

LAB - Copy data from Azure Blob storage to a database in Azure SQL Database by using Azure Data Factory

In this tutorial, you create a data factory by using the Azure Data Factory user interface (UI). The pipeline in this data factory copies data from Azure Blob storage to a database in Azure SQL Database. The configuration pattern in this tutorial applies to copying from a file-based data store to a relational data store. For a list of data stores supported as sources and sinks, see the [supported data stores](#) table.

In this tutorial, you perform the following steps:

- Create a data factory.
- Create a pipeline with a copy activity.
- Test run the pipeline.
- Trigger the pipeline manually.
- Trigger the pipeline on a schedule.
- Monitor the pipeline and activity runs.

Prerequisites

- **Azure subscription.** If you don't have an Azure subscription, create a [free Azure account](#) before you begin.
- **Azure storage account.** You use Blob storage as a *source* data store. If you don't have a storage account, see [Create an Azure storage account](#) for steps to create one.
- **Azure SQL Database.** You use the database as a *sink* data store. If you don't have a database in Azure SQL Database, see the [Create a database in Azure SQL Database](#) for steps to create one.

Create a blob and a SQL table

Now, prepare your Blob storage and SQL database for the tutorial by performing the following steps.

Create a source blob

1. Launch Notepad. Copy the following text, and save it as an **emp.txt** file on your disk:

```
FirstName,LastName
John,Doe
Jane,Doe
```

2. Create a container named **adftutorial** in your Blob storage. Create a folder named **input** in this container. Then, upload the **emp.txt** file to the **input** folder. Use the Azure portal or tools such as [Azure Storage Explorer](#) to do these tasks.

Create a sink SQL table

1. Use the following SQL script to create the **dbo.emp** table in your database:

```
CREATE TABLE dbo.emp
(
    ID int IDENTITY(1,1) NOT NULL,
    FirstName varchar(50),
    LastName varchar(50)
)
GO

CREATE CLUSTERED INDEX IX_emp_ID ON dbo.emp (ID);
```

2. Allow Azure services to access SQL Server. Ensure that **Allow access to Azure services** is turned **ON** for your SQL Server so that Data Factory can write data to your SQL Server. To verify and turn on this setting, go to logical SQL server > Overview > Set server firewall> set the **Allow access to Azure services** option to **ON**.

Create a data factory

In this step, you create a data factory and start the Data Factory UI to create a pipeline in the data factory.

1. Open **Microsoft Edge** or **Google Chrome**. Currently, Data Factory UI is supported only in Microsoft Edge and Google Chrome web browsers.
2. On the left menu, select **Create a resource > Integration > Data Factory**.
3. On the **Create Data Factory** page, under **Basics** tab, select the Azure **Subscription** in which you want to create the data factory.
4. For **Resource Group**, take one of the following steps:
 - a. Select an existing resource group from the drop-down list.
 - b. Select **Create new**, and enter the name of a new resource group.

To learn about resource groups, see [Use resource groups to manage your Azure resources](#).

5. Under **Region**, select a location for the data factory. Only locations that are supported are displayed in the drop-down list. The data stores (for example, Azure Storage and SQL Database) and computes (for example, Azure HDInsight) used by the data factory can be in other regions.
6. Under **Name**, enter **ADFTutorialDataFactory**.

The name of the Azure data factory must be *globally unique*. If you receive an error message about the name value, enter a different name for the data factory. (for example, yournameADFTutorialDataFactory). For naming rules for Data Factory artifacts, see [Data Factory naming rules](#).

Create Data Factory ...

Basics Git configuration 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

Region * ⓘ

Name * ⓘ

✖ The Data Factory name is already taken. Choose a different name.

Version * ⓘ

7. Under **Version**, select **V2**.
8. Select **Git configuration** tab on the top, and select the **Configure Git later** check box.
9. Select **Review + create**, and select **Create** after the validation is passed.
10. After the creation is finished, you see the notice in Notifications center. Select **Go to resource** to navigate to the Data factory page.
11. Select **Open** on the **Open Azure Data Factory Studio** tile to launch the Azure Data Factory UI in a separate tab.

Create a pipeline

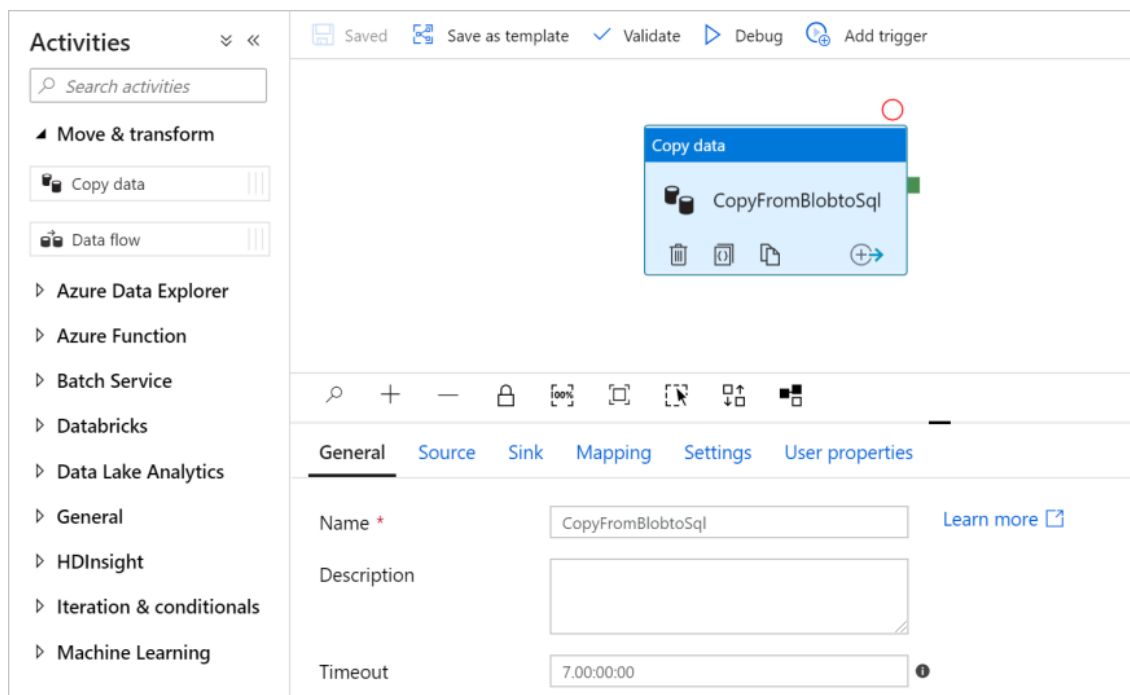
In this step, you create a pipeline with a copy activity in the data factory. The copy activity copies data from Blob storage to SQL Database.

In this tutorial, you start with creating the pipeline. Then you create linked services and datasets when you need them to configure the pipeline.

1. On the home page, select **Orchestrate**.



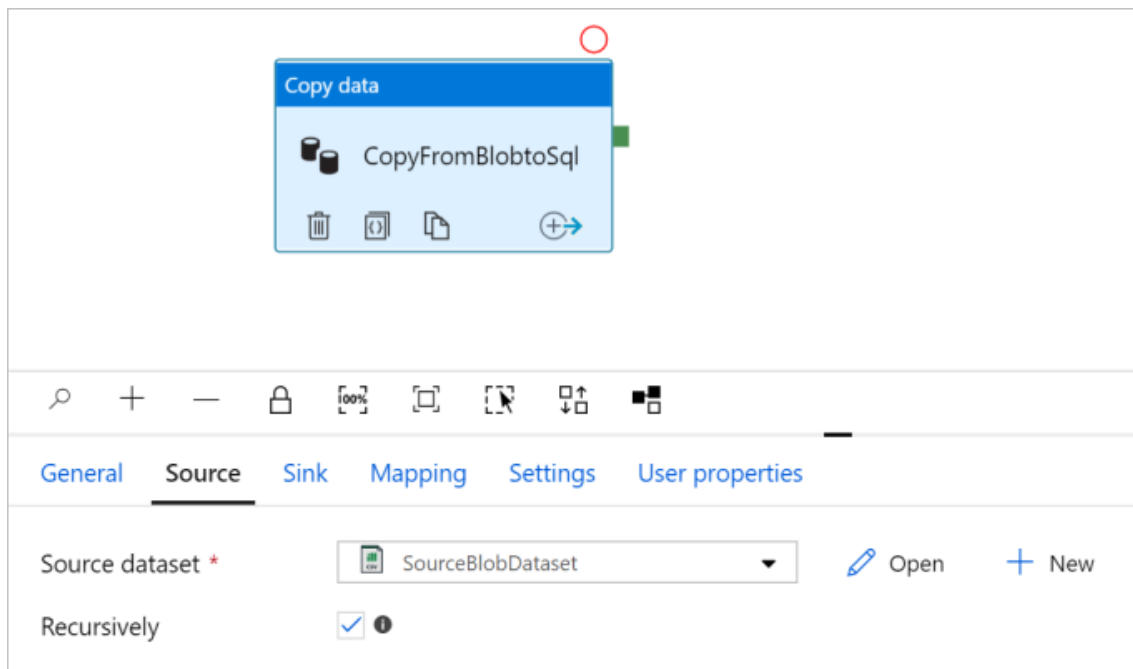
2. In the General panel under **Properties**, specify **CopyPipeline** for **Name**. Then collapse the panel by clicking the Properties icon in the top-right corner.
3. In the **Activities** tool box, expand the **Move and Transform** category, and drag and drop the **Copy Data** activity from the tool box to the pipeline designer surface. Specify **CopyFromBlobToSql** for **Name**.



Configure source

In this tutorial, you use *Account key* as the authentication type for your source data store, but you can choose other supported authentication methods: *SAS URI*, *Service Principal* and *Managed Identity* if needed. Refer to corresponding sections in [this article](#) for details. To store secrets for data stores securely, it's also recommended to use an Azure Key Vault. Refer to [this article](#) for detailed illustrations.

1. Go to the **Source** tab. Select **+ New** to create a source dataset.
2. In the **New Dataset** dialog box, select **Azure Blob Storage**, and then select **Continue**. The source data is in Blob storage, so you select **Azure Blob Storage** for the source dataset.
3. In the **Select Format** dialog box, choose the format type of your data, and then select **Continue**.
4. In the **Set Properties** dialog box, enter **SourceBlobDataset** for Name. Select the checkbox for **First row as header**. Under the **Linked service** text box, select **+ New**.
5. In the **New Linked Service (Azure Blob Storage)** dialog box, enter **AzureStorageLinkedService** as name, select your storage account from the **Storage account name** list. Test connection, select **Create** to deploy the linked service.
6. After the linked service is created, it's navigated back to the **Set properties** page. Next to **File path**, select **Browse**.
7. Navigate to the **adftutorial/input** folder, select the **emp.txt** file, and then select **OK**.
8. Select **OK**. It automatically navigates to the pipeline page. In **Source** tab, confirm that **SourceBlobDataset** is selected. To preview data on this page, select **Preview data**.



Configure sink

In this tutorial, you use *SQL authentication* as the authentication type for your sink data store, but you can choose other supported authentication methods: *Service Principal* and *Managed Identity* if needed. Refer to corresponding sections in [this article](#) for details. To store secrets for data stores securely, it's also recommended to use an Azure Key Vault. Refer to [this article](#) for detailed illustrations.

1. Go to the **Sink** tab, and select **+ New** to create a sink dataset.
2. In the **New Dataset** dialog box, input "SQL" in the search box to filter the connectors, select **Azure SQL Database**, and then select **Continue**. In this tutorial, you copy data to a SQL database.
3. In the **Set Properties** dialog box, enter **OutputSqlDataset** for Name. From the **Linked service** dropdown list, select **+ New**. A dataset must be associated with a linked service. The linked service has the connection string that Data Factory uses to connect to SQL Database at runtime. The dataset specifies the container, folder, and the file (optional) to which the data is copied.
4. In the **New Linked Service (Azure SQL Database)** dialog box, take the following steps:
 - a. Under **Name**, enter **AzureSqlDatabaseLinkedService**.
 - b. Under **Server name**, select your SQL Server instance.
 - c. Under **Database name**, select your database.
 - d. Under **User name**, enter the name of the user.
 - e. Under **Password**, enter the password for the user.
 - f. Select **Test connection** to test the connection.
 - g. Select **Create** to deploy the linked service.

New linked service (Azure SQL Database)

Name *

Description

Connect via integration runtime *

☒ Connection string
 ☐ Azure Key Vault

Account selection method
☒ From Azure subscription ☐ Enter manually

Azure subscription

Server name *

Database name *

Authentication type *

User name *

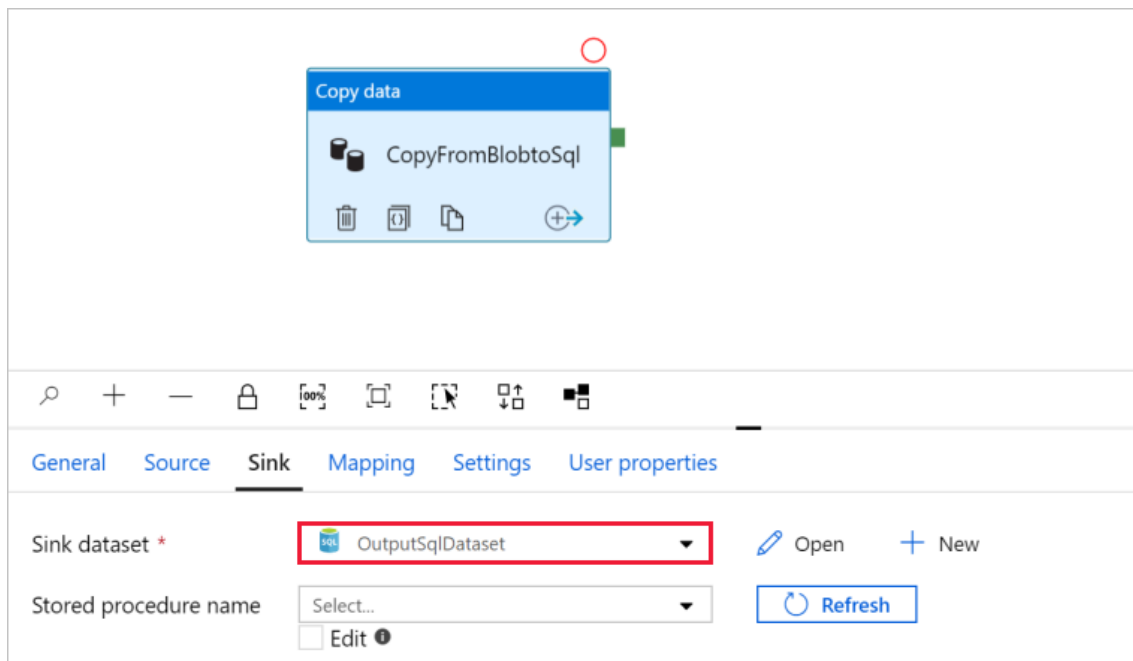
☒ Password
 ☐ Azure Key Vault

Password *

Additional connection properties

☒ Connection successful

- It automatically navigates to the **Set Properties** dialog box. In **Table**, select **[dbo].[emp]**. Then select **OK**.
- Go to the tab with the pipeline, and in **Sink Dataset**, confirm that **OutputSqlDataset** is selected.



You can optionally map the schema of the source to corresponding schema of destination by following [Schema mapping in copy activity](#).

Validate the pipeline

To validate the pipeline, select **Validate** from the tool bar.

You can see the JSON code associated with the pipeline by clicking **Code** on the upper right.

Debug and publish the pipeline

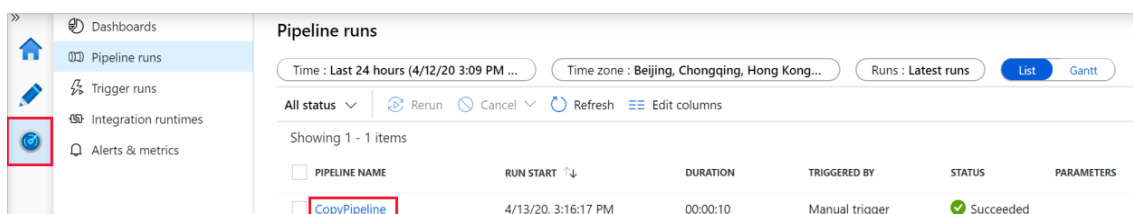
You can debug a pipeline before you publish artifacts (linked services, datasets, and pipeline) to Data Factory or your own Azure Repos Git repository.

1. To debug the pipeline, select **Debug** on the toolbar. You see the status of the pipeline run in the **Output** tab at the bottom of the window.
2. Once the pipeline can run successfully, in the top toolbar, select **Publish all**. This action publishes entities (datasets, and pipelines) you created to Data Factory.
3. Wait until you see the **Successfully published** message. To see notification messages, click the **Show Notifications** on the top-right (bell button).

Trigger the pipeline manually

In this step, you manually trigger the pipeline you published in the previous step.

1. Select **Trigger** on the toolbar, and then select **Trigger Now**. On the **Pipeline Run** page, select **OK**.
2. Go to the **Monitor** tab on the left. You see a pipeline run that is triggered by a manual trigger. You can use links under the **PIPELINE NAME** column to view activity details and to rerun the pipeline.



- To see activity runs associated with the pipeline run, select the **CopyPipeline** link under the **PIPELINE NAME** column. In this example, there's only one activity, so you see only one entry in the list. For details about the copy operation, select the **Details** link (eyeglasses icon) under the **ACTIVITY NAME** column. Select **All pipeline runs** at the top to go back to the Pipeline Runs view. To refresh the view, select **Refresh**.

The screenshot shows the 'CopyPipeline' activity runs view. At the top, there's a breadcrumb 'All pipeline runs > CopyPipeline - Activity runs'. Below this, the 'CopyPipeline' title is followed by 'List' and 'Gantt' tabs. A toolbar contains 'Rerun', 'Rerun from activity', 'Rerun from failed activity', and 'Refresh' buttons. A 'Copy data' window is open, showing a green checkmark and the activity 'CopyFromBlobtoSql'. Below this, the 'Activity runs' section shows a table with one row: 'CopyFromBlo...' with status 'Succeeded'.

ACTIVITY NAME	ACTIVITY TYPE	RUN START	DURATION	STATUS	INTEGRATION RUNTIME
CopyFromBlo...	Copy	4/13/20, 3:16:20 PM	00:00:07	Succeeded	DefaultIntegrationRuntime (West US 2)

- Verify that two more rows are added to the **emp** table in the database.

Trigger the pipeline on a schedule

In this schedule, you create a schedule trigger for the pipeline. The trigger runs the pipeline on the specified schedule, such as hourly or daily. Here you set the trigger to run every minute until the specified end datetime.

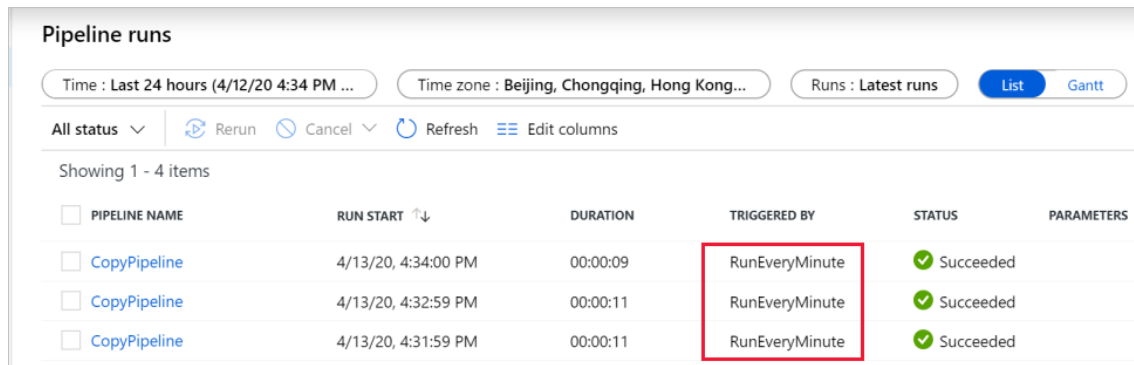
- Go to the **Author** tab on the left above the monitor tab.
- Go to your pipeline, click **Trigger** on the tool bar, and select **New/Edit**.
- In the **Add triggers** dialog box, select **+ New** for **Choose trigger** area.
- In the **New Trigger** window, take the following steps:
 - Under **Name**, enter **RunEveryMinute**.
 - Update the **Start date** for your trigger. If the date is before current datetime, the trigger will start to take effect once the change is published.
 - Under **Time zone**, select the drop-down list.
 - Set the **Recurrence** to **Every 1 Minute(s)**.
 - Select the checkbox for **Specify an end date**, and update the **End On** part to be a few minutes past the current datetime. The trigger is activated only after you publish the changes. If you set it to only a couple of minutes apart, and you don't publish it by then, you don't see a trigger run.
 - For **Activated** option, select **Yes**.
 - Select **OK**.

Important

A cost is associated with each pipeline run, so set the end date appropriately.

- On the **Edit trigger** page, review the warning, and then select **Save**. The pipeline in this example doesn't take any parameters.

6. Click **Publish all** to publish the change.
7. Go to the **Monitor** tab on the left to see the triggered pipeline runs.



Pipeline runs

Time : Last 24 hours (4/12/20 4:34 PM ...) Time zone : Beijing, Chongqing, Hong Kong... Runs : Latest runs [List](#) [Gantt](#)

All status ☐ Rerun ☐ Cancel ☐ Refresh ☐ Edit columns

Showing 1 - 4 items

<input type="checkbox"/> PIPELINE NAME	RUN START <input type="checkbox"/>	DURATION	TRIGGERED BY	STATUS	PARAMETERS
<input type="checkbox"/> CopyPipeline	4/13/20, 4:34:00 PM	00:00:09	RunEveryMinute	✓ Succeeded	
<input type="checkbox"/> CopyPipeline	4/13/20, 4:32:59 PM	00:00:11	RunEveryMinute	✓ Succeeded	
<input type="checkbox"/> CopyPipeline	4/13/20, 4:31:59 PM	00:00:11	RunEveryMinute	✓ Succeeded	

8. To switch from the **Pipeline Runs** view to the **Trigger Runs** view, select **Trigger Runs** on the left side of the window.
9. You see the trigger runs in a list.
10. Verify that two rows per minute (for each pipeline run) are inserted into the **emp** table until the specified end time.

Conclusion

The pipeline in this sample copies data from one location to another location in Blob storage. You learned how to:

- Create a data factory.
- Create a pipeline with a copy activity.
- Test run the pipeline.
- Trigger the pipeline manually.
- Trigger the pipeline on a schedule.
- Monitor the pipeline and activity runs.