



Azure Data Factory

Azure Data Factory is a cloud-based data integration service provided by Microsoft Azure. It allows you to create, schedule, and manage data pipelines that move data between various supported data stores and processing systems. Here's an overview of its key features and components:

1. **Data Pipelines:** Azure Data Factory enables you to create data pipelines that orchestrate and automate the movement and transformation of data from various sources to destinations. These pipelines can include activities like data ingestion, data transformation, and data loading.
2. **Integration with Various Data Sources:** It supports integration with a wide range of data sources, including Azure Blob Storage, Azure SQL Database, Azure Data Lake Storage, Azure Synapse Analytics (formerly SQL Data Warehouse), on-premises data sources, and many others.
3. **Data Transformation:** Azure Data Factory allows you to transform data using various built-in activities such as mapping data flows, SQL transformations, and custom code activities. Data flows provide a visual interface for building data transformation logic using a drag-and-drop approach.
4. **Monitoring and Management:** It provides monitoring and management capabilities through Azure Monitor and Azure Data Factory's monitoring dashboard. You can monitor pipeline runs, activity runs, and trigger events, and set up alerts for monitoring pipeline health and performance.
5. **Integration with Other Azure Services:** Azure Data Factory integrates with other Azure services such as Azure Databricks, Azure Machine Learning, Azure Functions, and more, enabling you to build end-to-end data integration and analytics solutions.
6. **Scalability and Reliability:** It offers scalability and reliability features such as auto-scaling, fault tolerance, and retry policies to ensure that your data integration processes can handle varying workloads and maintain high availability.
7. **Security:** Azure Data Factory provides robust security features including encryption of data in transit and at rest, role-based access control (RBAC), integration with Azure Active Directory for authentication, and network security controls.



Use cases of Azure data factory:

Azure Data Factory is a versatile tool with a wide range of use cases across industries. Here are some common scenarios where organizations leverage Azure Data Factory:

1. **Data Warehousing:** Companies often use Azure Data Factory to populate and manage data warehouses. It can extract data from various sources, transform it to fit the schema of the data warehouse, and load it into Azure Synapse Analytics or other data warehousing solutions.
2. **ETL (Extract, Transform, Load) Processes:** Azure Data Factory is commonly used for ETL processes, where data needs to be extracted from multiple sources, transformed according to business rules, and loaded into a target destination. This could include

aggregating data from sales systems, CRM databases, and web logs, then loading it into a data lake or data warehouse.

3. **Real-time Data Ingestion:** Organizations dealing with real-time data streams, such as IoT sensor data, social media feeds, or clickstream data, can use Azure Data Factory to ingest, process, and analyze streaming data in near real-time. It can connect to Azure Event Hubs, Azure IoT Hub, or other streaming sources for continuous data ingestion.
4. **Data Migration:** Azure Data Factory simplifies the process of migrating data from on-premises data sources to the cloud or between different cloud platforms. It supports various data migration scenarios, including database migration, file migration, and big data migration.
5. **Data Integration for Analytics:** Businesses often use Azure Data Factory to integrate data from multiple sources for analytics purposes. It can combine structured and unstructured data from sources such as databases, files, and cloud services, and transform it into formats suitable for analysis using tools like Azure Synapse Analytics, Azure Databricks, or Power BI.
6. **Data Archiving and Backup:** Azure Data Factory can automate the process of archiving and backing up data to ensure data durability and compliance. It can schedule regular backups of on-premises or cloud-based data stores to Azure Blob Storage, Azure Data Lake Storage, or Azure SQL Database.
7. **Machine Learning Pipelines:** Organizations leveraging machine learning for predictive analytics or AI applications can use Azure Data Factory to build and manage end-to-end machine learning pipelines. It can orchestrate the flow of data between data storage, data preprocessing, model training, and model deployment stages.
8. **Hybrid Data Integration:** Azure Data Factory supports hybrid data integration scenarios where data resides both in the cloud and on-premises. It can securely connect to on-premises data sources using self-hosted integration runtimes, enabling seamless integration between cloud-based and on-premises data environments.

In this walkthrough, we're setting up Azure Data Factory to orchestrate the movement and transformation of data between various sources and destinations. The end goal is to establish a data integration pipeline that extracts data from a CSV file stored in Azure Data Lake Storage Gen2, transforms it, and loads it into a table in Azure Synapse Analytics. This process allows organizations to automate data workflows, ensuring data consistency, efficiency, and accessibility for analytics and decision-making purposes.

To begin with the Lab:

1. There are some prerequisites for this lab and they are, that you should have resources from previous labs.
2. You should have an external storage account created separately and it should have two containers one for CSV and the other one for Parquet, then the other storage account should be created while creating your Azure Synapse Workspace.
3. Then you should also have an SQL Database on which sample data should be loaded while creating it.

Name ↑↓	Type ↑↓	Location ↑↓	
<input type="checkbox"/> datalake121	Storage account	North Europe	...
<input type="checkbox"/> dataserver121	SQL server	North Europe	...
<input type="checkbox"/> datasynapse1234	Synapse workspace	North Europe	...
<input type="checkbox"/> demodb (dataserver121/demodb)	SQL database	North Europe	...
<input type="checkbox"/> demopooldb (datasynapse1234/demopooldb)	Dedicated SQL pool	North Europe	...
<input type="checkbox"/> sqlstorage1010	Storage account	North Europe	...

4. Now you are going to create Azure Data Factory. For that go to create resources section and search for it. Then choose this service accordingly.

The screenshot shows the Microsoft Azure portal interface. At the top, there's a navigation bar with a search bar and several other service icons. Below the bar, the main content area has a title 'Data Factory'. To the left is a blue icon representing the service. To the right of the icon, the word 'Data Factory' is displayed in large, bold, black font, followed by a small blue heart icon and the text 'Add to Favorites'. Below this, the text 'Microsoft | Azure Service' is shown. Underneath, there's a section titled 'Plan' with a dropdown menu containing the option 'Data Factory'. To the right of the dropdown is a large blue 'Create' button.

5. After that you have to choose your resource group where your all resources reside. Then give a unique name to Azure data factory. After that choose your location. Then click on next.

Basics Git configuration Networking Advanced Tags Review + create

One-click to create data factory with sample pipeline and datasets. [Try it](#)

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

Name * ⓘ ✓

Region * ⓘ ▼

Version * ⓘ ▼

6. On the next page click on configure git later on and move directly to review page and create your resources.

Basics **Git configuration** Networking Advanced Tags Review + create

Azure Data Factory allows you to configure a Git repository with either Azure DevOps or GitHub. Git is a version control system that allows for easier change tracking and collaboration.

[Learn more about Git integration in Azure Data Factory](#)

Configure Git later ⓘ

7. Now you have to wait for some time and wait until the deployment gets complete. After that click on go to resources.

✓ Your deployment is complete

 Deployment name : Microsoft.DataFactory-20240422193152 Start time : 4/22/2024, 7:35:16 PM
Subscription : Azure Pass - Sponsorship (e41df6f3-2d66-416f-9924...) Correlation ID : 622d230a-5fcf-4529-96e3-39a5dbf9c9c6
Resource group : demo-resource-group

› Deployment details

∨ Next steps

[Go to resource](#)

8. From its dashboard click on launch studio.

Home > Microsoft.DataFactory-20240422193152 | Overview >

demodatafactory121 Data factory (V2)

Search

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Settings

Networking

Managed identities

Properties

Locks

Getting started

Quick start

Monitoring

Alerts

Metrics

Diagnostic settings

Logs

Essentials

Resource group (move) : demo-resource-group

Status : Succeeded

Location : North Europe

Subscription (move) : Azure Pass - Sponsorship

Subscription ID : e41df6f3-2d66-416f-9924-552b6cda27ec

Type : Data factory

Copy to clipboard

Getting started : Quick start

JSON View

Azure Data Factory Studio

Launch studio

Quick Starts

Tutorials

Template Gallery

Training Modules

Monitoring

9. This will give us an entire interface in which we can start working with our pipelines, activities, our data sets, link services, etcetera.
10. Below you can see that it is similar to Azure Synapse Studio.

Microsoft Azure | Data Factory > demodatafactory121 | Search factory and documentation

Home

Author

Monitor

Manage

Learning Center

Collapse

demodatafactory121

New

Ingest

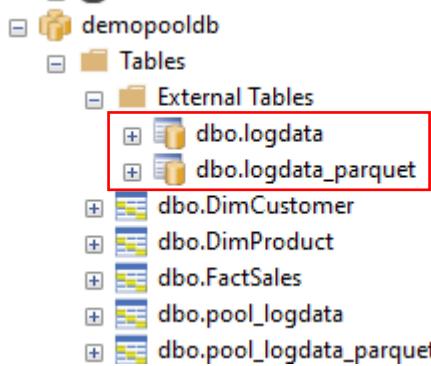
Orchestrate

Transform data

Configure SSIS

Recent resources

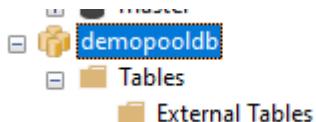
11. Now we need to get into SSMS (SQL Server Management Studio). Connect your SQL Database and Synapse Workspace.
12. Here you will expand your Pool db tables. Then you will find your external tables and normal tables which you have created while working with synapse.
13. Or if you have already deleted all your tables and starting again then you just have to create synapse workspace, SQL Database and pool db in your synapse workspace. Along with that you should have your storage account in place in which you should be having two containers one for csv and other for parquet with respective files in it.



14. Now we are going to Drop these tables. We did this so that we could have a clean slate for ourselves.

```
SQLQuery1.sql - dat...(sqladminuser (0))* ↴ ×  
DROP EXTERNAL TABLE logdata  
DROP EXTERNAL TABLE logdata_parquet  
  
DROP TABLE DimCustomer  
DROP TABLE DimProduct  
DROP TABLE FactSales  
DROP TABLE pool_logdata  
DROP TABLE pool_logdata_parquet
```

15. Now you can see that we have cleaned all of the tables.

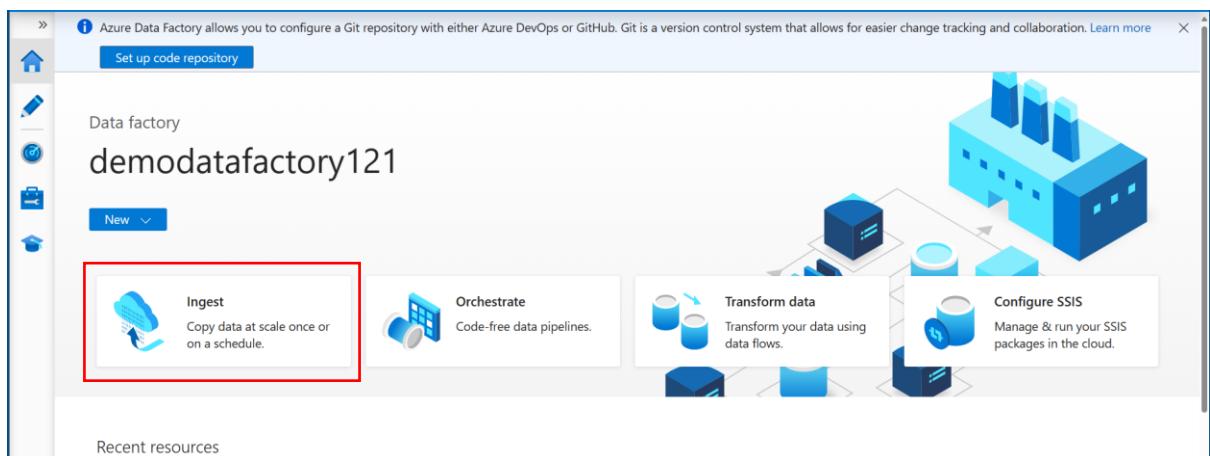


16. Now we are going to create a new table in our demo pool db. For that right click on your pool db and start a new query.
17. Then you have to paste the statement here in SSMS. There you can see that you have created a new table.

```
SQLQuery3.sql - dat... (sqladminuser (0))* ↻ X
CREATE TABLE [logdata]
(
    [Correlation id] [varchar](200) NULL,
    [Operation name] [varchar](200) NULL,
    [Status] [varchar](100) NULL,
    [Event category] [varchar](100) NULL,
    [Level] [varchar](100) NULL,
    [Time] [datetime] NULL,
    [Subscription] [varchar](200) NULL,
    [Event initiated by] [varchar](1000) NULL,
    [Resource type] [varchar](1000) NULL,
    [Resource group] [varchar](1000) NULL,
    [Resource] [varchar](2000) NULL
)
150 % ←
Messages
Commands completed successfully.

Completion time: 2024-04-24T12:41:55.9364626+05:30
```

18. Now you have to navigate to Azure Data Factory. From the home section you have to click on Ingest.



19. Here we are going to start our build-in copy task which we have seen in Synapse also. Just click on next.

Copy Data tool

① Properties

Use Copy Data Tool to perform a one-time or scheduled data load from 90+ data sources. Follow the wizard experience to specify your data loading settings, and let the Copy Data Tool generate the artifacts for you, including pipelines, datasets, and linked services. [Learn more](#)

② Source

③ Destination

④ Settings

⑤ Review and finish

Properties

Select copy data task type and configure task schedule

Task type

Built-in copy task
You will get a single pipeline which is capable of smoothly copying data from over 100 different data sources.

Metadata-driven copy task
You will get parameterized pipelines which can read metadata from an external store to load data at a large scale.

You will get single pipeline to quickly copy objects from data source store to destination in a very intuitive manner.

Task cadence or task schedule *

Run once now Schedule Tumbling window

[< Previous](#) [Next >](#) [Cancel](#)

20. After that you have to choose New connection.

Copy Data tool

① Properties

Source data store

Specify the source data store for the copy task. You can use an existing data store connection or specify a new data store.

② Source

Source type

Connection * [+ New connection](#)

③ Dataset

④ Configuration

21. Then choose Azure after that choose Azure data lake storage Gen2. Select it and click on continue.

New connection

Search

All **Azure** Database File Generic protocol NoSQL Services and apps

 Azure Blob Storage	 Azure Cosmos DB for MongoDB	 Azure Cosmos DB for NoSQL
 Azure Data Explorer (Kusto)	 Azure Data Lake Storage Gen2	 Azure Database for MariaDB

22. Then you have to give it a name, you can copy your storage account name in which your containers are.

New connection

 Azure Data Lake Storage Gen2 [Learn more](#) 

Name *

Description

Connect via integration runtime * 



Authentication type



23. After that choose your subscription and choose your storage account in which you have your contains. Test connection should be linked service. Then hit on create.

Account selection method 

From Azure subscription Enter manually

Azure subscription 



Storage account name *



Test connection 

To linked service To file path

24. Once you are connected then you have to browse for you container in which you have CSV file.

Source data store

Specify the source data store for the copy task. You can use an existing data store connection or specify a new data store.

Source type



Connection *

[!\[\]\(d3b37595da467182627cf445f6121bf9_img.jpg\) Edit](#) [!\[\]\(f80d926d1b0cbf6f8225b43b8466a488_img.jpg\) New connection](#)

File or folder

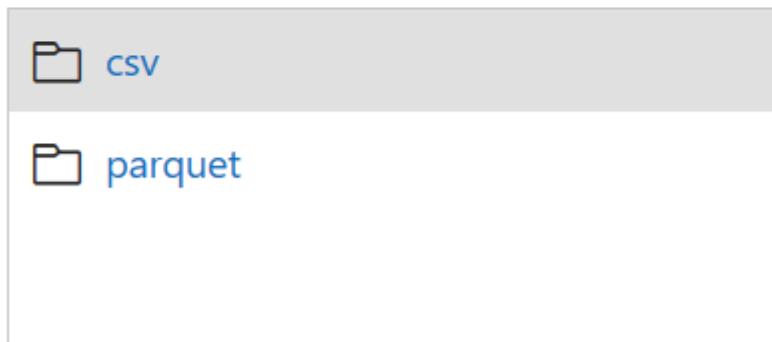
If the identity you use to access the data store only has permission to subdirectory instead of the entire account, specify the path to browse. Append a slash (/) at the end if the path refers to a folder.

25. Choose your CSV container and go inside of it then choose your Log.csv file.

Browse

Select a file or folder.

Root folder



26. Below you can see that we have our file in place.

Source data store

Specify the source data store for the copy task. You can use an existing data store connection or specify a new data store.

Source type	All
Connection *	<input type="button" value="sq storage1010_service"/> Edit <input type="button" value="New connection"/>
File or folder	<input type="text" value="csv/Log.csv"/> <input type="button" value="Browse"/>

27. Then click on the next button here you will see that the wizard has identified the underlying format for our CSV file. You can also preview your data to check whether the tool has your proper data or not. After that click on next.

Copy Data tool

Properties

Source

Dataset

Configuration

Destination

Settings

Review and finish

File format settings

File format: DelimitedText

Column delimiter: Comma (,)

Row delimiter: Default (\r\n, or \n)

First row as header

Compression type: Select...

Additional columns: New

< Previous Next >

Cancel

28. Now for the destination you have to click on the new connection again and this time you will choose your Synapse Workspace.

Destination data store

Specify the destination data store for the copy task. You can use an existing data store connection or specify a new data store.

Destination type	All
Connection *	Select... <input type="button" value="New connection"/>

29. Here in all section you have to scroll down and find Azure Synapse Analytics

New connection

Search

All Azure Database File Generic protocol NoSQL Services and apps

 Azure Synapse Analytics	 Azure Table Storage	 Dataverse (Common Data Service for Apps)
--	--	---

30. First you are going to give it the same name as your Synapse Workspace.

New connection

 Azure Synapse Analytics [Learn more](#) 

Name *

Description

Connect via integration runtime * 

Connection string

[Azure Key Vault](#)

31. Now here you have to choose your subscription then your server name of Synapse workspace then choose your Pool db.

Account selection method 

From Azure subscription Enter manually

Azure subscription

Server name *



Database name *



SQL pool *

 demopooldb

32. Then in the authentication type choose SQL Authentication and then give username of your Synapse Analytics and then give your password.

Authentication type *

SQL authentication
 ▼

User name *

sqladminuser

Password
Azure Key Vault

Password *

33. After that we will test our connection. Once your connection is successful then click on create.



34. Now you have to choose your existing table for that click on use existing table.

Destination data store

Specify the destination data store for the copy task. You can use an existing data store connection or specify a new data store.

Destination type	<div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;">All</div> <div style="border: 1px solid #ccc; padding: 2px; display: flex; align-items: center;"> 🔗 datasynapse1234_pooldb Edit + New connection </div>						
<table border="0" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 30%;">Source</th> <th style="width: 40%; text-align: center;">Destination</th> <th style="width: 30%;"></th> </tr> </thead> <tbody> <tr> <td style="vertical-align: middle;">Log</td> <td style="text-align: center; vertical-align: middle;"> → <div style="border: 1px solid #ccc; padding: 2px; margin-right: 10px;">dbo</div> . Log (auto-create) </td> <td style="text-align: center; vertical-align: middle;"> * Use existing table </td> </tr> </tbody> </table>		Source	Destination		Log	→ <div style="border: 1px solid #ccc; padding: 2px; margin-right: 10px;">dbo</div> . Log (auto-create)	* Use existing table
Source	Destination						
Log	→ <div style="border: 1px solid #ccc; padding: 2px; margin-right: 10px;">dbo</div> . Log (auto-create)	* Use existing table					

35. Below you can see that we have chosen our existing table. Then click on next.

Source	Destination	
▼ Log	→ <div style="border: 1px solid #ccc; padding: 2px; margin-right: 10px;">dbo.logdata</div> <div style="border: 1px solid #ccc; padding: 2px; margin-right: 10px;">↻</div>	Auto-create a destination table with the source schema

36. Now it will do a column mapping for your table. Click on next.

The screenshot shows the 'Copy Data tool' interface. On the left, a vertical navigation bar lists steps: Properties, Source, Destination (highlighted in blue), Dataset, Configuration, Settings (with a gear icon), and Review and finish. The main area is titled 'Column mapping' with the sub-section 'Table mappings (1)'. It shows a mapping from 'Source' (Azure Data Lake Storage Gen2 file) to 'Destination' (dbo.logdata). Below this, a table lists columns being mapped:

Source	Type	Destination
Correlation id	abc String	Correlation id
Operation name	abc String	Operation nar
Status	abc String	Status
Event category	abc String	Event categor
Level	abc String	Level
Time	abc String	Time

A warning message at the bottom states: '⚠️ Copying from column Time to column Time may have data'.

At the bottom of the screen are buttons for '< Previous' and 'Next >'.

37. Then give it a name and scroll down.

Settings

Enter name and description for the copy data task, more options for data movement

Task name *	<input type="text" value="01-Copy-data-to-logdata"/>
Task description	<input type="text"/>

38. Here you have to disable staging and choose bulk insert.

Settings

Enter name and description for the copy data task, more options for data movement

Fault tolerance ⓘ	<input type="text"/>
Enable logging ⓘ	<input type="checkbox"/>
Enable staging ⓘ	<input type="checkbox"/>
✓ Advanced	
Copy method	<input type="radio"/> Copy command ⓘ <input type="radio"/> PolyBase ⓘ <input checked="" type="radio"/> Bulk insert <input type="radio"/> Upsert
Bulk insert table lock ⓘ	<input type="radio"/> Yes <input checked="" type="radio"/> No

39. After that click on next and move to the review page and deploy your pipeline. Once it is successfully deployed then click on finish.

The screenshot shows the 'Copy Data tool' interface. On the left, a vertical navigation bar lists steps: Properties (checkmark), Source (checkmark), Destination (checkmark), Settings (checkmark), Review and finish (highlighted with a blue circle), Review (blue dot), and Deployment (blue dot). The main area displays a flow diagram from 'Azure Data Lake Storage Gen2' to 'Azure Synapse Analytics'. Below the diagram, the text 'Deployment complete' is centered. To the right, a table titled 'Deployment step' shows four steps all in 'Succeeded' status: 'Validating copy runtime environment', 'Creating datasets', 'Creating pipelines', and 'Running pipelines'. A note at the bottom states: 'Datasets and pipelines have been created. You can now monitor and edit the copy pipelines or click finish to close Copy Data Tool.' At the bottom are three buttons: 'Finish' (blue), 'Edit pipeline' (white), and 'Monitor' (white).

40. Now move to the monitor section and here you can see that your pipeline run was successful.

The screenshot shows the Azure Data Factory 'Monitor' section. The left sidebar has 'Monitor' selected. The main area shows a list of 'Pipeline runs'. One run is highlighted: '01-Copy-data-to-logdata - Activity runs'. The run details show a 'Copy data' activity named 'Copy_sl6' with a green checkmark indicating success. Below this, the 'Activity runs' table shows one entry: 'Copy_sl6' with status 'Succeeded', activity type 'Copy data', run start '4/24/2024, 1:05:41 PM', and duration '21s'. Navigation buttons like 'Rerun', 'Cancel', 'Refresh', 'Update pipeline', 'List', and 'Gantt' are visible at the top of the run details.

41. Once it is done now move to SSMS there you have to run the Select query and check for your data.

42. Below you can see your data.

SELECT * FROM [logdata]

	Correlation id	Operation name	Status	Event category	Level	Time
1	96600577-d906-4df2-8cd0-69c2872c7fd	Create or Update Virtual Machine Extension	Started	Administrative	Informational	2023-04-17 06:31:00.723
2	c43307c3-015b-4c18-9556-456410c313de	Delete Network Security Group	Succeeded	Administrative	Informational	2023-04-17 06:07:51.210
3	70e3851e-079b-40d7-bb87-31f5386befcc	Delete Disk	Succeeded	Administrative	Informational	2023-04-17 06:04:33.950
4	cbbe3740-5009-4821-883c-cac8fd3c3659	'audit' Policy action.	Succeeded	Policy	Warning	2023-04-17 05:46:03.813
5	cbbe3740-5009-4821-883c-cac8fd3c3659	'audit' Policy action.	Succeeded	Policy	Warning	2023-04-17 05:46:02.207
6	f4c54d30-0b68-4c60-94fc-69349e63b788	'auditIfNotExists' Policy action.	Succeeded	Policy	Informational	2023-04-17 05:44:07.787
7	304f186a-198f-4c52-8a1c-1baaa98ceec9cd	'audit' Policy action.	Succeeded	Policy	Warning	2023-04-17 05:41:07.000
8	304f186a-198f-4c52-8a1c-1baaa98ceec9cd	'auditIfNotExists' Policy action.	Succeeded	Policy	Informational	2023-04-17 05:41:05.563
9	bc53c993-3e88-4a7e-a4dd-30f01ca14790	Validate Deployment	Succeeded	Administrative	Informational	2023-04-17 05:33:32.583
10	a56c2662-433b-4cd6-ba1a-8f18c024271a	Delete Network Interface	Succeeded	Administrative	Informational	2023-04-16 14:13:41.823
11	77903c28-3455-4d4f-8e03-fb0dd13d6181	Create Flow Log	Succeeded	Administrative	Informational	2023-04-16 12:39:19.997
12	530772f7-90a4-4741-89ac-4941e02b2ca2	Create Connection Monitor	Started	Administrative	Informational	2023-04-16 11:47:31.090
13	a265ddd3-d275-46ec-9fb1-f2633cf37b90	'auditIfNotExists' Policy action.	Succeeded	Policy	Informational	2023-04-16 06:18:40.863
14	420e3dc7-c154-49b1-9e49-f2f44de74cd6	'audit' Policy action.	Succeeded	Policy	Warning	2023-04-16 06:14:42.030
15	420e3dc7-c154-49b1-9e49-f2f44de74cd6	Create or Update Virtual Machine Extension	Succeeded	Administrative	Informational	2023-04-16 06:07:46.280
16	7907a8a7-2b1c-4f22-b1f0-2c580435fde4	Health Event Updated	Updated	Resource Health	Informational	2023-04-16 06:04:43.103
17	cf800f9f-2782-4e7e-8060-3302d96fadf8	Delete Disk	Started	Administrative	Informational	2023-04-14 15:40:14.787
18	78e56a68-af47-42db-8d63-a0afa9930652	Get Network Interface Effective Security Groups	Started	Administrative	Informational	2023-04-14 15:23:11.450
19	a3ff6da1-818b-4449-b23b-c7b14a944513	'audit' Policy action.	Succeeded	Policy	Warning	2023-04-14 15:07:56.157

Query executed successfully.

43. Now if you come back to the data factory and go to author you can see your pipeline which was created earlier and related to that there are several activities. Then in the datasets, you can see two datasets in place one is for the destination and the other is for the source.

The screenshot shows the Azure Data Factory Author interface. On the left sidebar, 'Author' is selected. In the main area, a pipeline named '01-Copy-data-to-ld...' is open. The pipeline contains a single activity named 'Copy data'. The 'Source' for this activity is set to 'Copy_sl6' and the 'Sink' is also set to 'Copy_sl6'. Below the pipeline, the 'Activities' pane shows a list of available activities, including 'Move and transform', 'Synapse', 'Azure Data Explorer', etc. On the left, under 'Factory Resources', the 'Datasets' section is highlighted with a red box, showing 'DestinationDataset_sl6' and 'SourceDataset_sl6'. Other sections like 'Pipelines', 'Data flows', and 'Power Query' are also listed.