

Data Ingestion Patterns – Hands On: Copy data from on-premises to ADLS

Updated on: 12/27/2016

Introduction

As customer adoption of Azure Data Lake Store increases, there will be many who will have data stored on on-premises solutions. These customers will want that on-premises data to be moved to Azure Data Lake Store in a one-time or recurring fashion. **Azure Data Factory** provides you an easy way to achieve this data copy in an intuitive, scalable, recurring fashion with appropriate data transformation.

In this lab, learn the basics of how to copy data from file shares stored on-premises to ADLS in a one-time fashion using Azure Data Factory and Data Management Gateway.

You will start with gzipped files stored on-premises (on your own machine) and copy them to Azure Data Lake Store via **Azure Data Factory** and **Data Management Gateway** in an uncompressed form for further analysis.

Prerequisites

- Machine running the latest version of Windows.
- To be logged-in to the machine using your domain credentials
- Access to an Azure subscription that you can create resources in.
- Access to an ADLA Account that you can submit jobs to.
- Access to an ADLS Account that you can submit jobs to.

Things to Keep in Mind

- File and folder paths in Azure Data Lake Store are case-sensitive
- If you are logging in with Microsoft credentials, please go to <http://aka.ms/AzurePortalProd> to get the same version of the portal that customers see.
- You can use any browser, but for the purposes of this lab using Microsoft Edge or Microsoft Internet Explorer will be simpler since they natively support the ClickOnce feature

Step 1: Get the source files ready on your machine

- Create a folder at the root of a drive called **C:\ADLSTutorial**.
 - Right click on the **ADLSTutorial** folder, then select **Share With > Specific People**.
 - Choose your own user account
 - click **Share**.

- Inside that shared folder, create a folder called **Source** that will contain the files to be copied.
- Download the SmallDataset.zip file from [here](#)
- Extract the ZIP. It will contain two .gz files.
- Put the contents of the zip into the Source folder.

Your folder structure should look like

```
C:\
  ADLSTutorial\
    Souce\
      commitssmall.csv.gz
      userssmall.csv.gz
```

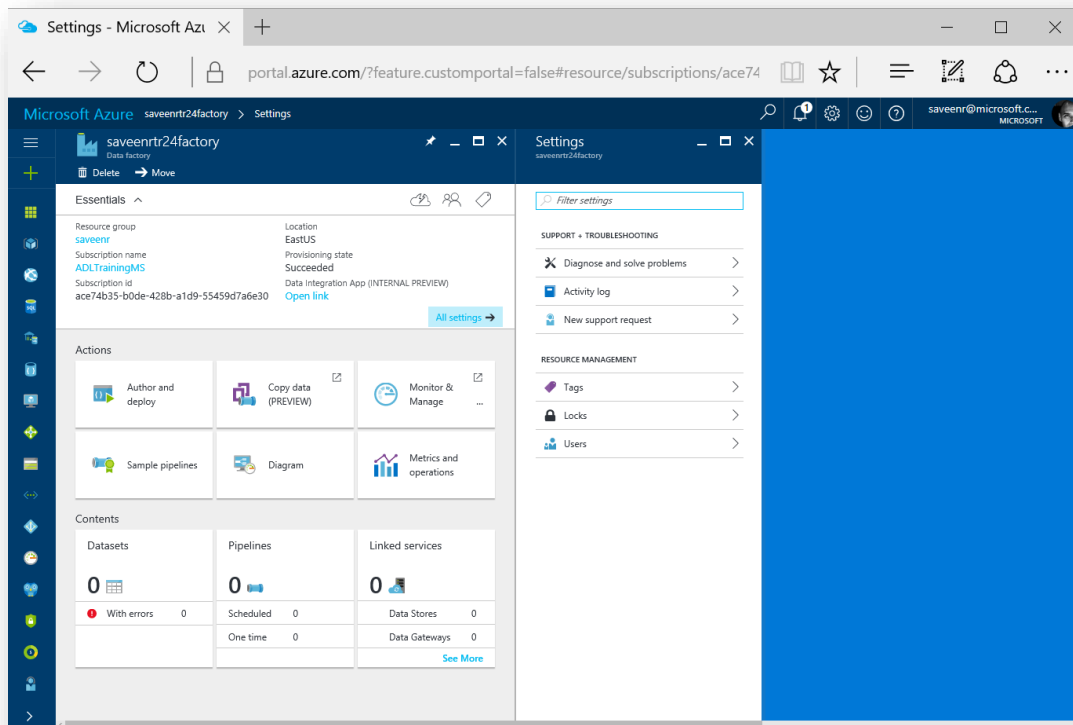
Step 2: Create the destination folder in ADLS

- Log into the Azure portal <http://portal.azure.com>
- Open your ADLS Account and create a folder called **/TR24**

Step 3: Create the Azure Data Factory

NOTE: The name of the Azure data factory must be globally unique. Try using something like "ADLSTutorial-
<username>"

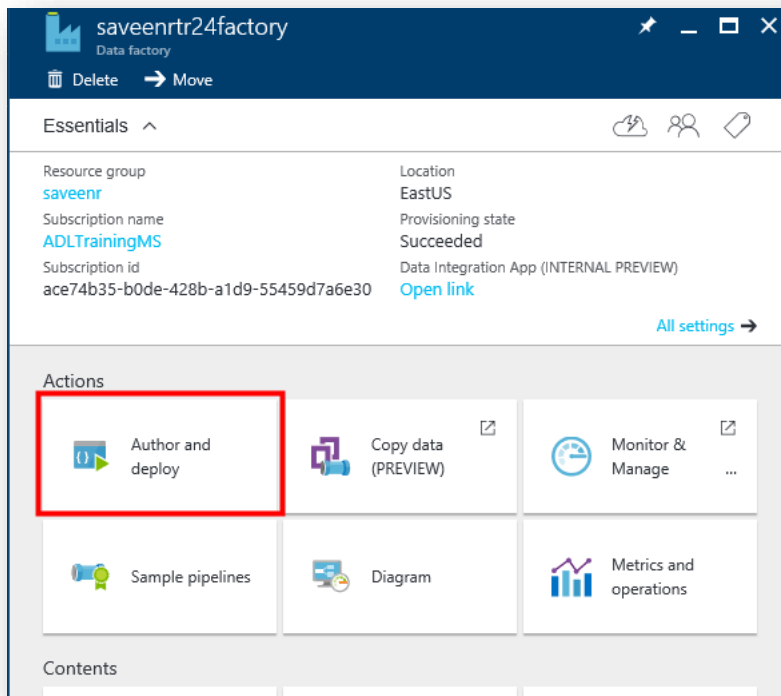
- Log into the Azure portal <http://portal.azure.com>
- click **New > Intelligence + Analytics > Data Factory** on the **Data + Analytics** blade. The **New data factory** blade will appear.
- In the **New data factory** blade:
 - Enter **ADLSTutorial-<username>** for the **name**.
 - For **Resource Group**, use whatever you want
 - For **Location**, leave with the default selection.
- Enable **Pin to dashboard** checked box.
- Click **Create**
- Once creation is complete, you will see the **Data Factory** blade as shown below:



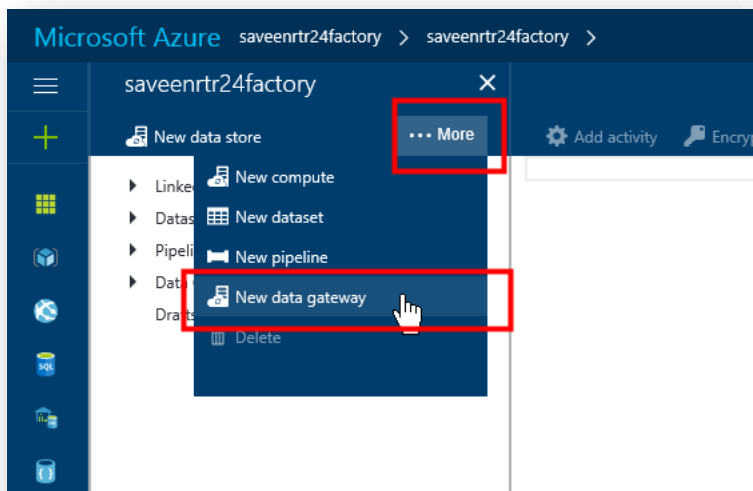
Step 4: Create a Data Management Gateway

IMPORTANT: if Data Management Gateway is already installed on your machine, uninstall it

- Navigate to **Data Factory** that you created
- Click **Author and deploy** tile to launch the **Editor** for the data factory.

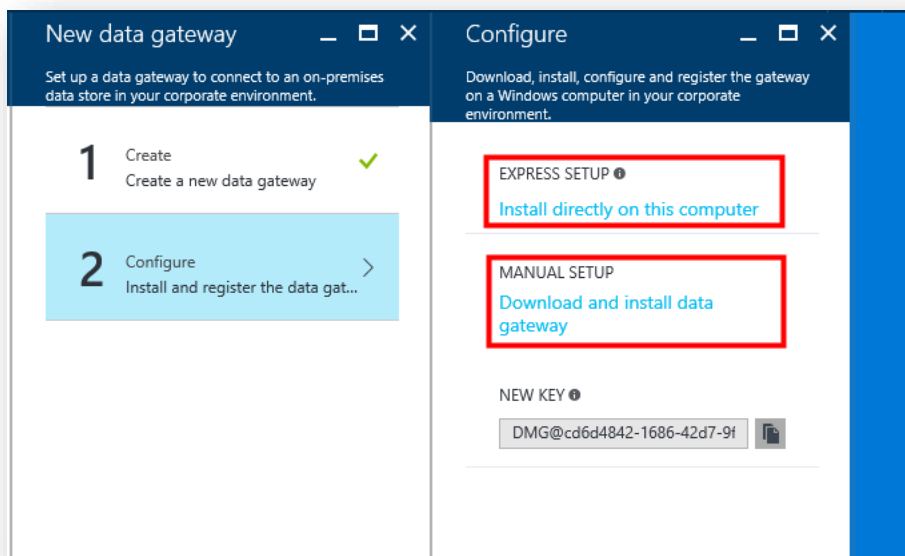


- In the Data Factory Editor, click ... (ellipsis) on the toolbar
- Click **New data gateway**



- The **New data gateway** blade will appear
- In the **Create** step, for the **Data gateway name**, enter **ADLSGateway-<username>-TR24**
- Click **OK**.

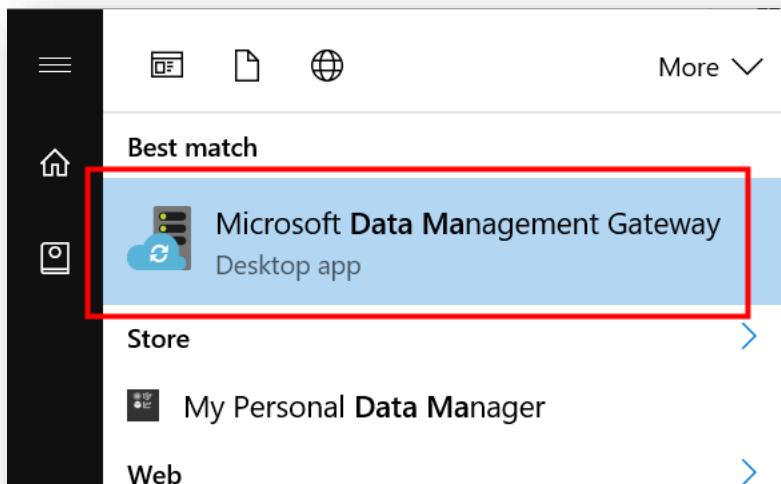
- In the **Configure** step, you'll have two ways of installing DMG: EXPRESS and MANUAL.
 - a. EXPRESS: Under **EXPRESS SETUP** click **Install directly on this computer**. This will download the installation package for the gateway, install, configure, and register the gateway on the computer using ClickOnce. This is the easiest way of installing DMG but please verify you have met the requirements for ClickOnce as defined in the appendix of this document.
 - b. MANUAL: Under MANUAL SETUP click on **Download and install data gateway**.
- NOTE: You must be an administrator on the local computer to install and configure the Data Management Gateway successfully. You can add additional users to the Data Management Gateway Users local Windows group. The members of this group will be able to use the Data Management Gateway Configuration Manager tool to configure the gateway.



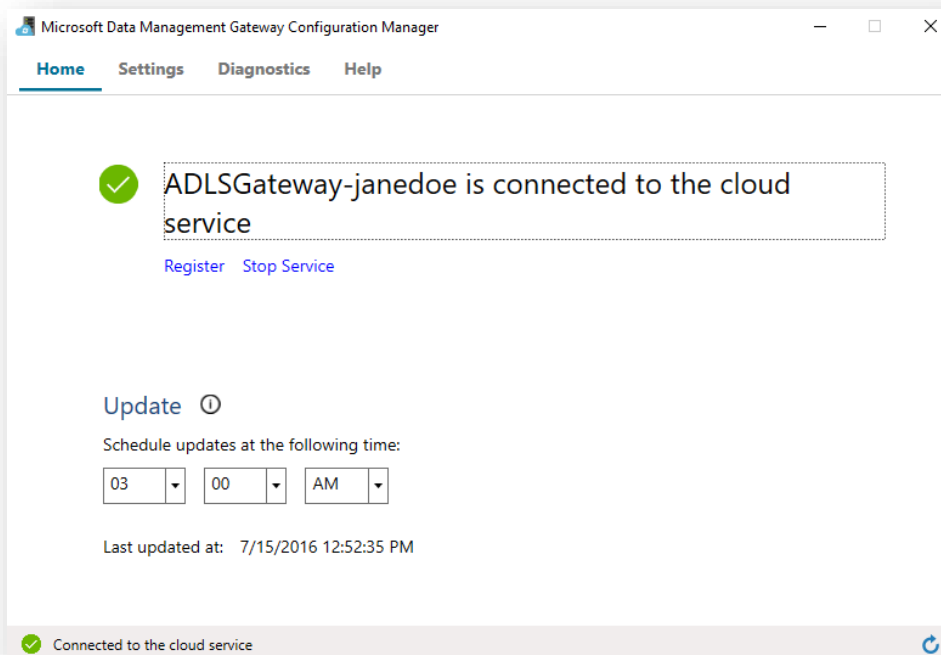
- To verify that DMG has been installed, check that the executable **ConfigManager.exe** exists in the folder: **C:\Program Files\Microsoft Data Management Gateway\1.0\Shared**.

NOTE: You must be an administrator on the local computer to install and configure the Data Management Gateway successfully. You can add additional users to the Data Management Gateway Users local Windows group. The members of this group will be able to use the Data Management Gateway Configuration Manager tool to configure the gateway.

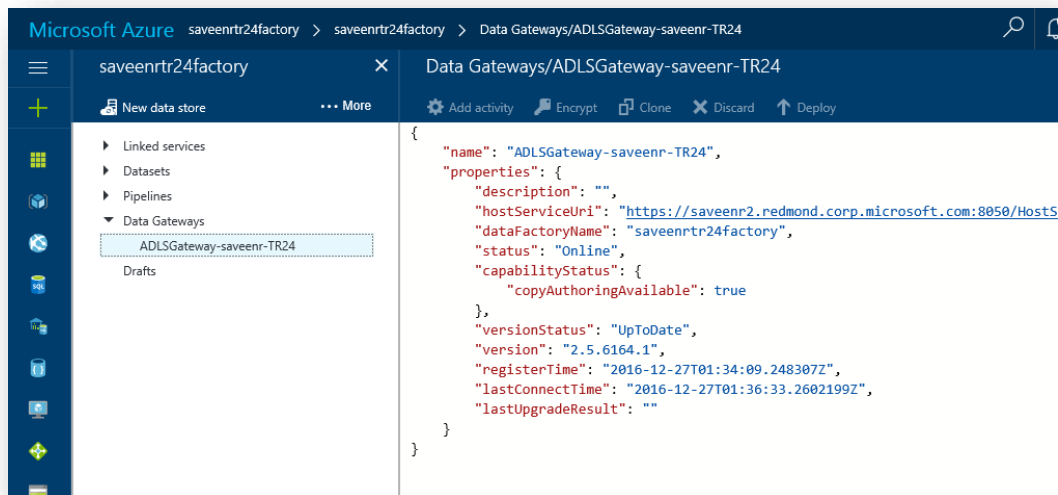
- Wait for a couple of minutes and ensure that Data Management Gateway is installed correctly. Once you see it in the start menu, click on it



- When DMG Starts

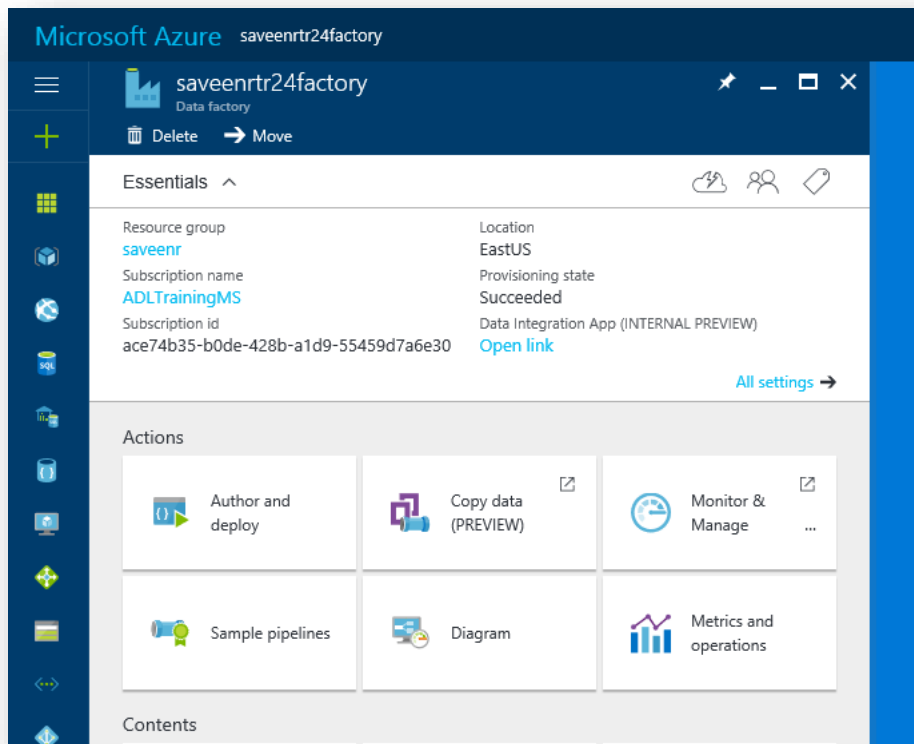


- Wait until the you see the message that it is **connected to the cloud service** as shown above.
- Go back to the Azure Portal, click **OK** on the **Configure** blade and then on the **New data gateway** blade.
- You should see **ADLSGateway- <your alias>** under **Data Gateways** in the tree view on the left. If you click on it, you should see the associated JSON.



Step 5: Create and run a pipeline using Copy Wizard

- Go to the Azure Portal
- Navigate to the ADF factory that you created
- Click **Copy data...** The **COPY DATA – Properties** page will appear.



Copy Data (saveenrtr24factory)

1 Properties
2 Source
3 Destination
4 Summary

Properties

Enter name and description for the copy data task and specify how often you want to run the task.

Task name (required) i

Task description

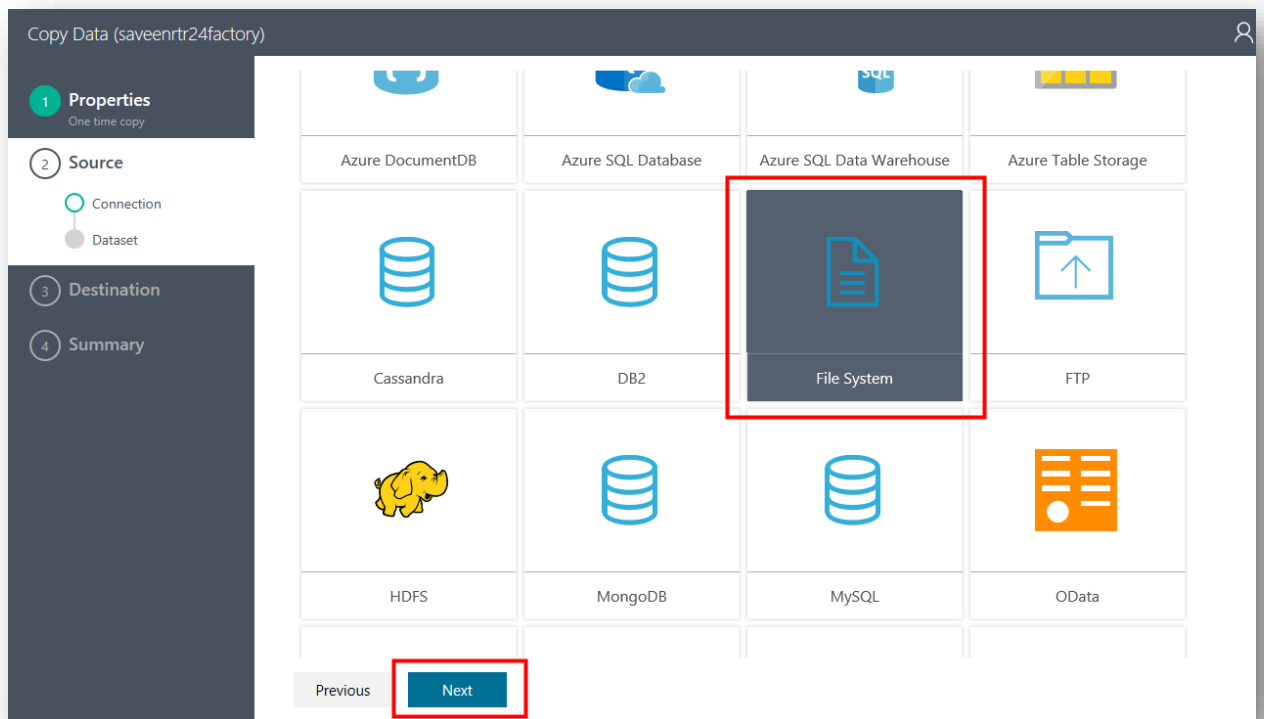
Task cadence (or) Task schedule
☐ Run once now
☒ Run regularly on schedule

Recurring pattern
 every day

Start date time (UTC)

End date time (UTC)

- Task name – Enter a name that you will be familiar with, otherwise note down the name that is prefilled e.g. CopyPipeline-l17 here. Note down this name on a piece of paper, you will need it later to monitor this pipeline.
- Click on **Run once now**. You can keep all the other values as the default ones.
- Click on **Next**.
- You will be taken to the **COPY DATA – Source data store** page.
- On the **COPY DATA – Source data store** page, click on **File System** on the Connect to a Data Store screen.
- Click on **Next**.



- You will be taken to the **COPY DATA – Specify File server share connection** page.

Copy Data (saveenrtr24factory)

1 Properties
One time copy

2 Source

3 Destination

4 Summary

Connection

Dataset

Specify File server share connection

Connection name (required)

Source-File-m72

Gateway (required)

ADLSGateway-saveenr-TR24

Create gateway

Path (required)

e.g. C:\[Folder], \\ServerName\SharedFolder\[Folder]

Credential encryption

None

User name (required)

Password (required)

Previous Next

- On the **COPY DATA - Specify File server share connection** page, as shown above, enter following values for the field.
- **Linked service name** – Leave it as is.
- **Path** – Full path for the share you created above i.e. [\\janedoelaptop\ADLSTutorial](#)
- **Credential encryption** – set to **Use credential manager**.
- **Gateway** – It will be populated automatically with the one you created above i.e. ADLSGateway-<username>
- Click on **Set credential**.
- Allow the CredentialManager.application to be installed.

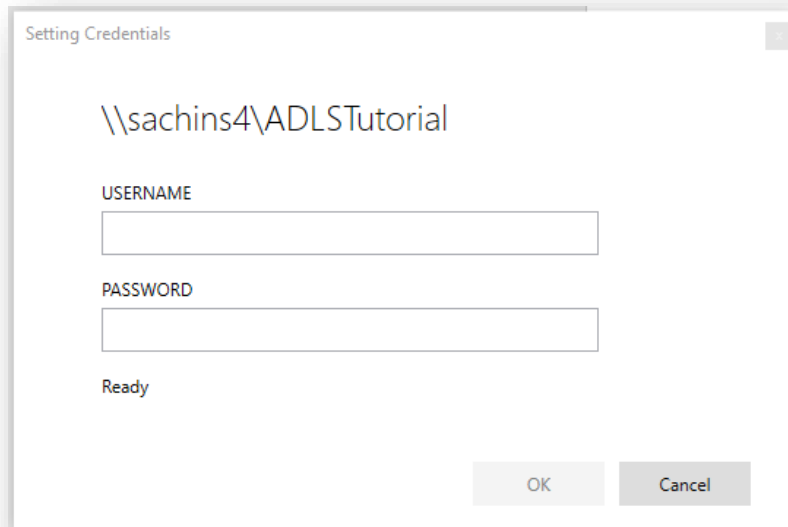
Encrypted credential (required)

Set credential

Do you want to open CredentialsManager.application (15.5 KB) from prod2.portal.clouddatahub.net?

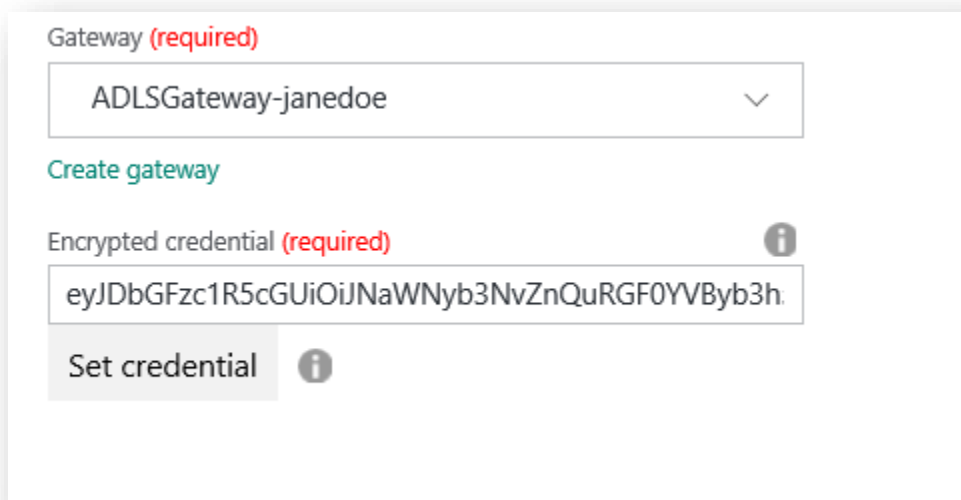
Open Cancel

- You will be shown the **Setting Credentials** screen, where you will enter the credentials that allow access to the share. Since you are using a domain joined redmond machine, your Microsoft credentials should work i.e. username= <domain>\<username> and appropriate password.



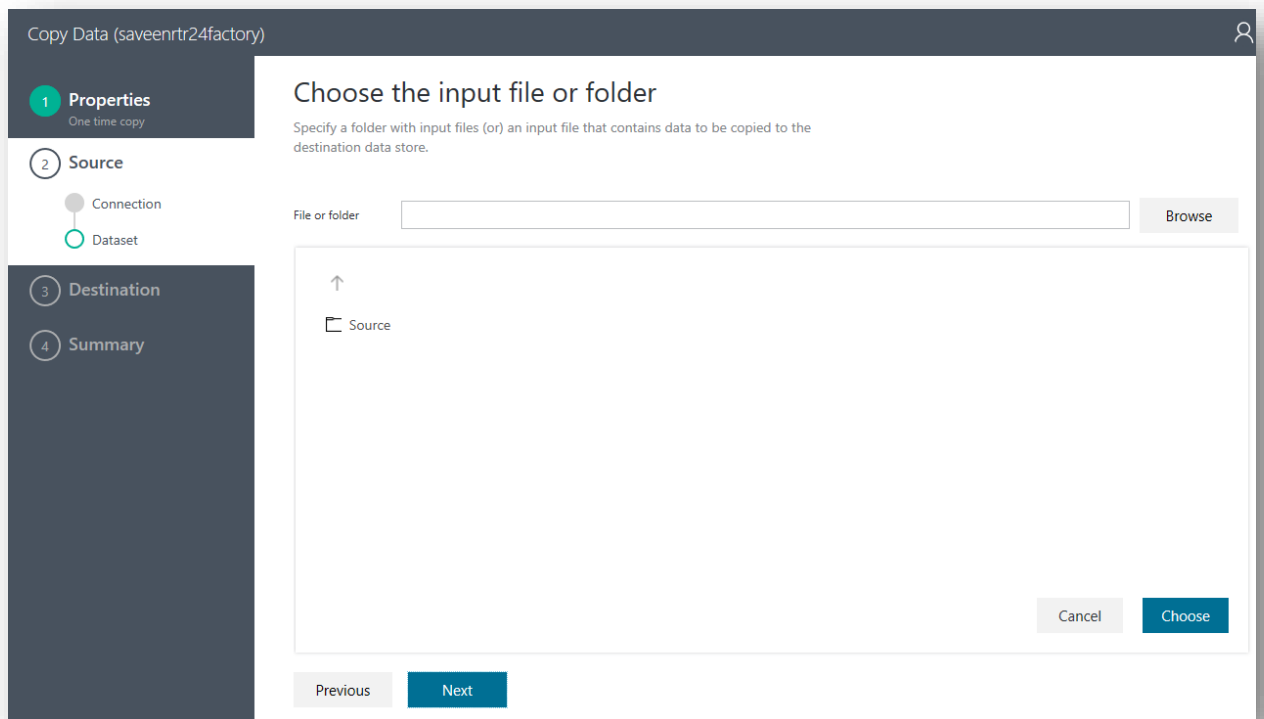
The image shows a Windows 'Setting Credentials' dialog box. At the top, it displays the path '\\sachins4\ADLSTutorial'. Below this, there are two input fields: 'USERNAME' and 'PASSWORD'. The 'Ready' status is shown at the bottom left. At the bottom right, there are 'OK' and 'Cancel' buttons.

- Click on **OK** after entering the credentials.
- If you entered your credentials correctly, you should see the **Encrypted credential** field appropriately filled.



The image shows a configuration screen for a gateway. At the top, it says 'Gateway (required)' in red. Below this is a dropdown menu showing 'ADLSGateway-janedoe'. There is a 'Create gateway' link in green. Below that, it says 'Encrypted credential (required)' in red, followed by an information icon. A text box contains the encrypted credential: 'eyJDbGFzc1R5cGUiOiJNaWNyb3NvZnQuRGF0YVByb3h.'. At the bottom, there is a 'Set credential' button with an information icon.

- Click on **Next**. You will be taken to the **COPY DATA – Choose the input file or folder** page.



- On the **Copy Data – Choose the input file or folder** page, as shown above, click on the **Source** folder that contains the gzip files and then click on **Choose**. On the same page, you will then be provided additional options you can set for the copy.

COPY DATA

Choose the input file or folder

Specify a folder with input files (or) an input file that contains data to be copied to the destination data store.

File or folder [Browse](#)

☒ Copy files recursively ⓘ
☐ Binary copy ⓘ

[Previous](#) [Next](#)

- For this tutorial, click on **Copy files recursively** option. This is to ensure that the all the data in the folder is copied over.
- Click on **Next**.
- You will be taken to the **COPY DATA – File Format settings** page.

Copy Data (saveenrtr24factory)

1 Properties
One time copy

2 Source

Connection

Dataset

3 Destination

4 Summary

File format settings

File format i
Text format

Column delimiter i
Comma (,)
☐ Use custom delimiter

Row delimiter i
Carriage Return + Line feed (r/n)
☐ Use custom delimiter

Skip line count i

☐ Column names in the first data row
☒ Treat empty column value as null

Advanced settings

Escape character i

Quote character i

Null value i
\\N

Encoding name i

Compression type i
None

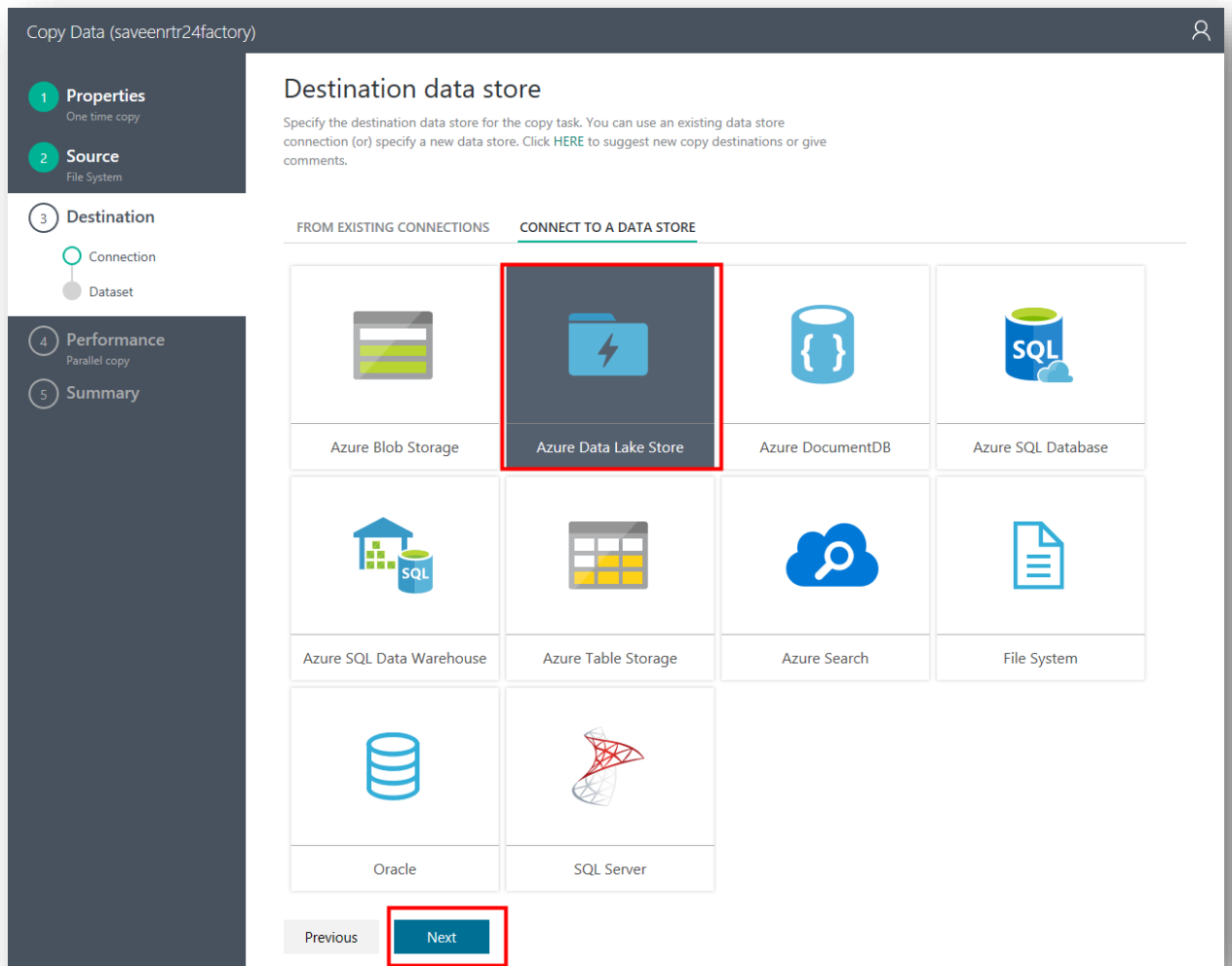
PREVIEW SCHEMA

Error when processing request: Error found when processing 'Csv/Tsv Format Text' source 'commitssmall.csv.gz' with row number 7: found more columns than expected column count: 1. activityId: 6dbbc5b3-46a5-4683-8553-90e52385f564

i Filename: Source/commitssmall.csv.gz [Browse](#)

Previous **Next**

- On the **COPY DATA – File format settings** page as shown above, expand the **Advanced settings** section. Scroll down and select the **Compression type** to **GZip**. You should see a preview of the data in the first source file as shown above.
- Click on **Next**.
- You will be taken to the **COPY DATA – Destination data store** page.



- Here you choose the **Azure Data Lake** as the item.
- Click on **Next**.
- You will be taken to the **COPY DATA – Specify Data Lake Store connection** page.

Copy Data (saveenrtr24factory)

1 Properties
One time copy

2 Source
File System

3 Destination
Connection
Dataset

4 Performance
Parallel copy

5 Summary

Specify Data Lake Store connection

Connection name (required) ?
Destination-DataLakeStore-m72

Azure subscription (required) ?
Select all

Data Lake store account name (required)
abcderfgadls

Previous Next

- On the **COPY DATA – Specify Data Lake Store connection** page, as shown above, enter the following fields:
 - Linked service name** – Leave as is.
 - Data Lake store account name** - Pull down and select the relevant entry.
 - Click on **Next**.
 - You will be taken to the **COPY DATA – Choose the output file or folder** page.
- On the **COPY DATA – Choose the output file or folder** page, shown above, first thing is to choose the **Folder path** where you will land the data. To do this, click on Browse and navigate to TR24/Users/janedoe/ADLSTutorial and select **Destination**. As you will recall, this is the directory that you created in step 2 earlier.

Copy Data (saveenrtr24factory)

1 Properties
One time copy

2 Source
File System

3 Destination
Connection
Dataset

4 Performance
Parallel copy

5 Summary

Choose the output file or folder

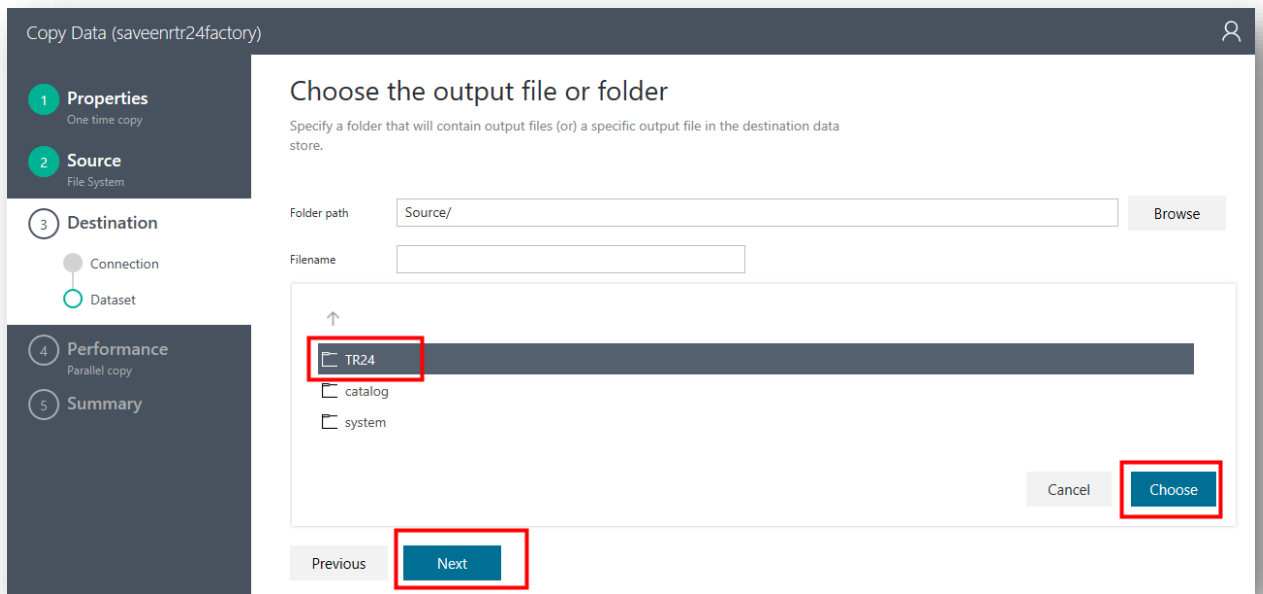
Specify a folder that will contain output files (or) a specific output file in the destination data store.

Folder path Source/ Browse

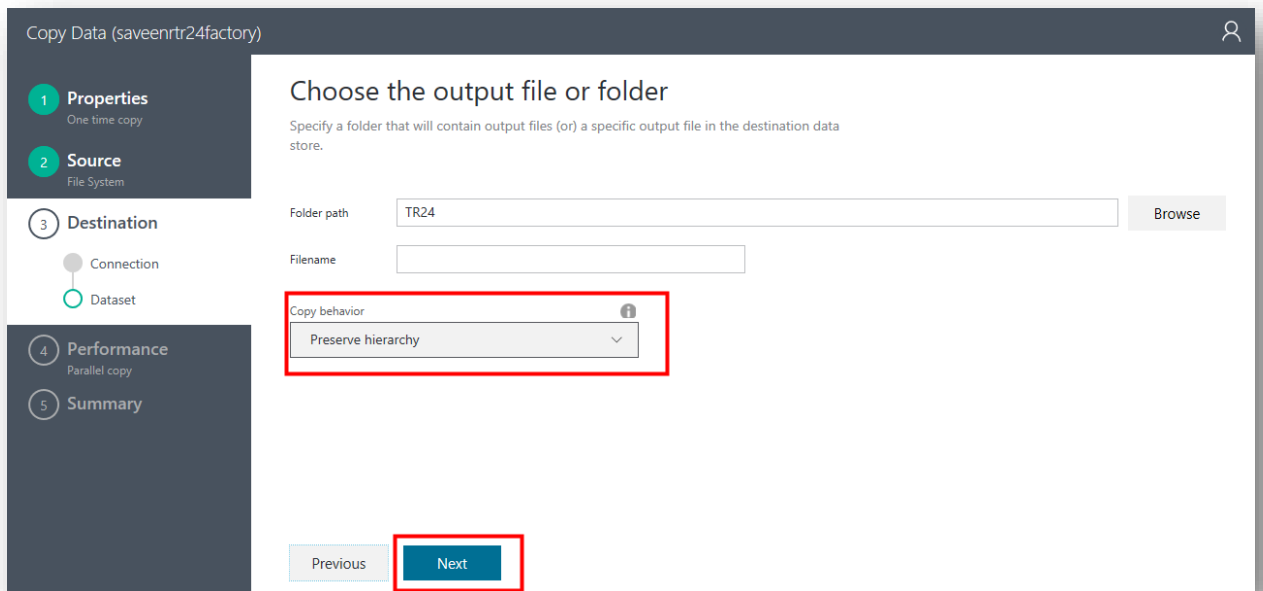
Filename

Copy behavior ?
Merge files

Previous Next



- Click on **Choose**. You will see the page below.



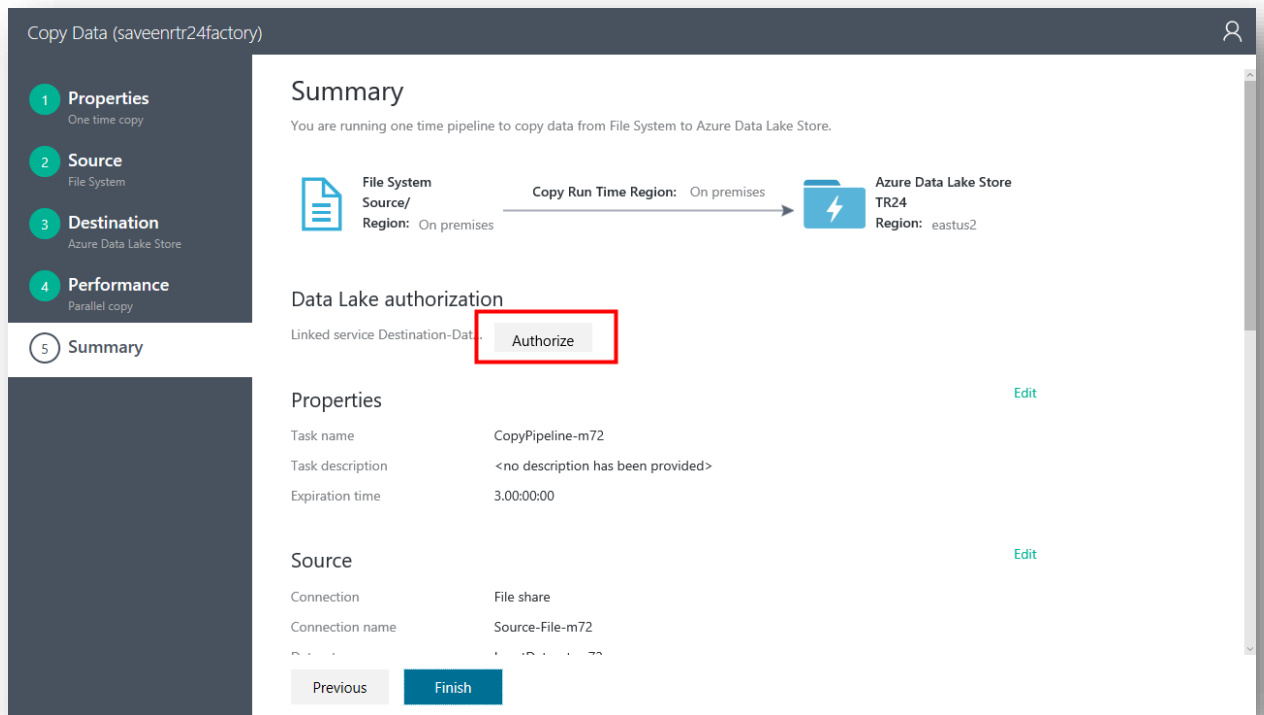
- Leave the **Filename** field as is.
- In the **Copy behavior** pull down, select the “**Preserve hierarchy**” option.
- Click on **Next**. You will be taken to the **COPY DATA – File format settings** page.

The screenshot shows the 'File format settings' page of the 'Copy Data (saveenrtr24factory)' wizard. On the left, a sidebar contains five steps: 1 Properties (One time copy), 2 Source (File System), 3 Destination (Connection and Dataset), 4 Performance (Parallel copy), and 5 Summary. The main area is titled 'File format settings' and includes dropdown menus for 'File format' (Text format), 'Column delimiter' (Comma (,)), and 'Row delimiter' (Carriage Return + Line feed (\r\n)). There are checkboxes for 'Use custom delimiter' and 'Add header to file'. An 'Advanced settings' section is collapsed. At the bottom, 'Previous' and 'Next' buttons are visible, with the 'Next' button highlighted by a red rectangle.

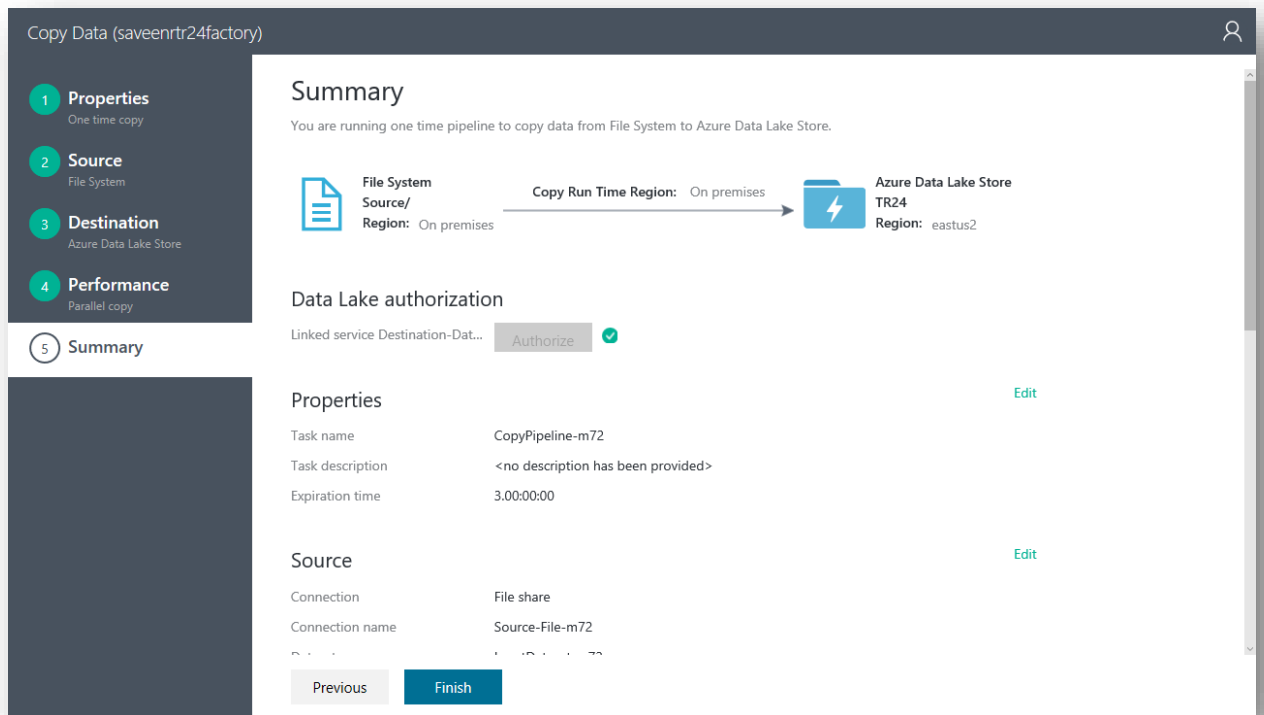
- On the **COPY DATA – File format settings** page, do not select anything. Just click on **Next**.

The screenshot shows the 'Performance settings' page of the 'Copy Data (saveenrtr24factory)' wizard. The sidebar on the left now highlights step 4 Performance (Parallel copy). The main area is titled 'Performance settings' with the subtitle 'Performance improvement options'. It features a collapsed 'Advanced settings' section. At the bottom, 'Previous' and 'Next' buttons are visible, with the 'Next' button highlighted by a red rectangle.

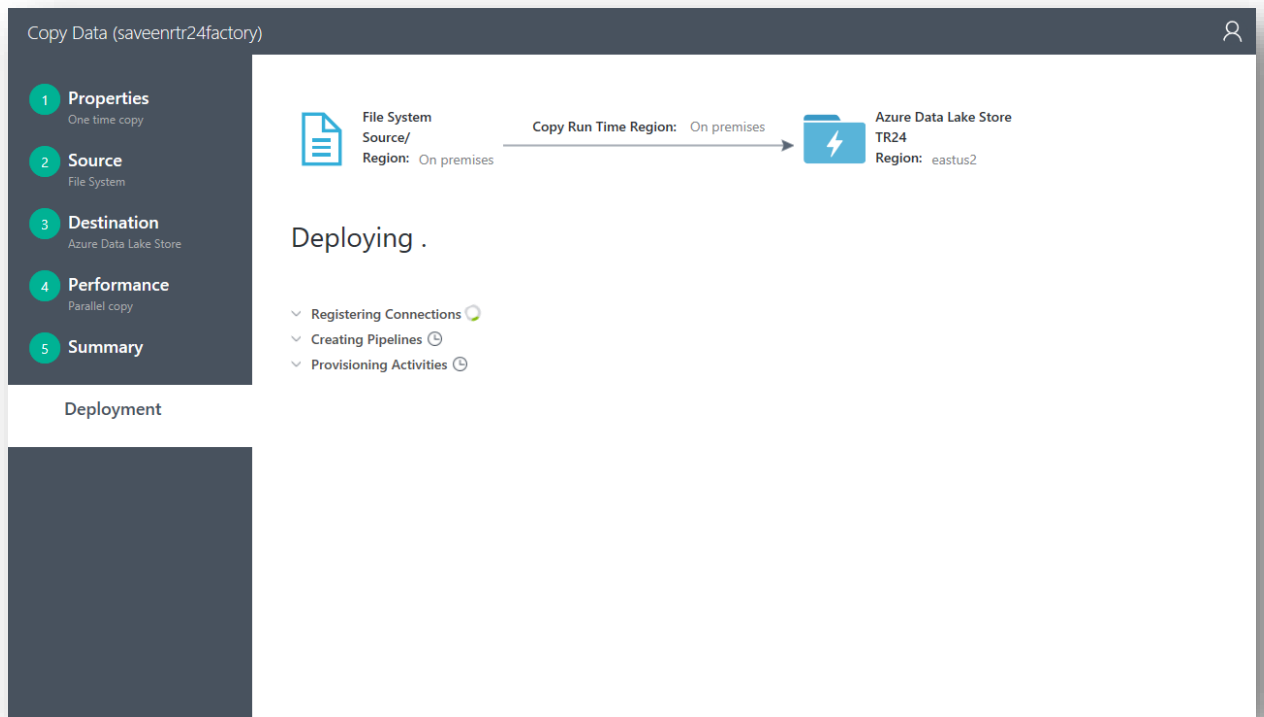
- You will be taken to the **COPY DATA – Summary** page.



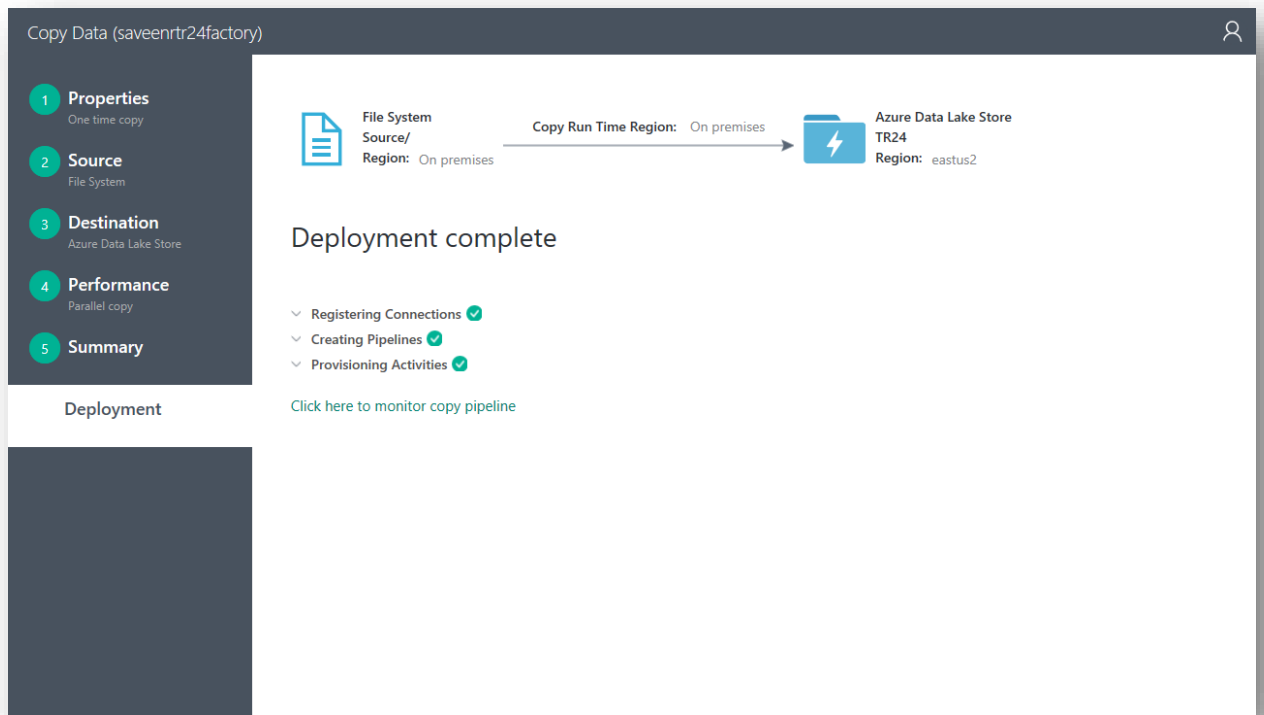
- On the **COPY DATA - Summary** page, as shown above, first click on the **Authorize** button in the **Data Lake authorization** section. This will authorize ADF access to ADLS. You will be prompted for credentials that have access to the ADLS account we are accessing.
- You should see a green checkmark once Authorization has completed successfully.



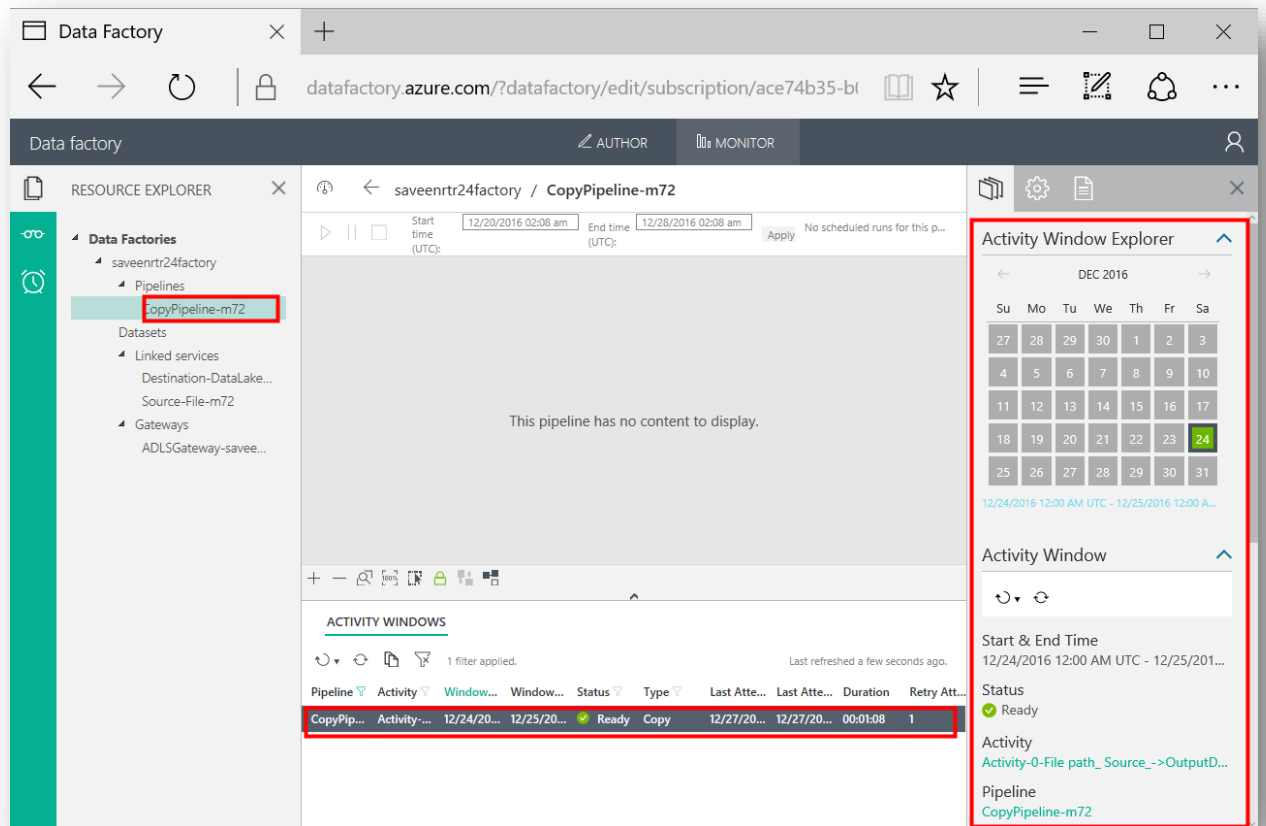
- Click on **Finish** and wait. You will be taken to the **COPY DATA - Deploying** page.



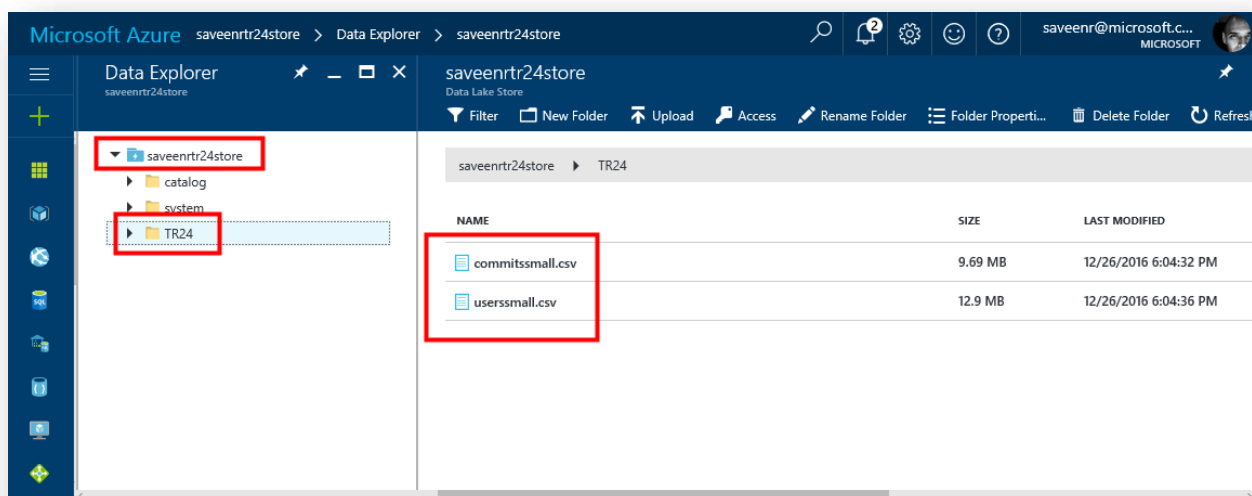
- In the **Deploying** screen where you will see **Registering Connections** and **Creating Pipelines**.
- Once there are green check marks next to those, you will see the title change to Deployment complete.



- At this point, recall the name of the Pipeline that you created earlier e.g. in this tutorial it was **CopyPipeline-I17**. Or you can find it by clicking on the down arrow next to **Creating Pipelines**.
- You will additionally see **Click here to monitor copy pipeline** link. Click on that link. You will be taken to the monitoring UI that monitors all your pipelines.



- Wait for the copy activity to be completed. Once you see it has succeeded in the right pane, go to the **MyStoreAccount** ADLS account via the Azure portal. Using Data explorer navigate to the TR24/users/janedoe/ADLSTutorial/Destination folder.



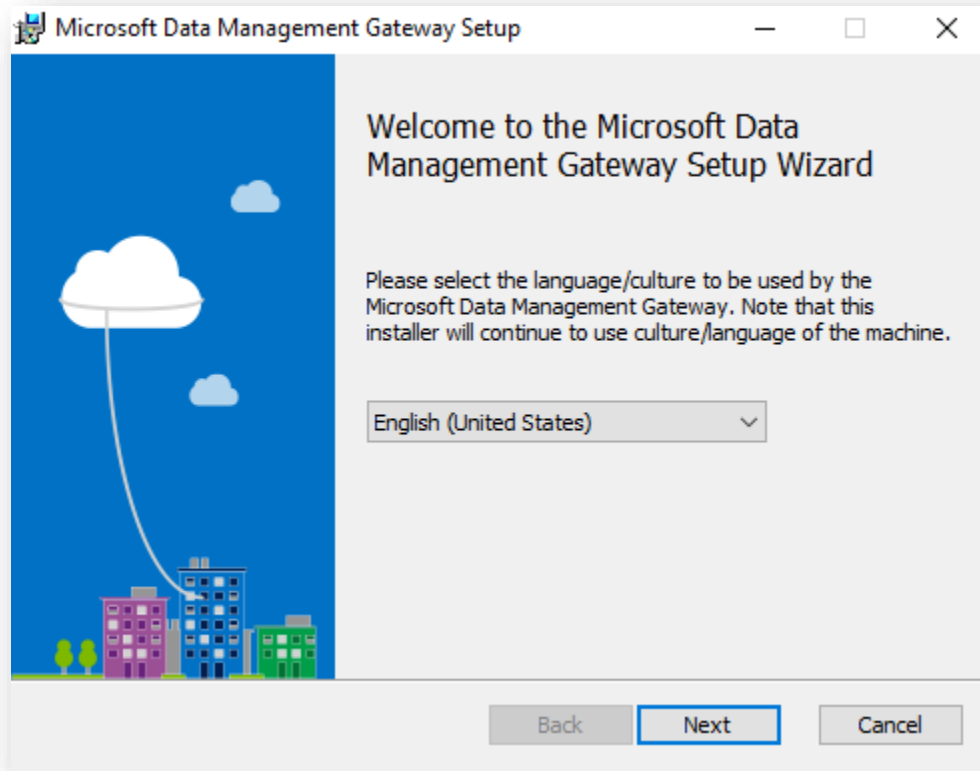
- You should see 2 files in the folder as shown above. Look at the **LAST MODIFIED** column and confirm that your file was recently copied over. Also confirm that the file does not have the “gz” suffix anymore. This shows that the gzip files were successfully unzipped and copied over from your machine.
- **Congratulations!** You have taken multiple gzip files stored in a folder from an on-premises location, and copied them over to Azure Data Lake Store in the cloud in an uncompressed form. You are now ready to process them using HDInsight, Azure Data Lake Analytics or an analytics tool of your choice.

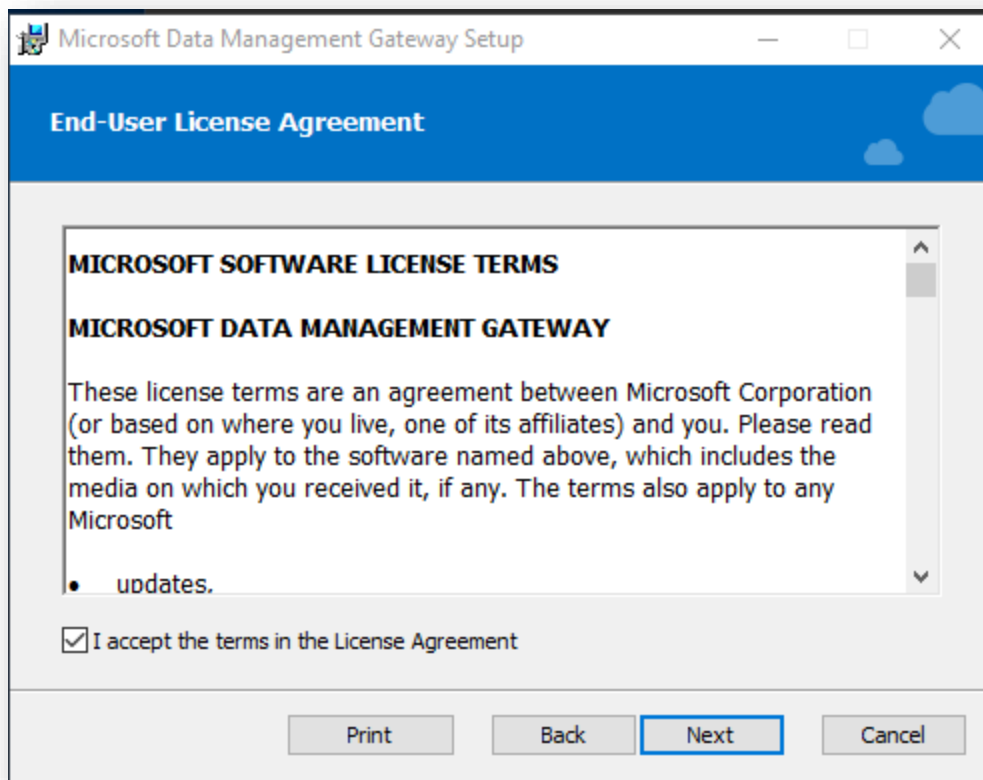
Conclusion

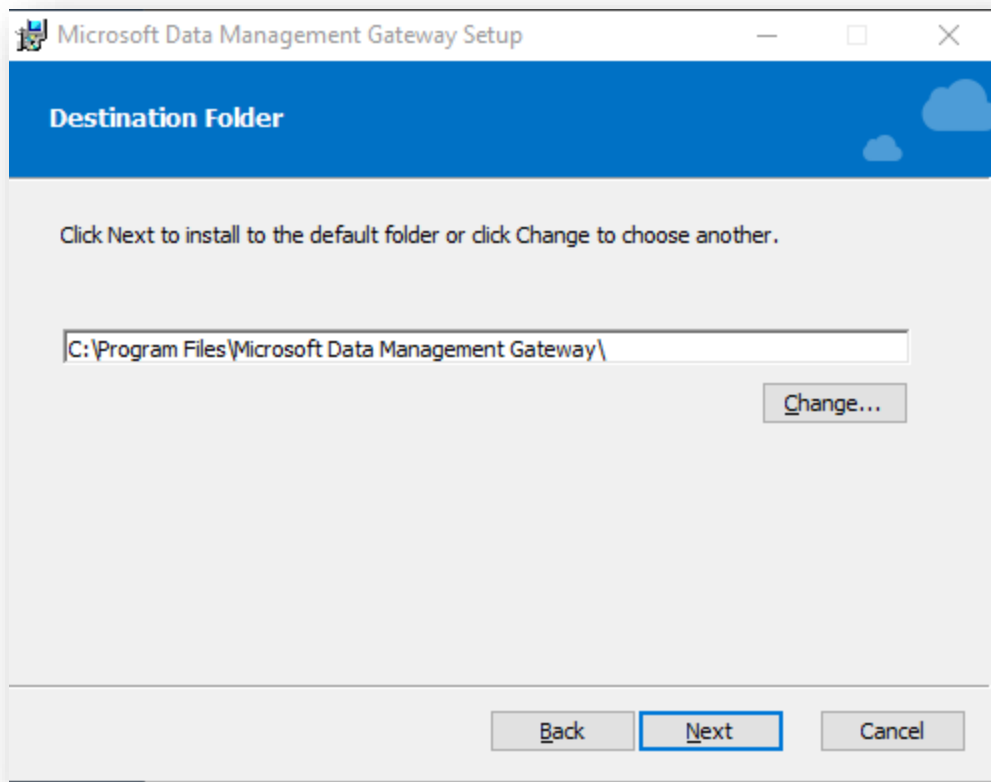
Through this lab, you have learnt how to copy over your files one-time from on-premises location using ADF over to Azure Data Lake Store. You can now expand this learning to try out different types of sources, do copying in a recurring fashion, do different transformations etc. by using Azure Data Factory as the tool of choice.

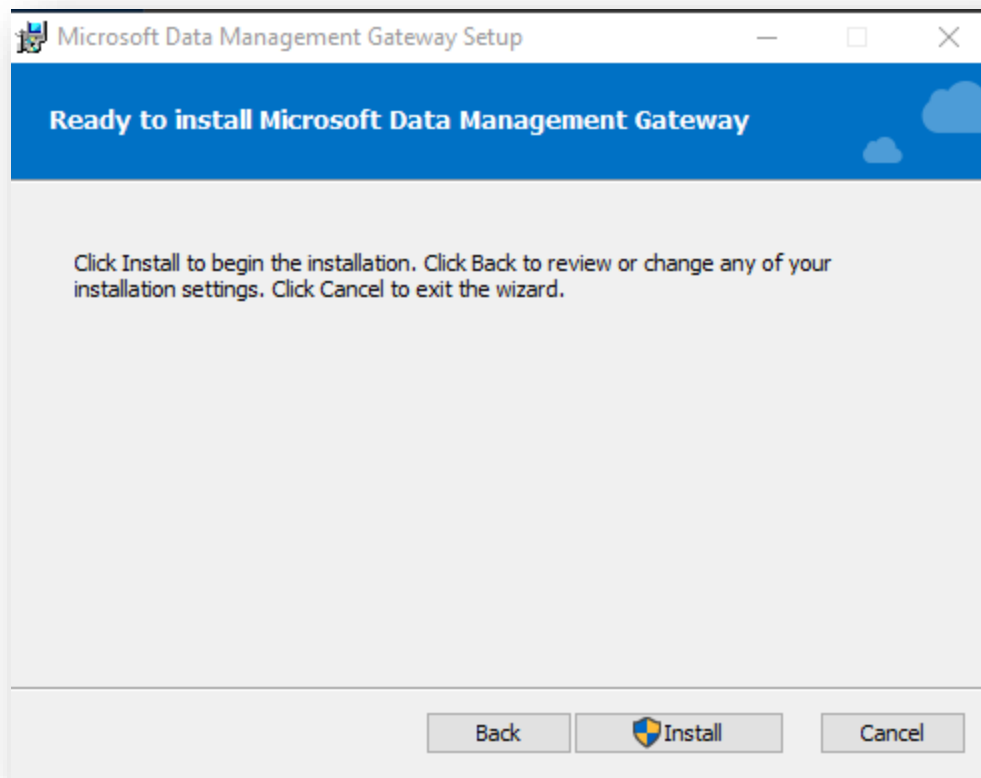
Appendix

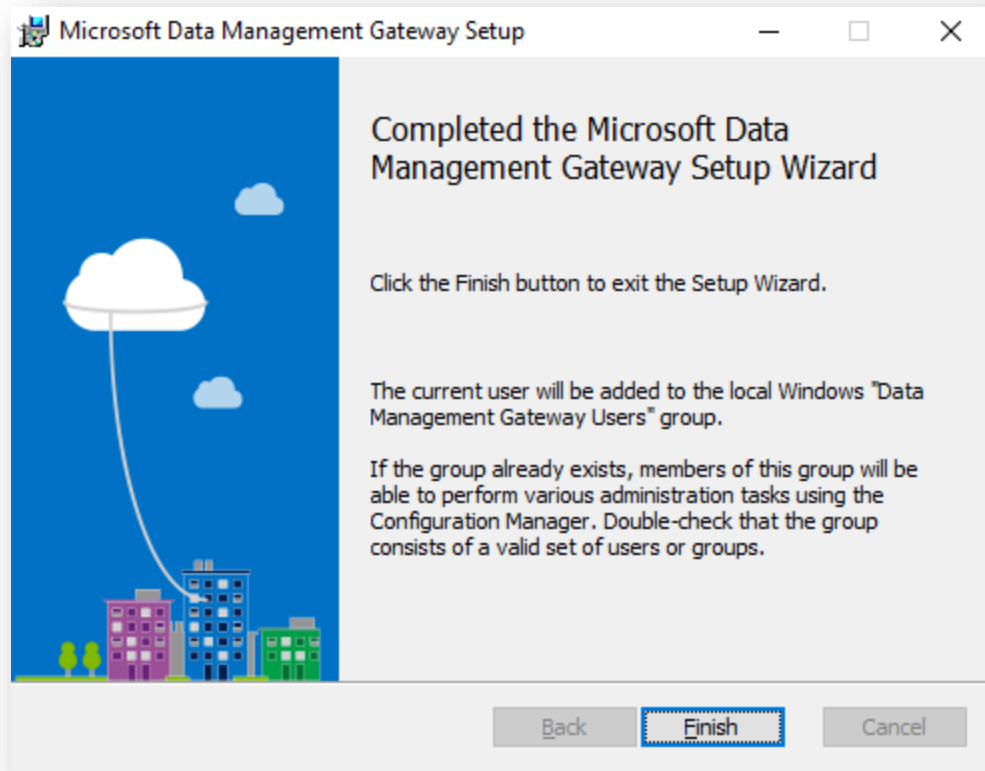
Installing ADF Data Management Gateway manually

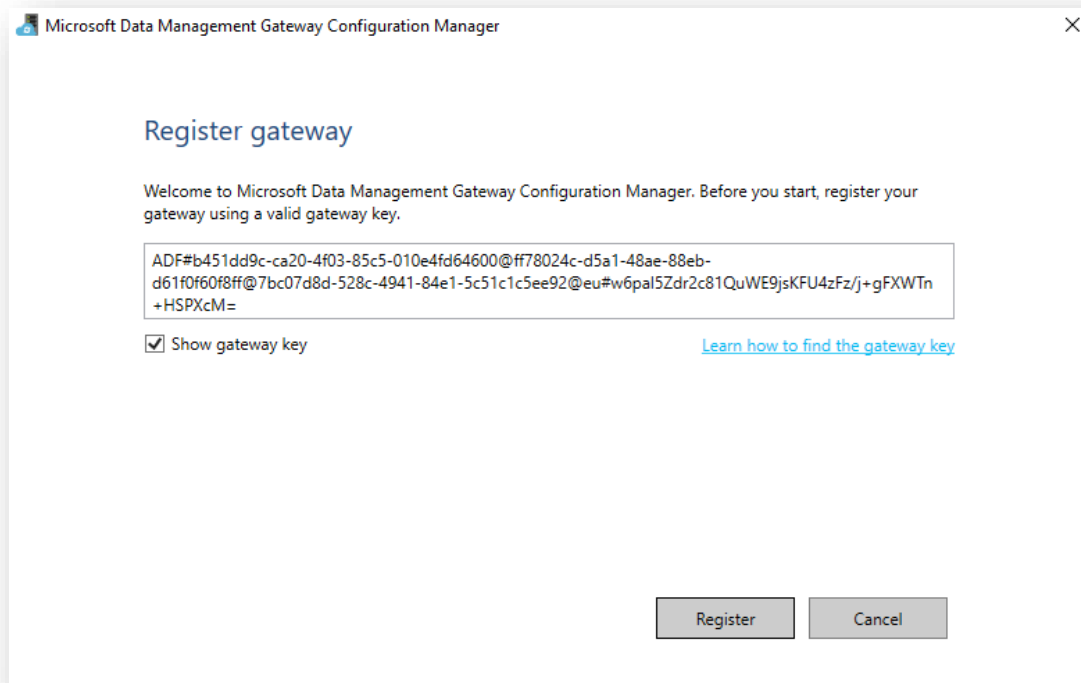












PowerShell Script to Create the ADLSTutorial Share

```
# NOTE: Run as Admin
$dir = "D:\ADLSTutorial"
$user = "redmond\saveenr"

New-Item $dir -type Directory
New-SMBShare -Name "ADLSTutorial" -Path $dir `
-FullAccess $user
```

ClickOnce Installation

NOTE: Installing via Browser.

Installing DMG from the portal uses the Microsoft ClickOnce feature.

- Internet Explorer -> Already supports ClickOnce
- Edge -> Already supports Click Once
- Chrome -> Install from the [Chrome web store](#), search with "ClickOnce" keyword, choose one of the ClickOnce extensions, and install it.
- Firefox -> Click **Open Menu** button on the toolbar (**three horizontal lines** in the top-right corner), click **Add-ons**, search with "ClickOnce" keyword, choose one of the ClickOnce extensions, and install it.

Changelog

- 2016/12/26 Cleaned up formatting
- 2016/12/24 Clarified text
- 2016/12/23 Updated for TR24