

Azure Synapse

Need to have a Data Warehouse – It helps us do analytics on the data that we have. The data is stored in a way where it is made to process high volumes of read requests.

Synapse initially was just a data warehouse but now it is known as Azure Synapse Analytics. Now we can create warehouses with the help of SQL, integrate the data using pipelines and also use data from data lakes.

We can also use Spark for processing and the data and services like Azure monitor and Azure Active Directory with synapse.

Creating an Azure Synapse Workspace

Go to your home screen > Create a Resource > Search for Azure Synapse Analytics > Create.

Enter : Your subscription, your resource group, *unique* workspace name, region, data lake gen2 details (new or old)

of your resources.

Subscription * ⓘ Azure subscription 1
The Synapse and SQL resource providers are now registered with this subscription.

Resource group * ⓘ data-grp
[Create new](#)

Managed resource group ⓘ Enter managed resource group name

Workspace details
Name your workspace, select a location, and choose a primary Data Lake Storage Gen2 file system to serve as the default location for logs and job output.

Workspace name * synapse203 ✓

Region * West US ✓

Select Data Lake Storage Gen2 * ⓘ ☒ From subscription ☐ Manually via URL

Account name * ⓘ (New) synapsedatalakedp203
[Create new](#)

File system name * (New) data
[Create new](#)

☒ Assign myself the Storage Blob Data Contributor role on the Data Lake Storage Gen2 account to interactively query it in the workspace.

On the next screen give the password for your SQL Administrator Credentials, make sure “allow pipelines” checkbox is ticked.

* Basics * **Security** Networking Tags Review + create

Configure security options for your workspace.

SQL administrator credentials
Provide credentials that can be used for administrator access to the workspace's SQL pools. If you don't provide a password, one will be automatically generated. You can change the password later.

SQL Server admin login * ⓘ sqladminuser ✓

SQL Password ⓘ ✓

Confirm password ✓

System assigned managed identity permission
Choose the permissions that you would like to assign to the workspace's system-assigned identity. [Learn more](#)

☒ Allow pipelines (running as workspace's system assigned identity) to access SQL pools. ⓘ

☐ Allow network access to Data Lake Storage Gen2 account. ⓘ

The selected Data Lake Storage Gen2 account does not restrict network access using any network access rules, or you selected a storage account manually via URL under Basics tab. [Learn more](#)

Leave everything else as is and create.

Synapse Compute Options

There are different compute options – Serverless SQL pool and SQL pool.

Serverless SQL Pool

- You can use this option to perform quick adhoc analysis of data
- Can use T-SQL
- Can only create external tables but cannot persist the data
- Charged based on how much you use the service and how much data your process

SQL Pool

- User to build your warehouse
- Can use T-SQL
- Used if you want to persist the data
- Charged based on the data warehousing units (which includes things like compute, memory, etc.)

External tables - Can be defined in the Serverless pool and the dedicated SQL pool. We use external tables when the table data is lying in an external source, but the table definition is lying in Azure synapse. This is useful when you don't want to load the table on to the server itself.

For example if there are tables that exist on an external source and there's data on the sql server, then to perform a join operation between the two, an external table can be used.

There are a few important checks that need to be done in order to access the external data:

- We first need to have authorization to use the external source of data
- We then need to define the format of the external file that we want to use as an external table
- Finally, create the external table

Using External Tables

Open up your synapse dashboard and click on Open Synapse Studio. In the synapse studio you can use SQL commands against your Serverless SQL pool as well as dedicated SQL pool. Can create pipelines to integrate your data and just view your data as well.

Executing a script on Azure Synapse Studio to create External tables

In the left-hand menu, click on develop, click on the plus icon in the develop screen and select SQL Script out of the options given.

Name the script on the left and copy and paste [this](#) SQL script on to the editor.

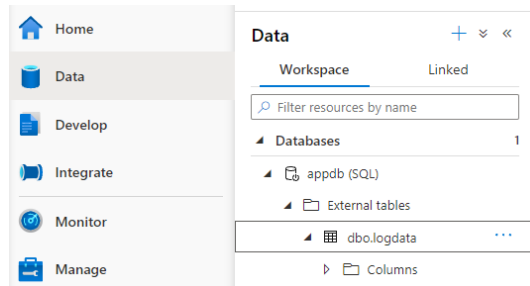
Now, we would be running a series of commands –

- Firstly, we would run the create data base command to create a database in the serverless pool
- Change the database from the top right of the editor where master is selected (refresh if the newly created database is not showing)
- Next, we create a master key that would be used to encrypt the database scope credentials which will allow ourselves to use the file that we would be using in our Data Lake Gen2 account.
- Now to create the scope credentials we need to get the shared access signature like follows and copy the SAS token –

The screenshot shows the 'Generate SAS and connection string' form in the Azure portal. The form is divided into several sections with expandable/collapsible headers. The 'Allowed services' section is expanded, showing checkboxes for Blob, File, Queue, and Table, with 'Blob' selected. The 'Allowed resource types' section is expanded, showing checkboxes for Service, Container, and Object, with 'Service' and 'Object' selected. The 'Allowed permissions' section is expanded, showing checkboxes for Read, Write, Delete, List, Add, Create, Update, Process, and Immutable storage, with 'List' selected. The 'Blob versioning permissions' section is expanded, showing a checkbox for 'enable deletion of versions', which is unchecked. The 'Start and expiry datetime' section is expanded, showing a table with columns for 'Start', 'Expiry', and 'Time zone'. The 'Allowed IP addresses' section is expanded, showing a text input field with the placeholder 'For example: 192.168.1.1 or 192.168.1.0/24'. The 'Allowed protocols' section is expanded, showing radio buttons for 'HTTPS only', 'HTTPS and HTTP', and 'HTTP', with 'HTTPS only' selected. The 'Preferred routing tier' section is expanded, showing radio buttons for 'Basic (default)', 'Microsoft network routing', and 'Internet routing', with 'Basic (default)' selected. The 'Signing key' section is expanded, showing a text input field with the placeholder 'key1'. At the bottom of the form, there is a blue button labeled 'Generate SAS and connection string'.

- Paste it in the SECRET variable and remove the '?' from the front. Run the command for creating scoped credential.
- Next, define the location of your data by giving the location of your file in the LOCATION variable like this : https://<datalake_name>.dfs.core.windows.net/<your_container_name>
- Run the create external data source command with the new location that you entered
- Next, we are giving the format of your file. We name the file format as TextFileFormat, and start reading from the second row as the first row is headers
- Next, we create the external table. We give the column name and the types that we want as all the data would be coming in the string format. Remove all the NULLs from the command and execute the script.
- Now you can run the select * command to see all your data
- To save this script for future use you can click on Publish all button at the top of the editor and it would be available in your Synapse Studio

Now if we click on Data in Synapse Studio, we would be able to see our external table.



Creating a dedicated SQL pool

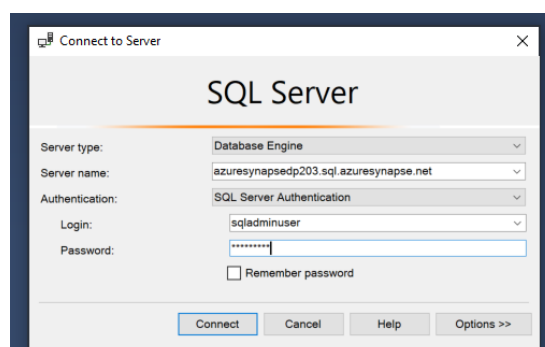
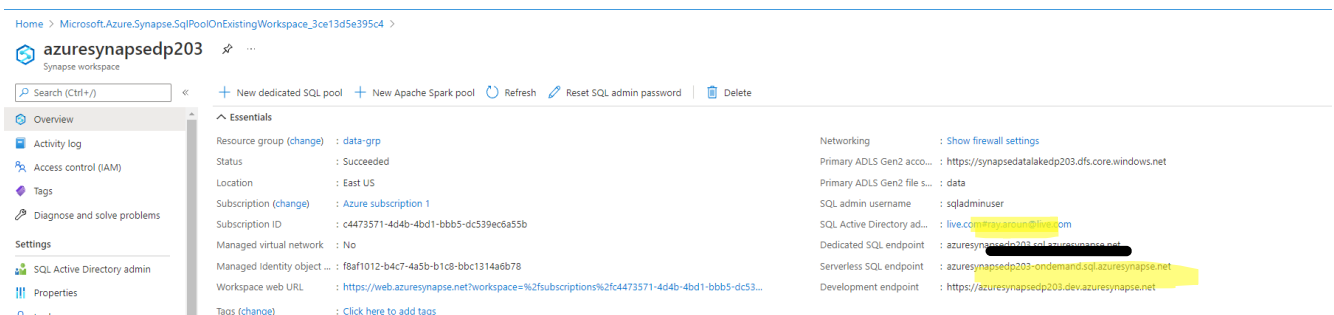
One of the main differences between the Serverless SQL pools and a Dedicated SQL pool is the ability to persist your data.

Go to Synapse workspace, from the left-hand menu select SQL pools. Click New and give a name, choose the performance level (for learning purposes, the lowest level is fine). Leave everything else as is and create the pool.

Once created this pool would also show up on synapse studio under databases.

Creating an external table in the Dedicated SQL pool using the Microsoft SQL Server Management Studio

Copy the Dedicated SQL Endpoint from your Synapse Workspace and paste it in the dialog box for a new connection in SSMS. For the login, copy the SQL admin username on the Synapse Workspace as well and give the password that you gave at the creation of the Synapse Workspace for the SQL admin.



Upon logging in, it would show your dedicated SQL pool. Right click on it and new query. Copy [this](#) script onto it.

We again go through the similar process of creating an external table:

- Select the first command to create a master key and execute it
- Copy your key for the Data lake gen2, and paste it in the SECRET variable. Execute it.
- Replace the name in the location variable with the name of your data lake and here we are specifying the driver (Hadoop) to source the external data. Execute it.
- Execute the external file format command to tell the format of external table
- Create the external table after removing all the NULLs

Now when we try to do a select * on logdata, it gives an error about converting varchar to datetime.

This error occurs because the date that we have in our file is not in the order in which Hadoop could infer it. So we would either have to change the format of the data in the file or we could clean the data or use a different driver to read the data.

After having cleaned the file, we again try to read the file. [Here](#) is the clean version of the file to skip the cleansing.

- Drop the table and create the table again with the new csv file

Upon doing the select * statement we can see the table contents now.

Creating external tables based on a parquet files

Go to your data lake gen2 account, create