

# Orchestrating Big Data Solutions with Azure Data Factory

## Lab 2 – Creating a Pipeline

### Overview

In this lab, you will use Azure Data Factory to implement a simple data pipeline that copies data from a file in Azure Blob Storage to a table in Azure SQL Database.

### What You'll Need

To complete the labs, you will need the following:

- A web browser
- A Microsoft account
- A Microsoft Azure subscription
- A Windows, Linux, or Mac OS X computer
- The lab files for this course
- The Azure resources created in the previous lab.

**Important:** Before starting this lab, make sure you have completed [Lab 1](#).

### Exercise 1: Create a Pipeline

You can create custom data flows by defining the linked services, datasets, and pipelines required to transfer and transform your data. In this exercise, you will create a simple pipeline that performs the same copy data task you accomplished with the Copy Data wizard in Lab 1.

#### Delete the copied transaction data

To avoid confusion, you will delete the database records inserted by the wizard in the previous lab, before recreating the pipeline manually.

1. Start your SQL Server client tool of choice, and connect to the **DataDB** database on your Azure SQL Database server (*server.database.windows.net*) using the server admin login credentials you specified when creating the Azure SQL database.

If you are using the cross-platform SQL Server command line interface, open a command line or console, and enter the following (case-sensitive) command, replacing *server* with your Azure SQL Database Server name, *login* with your server admin login name, and *password* with your server login password. Note that on Linux and Mac OS X operating systems, you may need to prefix any

special characters (such as \$) with a \ character (for example, if your password is *Pa\$\$w0rd*, enter *Pa\\$\\$w0rd*):

```
mssql -s server.database.windows.net -u login -p password -d DataDB -e
```

2. Enter the following command:

```
TRUNCATE TABLE dbo.transactions;
```

3. Enter the following command to verify that the table is empty again:

```
SELECT * FROM dbo.transactions;
```

4. Keep the SQL Server client tool open. You will return to it in a later procedure.

## View Existing Linked Services

Azure Data Factory uses *linked services* to access, transform, and store data. In the case of your pipeline, you require linked services for the blob store account where the source data is stored, and the Azure SQL Database containing the table you want to load.

1. In the Microsoft Azure portal, browse to the blade for your data factory, and click the **Author and deploy** tile.
2. In the pane on the left, expand **Linked Services** and note that linked services named **blob-store** and **sql-database** are already defined for the blob store and SQL database – these were created by the **Copy Data** wizard in the previous lab.
3. Click the **blob-store** linked service to view its JSON definition.
4. Click the **sql-database** linked service to view its JSON definition.

## Create Datasets

Datasets define schema for data flows, including data to be extracted from a source and data to be loaded into a destination (or *sink*). You must create datasets to define the text file data source and the database table to be used in your pipeline.

1. In the author and deploy page, in the pane on the left, click **More** and then click **New dataset**, and click **Azure Blob storage**. This creates a new draft JSON document to define the dataset for the source text data.
2. Replace the default JSON with the following code, (which you can copy and paste from **transactions-txt.json** in the folder where you extracted the lab files):

```
{
  "name": "transactions-txt",
  "properties": {
    "structure": [
      {
        "name": "tdate",
        "type": "Datetime"
      },
      {
        "name": "amount",
        "type": "Decimal"
      }
    ]
  },
  "type": "AzureBlob",
```

```

    "linkedServiceName": "blob-store",
    "typeProperties": {
      "folderPath": "adf-data/data/",
      "format": {
        "type": "TextFormat",
        "columnDelimiter": ",",
        "firstRowAsHeader": true
      }
    },
    "external": true,
    "availability": {
      "frequency": "Minute",
      "interval": 15
    }
  }
}

```

This JSON defines a schema for comma-delimited text containing a **tdate** column of **DateTime** values, and an **amount** column of **Decimal** values. The data is available in the **adf-data/data** folder of the blob storage account defined by the **blob-store** linked service, and new data will be available every fifteen minutes.

3. Click **Deploy** to deploy the dataset definition to your Azure Data Factory.
4. In the author and deploy page, in the pane on the left, click **More** and then click **New dataset**, and click **Azure SQL**. This creates a new draft JSON document to define the dataset for your SQL database table destination.
5. Replace the default JSON with the following code (which you can copy and paste from **dbo-transactions.json** in the folder where you extracted the lab files):

```

{
  "name": "dbo-transactions",
  "properties": {
    "type": "AzureSqlTable",
    "linkedServiceName": "sql-database",
    "structure": [
      {
        "name": "tdate",
        "type": "Datetime"
      },
      {
        "name": "amount",
        "type": "Decimal"
      }
    ],
    "typeProperties": {
      "tableName": "dbo.transactions"
    },
    "availability": {
      "frequency": "Minute",
      "interval": 15
    }
  }
}

```

This JSON defines a schema for a table named **dbo.transactions** containing a **tdate** column of **DateTime** values, and an **amount** column of **Decimal** values in the database defined by the **sql-database** linked service. The dataset can be updated every fifteen minutes.

6. Click **Deploy** to deploy the dataset definition to your Azure Data Factory.
7. In the pane on the left, expand the **Datasets** folder and verify that the **transactions-txt** and **dbo-transactions** datasets are listed.

## Create a Pipeline

Now that you have defined the linked services and datasets for your data flow, you can create a pipeline to encapsulate it. A pipeline defines a series of actions that use linked services to operate on datasets. In this case, your pipeline will consist of a single action to copy data from the **blob-store** linked service to the **sql-database** linked service.

1. In the author and deploy page, in the pane on the left, click **More** and then click **New pipeline**. This creates a new draft JSON document to define the pipeline for your data flow.
2. Replace the default JSON with the following code, which you can copy and paste from **copytrans.json** in the folder where you extracted the lab files:

```
{
  "name": "copytrans",
  "properties": {
    "activities": [
      {
        "type": "Copy",
        "typeProperties": {
          "source": {
            "type": "BlobSource",
            "recursive": false
          },
          "sink": {
            "type": "SqlSink",
            "writeBatchSize": 0,
            "writeBatchTimeout": "00:00:00"
          },
          "translator": {
            "type": "TabularTranslator",
            "columnMappings": "tdate:tdate,amount:amount"
          }
        },
        "inputs": [
          {
            "name": "transactions-txt"
          }
        ],
        "outputs": [
          {
            "name": "dbo-transactions"
          }
        ],
        "name": "CopyTransData"
      }
    ]
  },
  "name": "CopyTransData"
}
```

```
    "pipelineMode": "OneTime"
  }
}
```

This JSON defines a pipeline that includes a **Copy** action to copy the **transaction-txt** dataset in the blob store to the **dbo.transactions** table in the SQL database.

3. Click **Deploy** to deploy the pipeline to your Azure Data Factory.

### View Pipeline Status

Now that you have deployed your pipeline, you can use the Azure portal to monitor its status.

1. After the pipeline has been deployed, return to the blade for your Azure Data Factory, and click for the **Pipelines** tile - it may already indicate the **Wizard Copy** pipeline that you created with the **Copy Data** wizard previously – wait for it to also display the **copytrans** pipeline with a **Pipeline State** of **Running**.
2. Click the **copytrans** pipeline, and observe the state of the activities it contains – it may initially contain no activities, as the datasets are validated and the pipeline is prepared.
3. Close the **Pipelines** blade, and wait for around 15 minutes. Then reopen the **Pipelines** blade and the **copytrans** blade to view its state. When the pipeline has been run, its **State** in the **copytrans** blade will be indicated as **Ready**, and in the **Pipelines** blade its **Pipeline State** will be set to **Completed**.

### Verify that the Data Has Been Copied

Now that your pipeline has run, you can verify that the data has been copied to your SQL database.

1. In your SQL Server client tool, enter the following query:

```
SELECT * FROM dbo.transactions;
```

2. Verify that the table now contains transaction data, copied from the text file in your blob store.
3. Close the SQL Server client tool:

To exit *mssql*, enter the following command:

```
.quit
```

**Note:** You will use the resources you created in this lab when performing the next lab, so do not delete them.