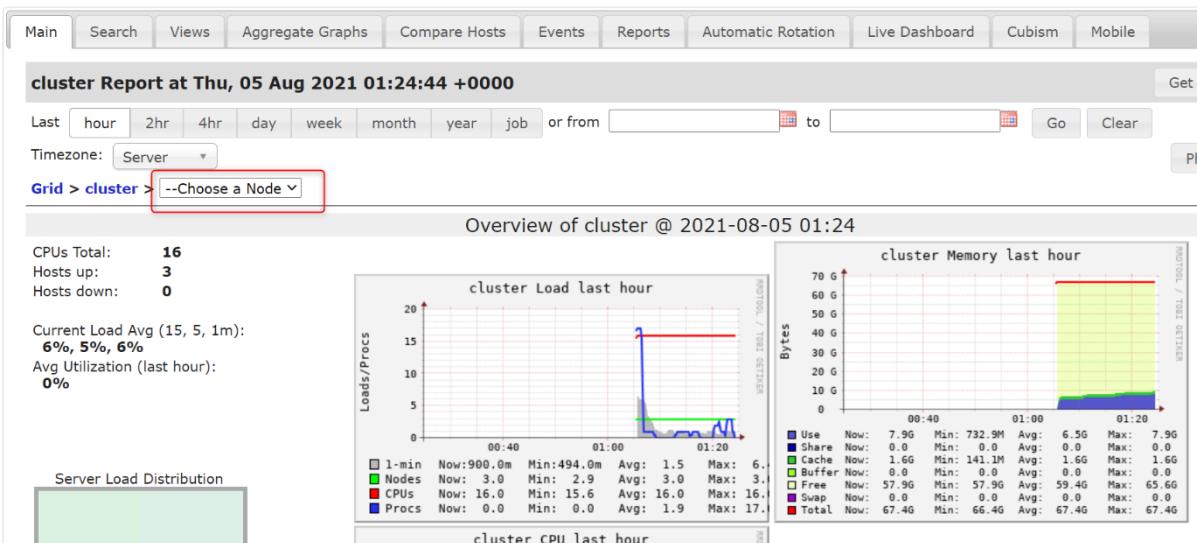
[Edit](#)[Permissions](#)[Clone](#)[Configuration](#)[Notebooks \(2\)](#)[Libraries](#)[Event Log](#)[Spark UI](#)[Driver Logs](#)[Metrics](#)[Ap](#)[Ganglia UI](#)

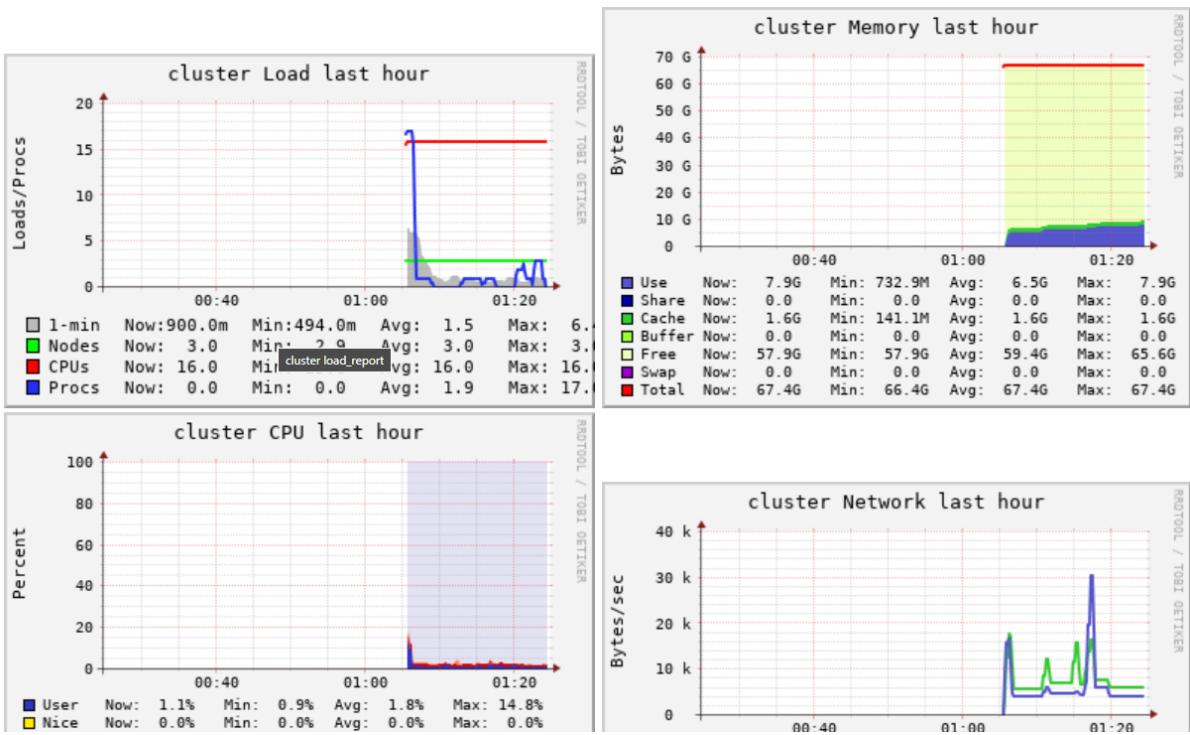
Historical Metrics Snapshots (5 files)

Name ↓

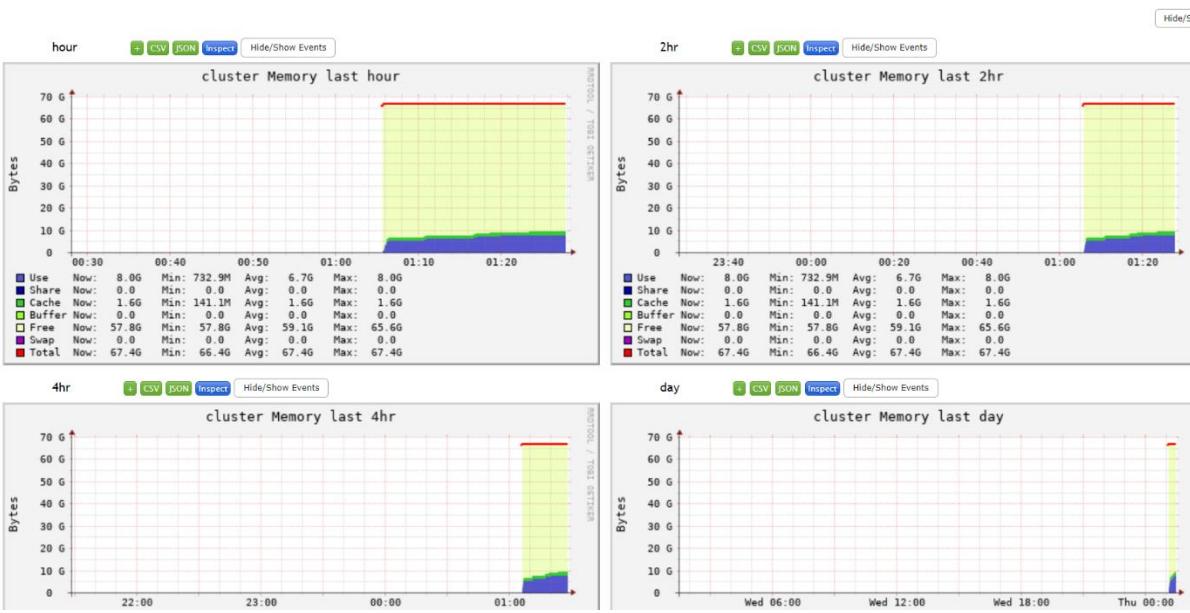
2021-08-04 18:30:01 IST

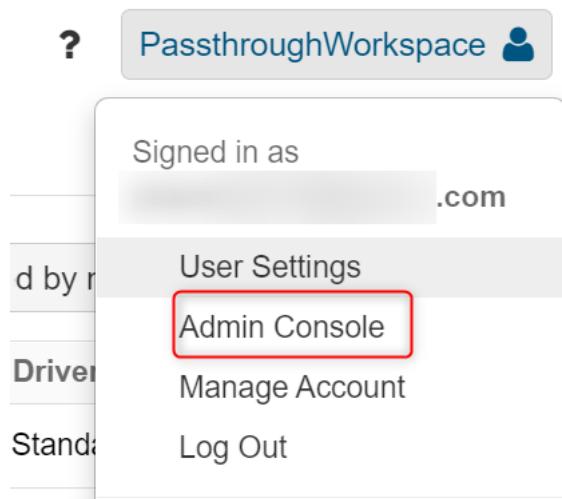
2021-08-04 18:15:01 IST





Cluster: cluster Graph: mem_report





Admin Console

Users Groups Global Init Scripts **Workspace Settings**

Filter

Access Control

Access Control

> Workspace Access Control: **Enabled**



> Cluster Visibility Control: **Enabled**



Who has access:

admins

Can Manage ▾



.com)

Can Manage ▾

No Permissions

Can Attach To

Can Restart

Can Manage

Add Users, Groups, and Service Principals:

Select User, Group or Service Principal... ▾