# Azure Data Factory - how to provision, access and development guide

Azure Data Factory is a cloud-based data integration service offered by Microsoft Azure, designed to empower organizations with the ability to collect, transform, and move data from diverse sources to various destinations. As a fully managed and scalable service, Azure Data Factory facilitates the creation, scheduling, and management of data pipelines, enabling businesses to gain valuable insights and make informed decisions based on their data.

**Key Features and Uses:**

- **Data Orchestration:** Azure Data Factory allows you to orchestrate and automate the movement and transformation of data across on-premises and cloud environments. This facilitates the creation of end-to-end data workflows.
- **Data Transformation:** With built-in data wrangling capabilities, Azure Data Factory supports the transformation of raw data into a usable format for analytics and reporting. This includes cleaning, enriching, and aggregating data as it moves through the pipeline.
- **Connectivity:** Seamlessly connect to a variety of on-premises and cloud-based data sources, such as Azure SQL Database, Azure Blob Storage, Azure Data Lake Storage, and more. This enables comprehensive integration with diverse data ecosystems.
- **Monitoring and Management:** Azure Data Factory provides robust monitoring and management tools, allowing users to track the performance of data pipelines, identify issues, and optimize data integration processes for enhanced efficiency.
- **Hybrid Data Movement:** Support for hybrid data movement ensures that organizations can integrate and synchronize data across both on-premises and cloud environments, providing flexibility in data management strategies.
- **Data Integration with Machine Learning:** Integrate machine learning models into data pipelines for advanced analytics and predictions. Azure Data Factory enables the seamless incorporation of machine learning capabilities to enhance data-driven decision-making.

**Benefits:**

- **Scalability:** Azure Data Factory offers a serverless and fully managed platform that scales dynamically based on the demands of your data integration processes. This ensures optimal performance and resource utilization without the need for manual intervention.
- **Cost Efficiency:** Pay only for the resources consumed during data integration processes, eliminating the need for upfront investments in infrastructure. The pay-as-you-go model allows organizations to manage costs efficiently and scale as needed.
- **Time Efficiency:** Streamline the development and deployment of data pipelines with a user-friendly visual interface and pre-built connectors. This accelerates the time-to-market for data-driven solutions and reduces the development lifecycle.
- **Data Governance and Security:** Azure Data Factory incorporates robust security measures, including encryption, identity and access management, and compliance certifications. This ensures that sensitive data is protected, and organizations can adhere to regulatory requirements.
- **Business Intelligence and Analytics:** Empower business users with timely and accurate data for analytics and reporting. Azure Data Factory enables organizations to leverage their data assets, driving actionable insights that contribute to informed decision-making and business intelligence.

In order to create your own Azure Data Factory , please follow the steps as mentioned below -

1. Go to Home on the Azure Portal and then search for "Create a Resource". Once you are in that page, click on "Integration" on the left hand pane. A page similar to what is shown below will come up. From here, select "Data Factory".

2. Click on "Create Data Factory", which will take you to a page like this. Select the resource group which was created earlier in ( 📄 Azure B lob Storage - how to provision, access and development guide ). Name your ADF instance following the naming standards which is mentioned on the creation page and select "Region" as "Australia East", since that is the closest Azure region for us. Select "Version" as "V2" and click "Create".

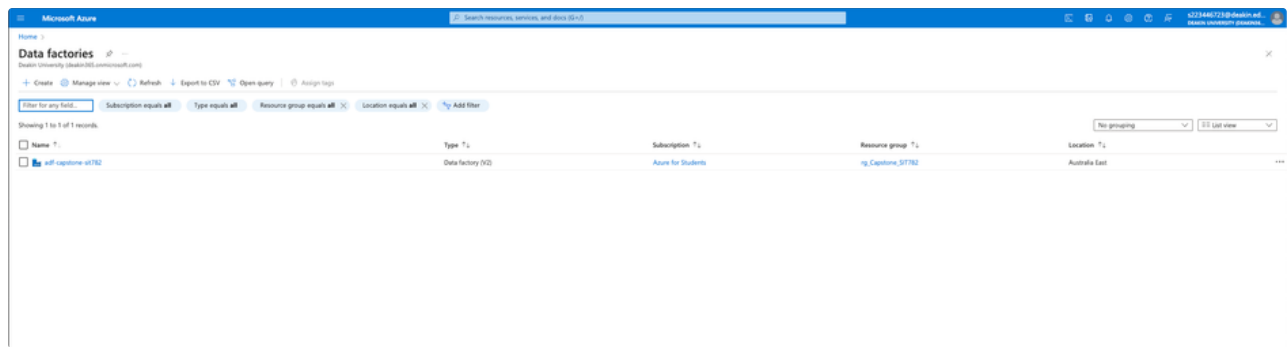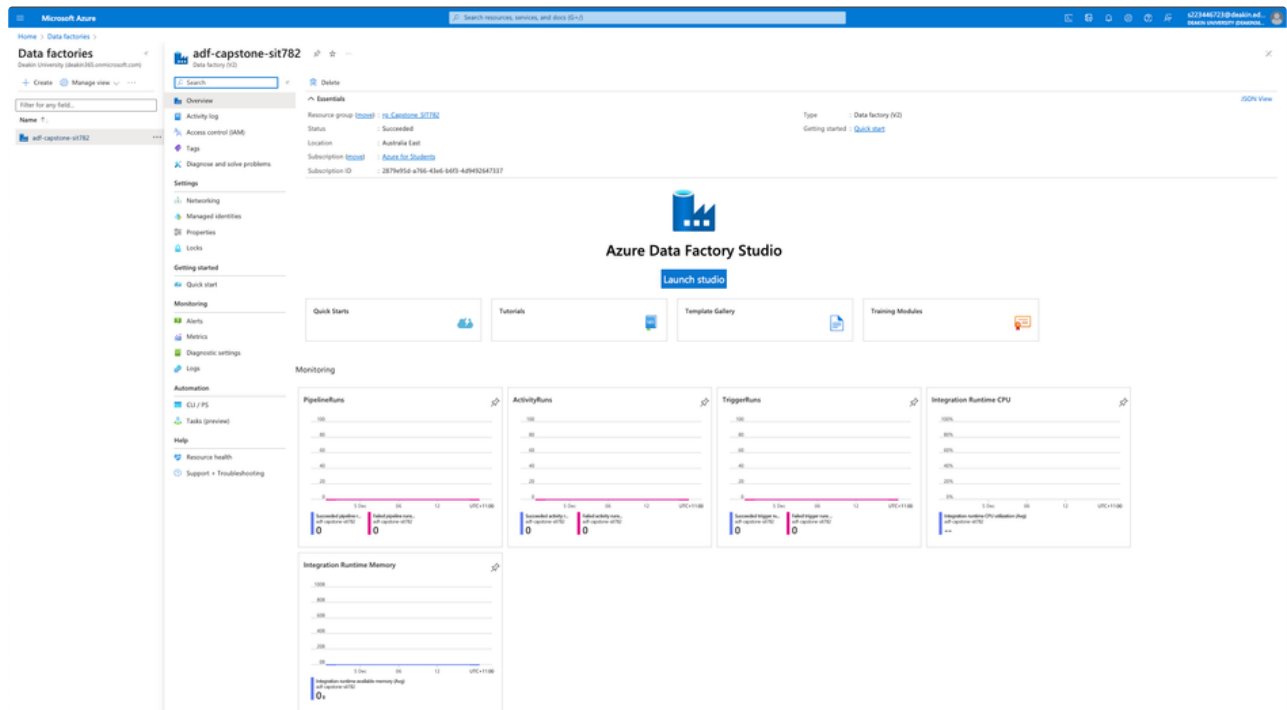3. Once the resource has been created, go to Home page and from there in the menu bar, search for "data factories". Select the Data factories services which comes up, as shown below -
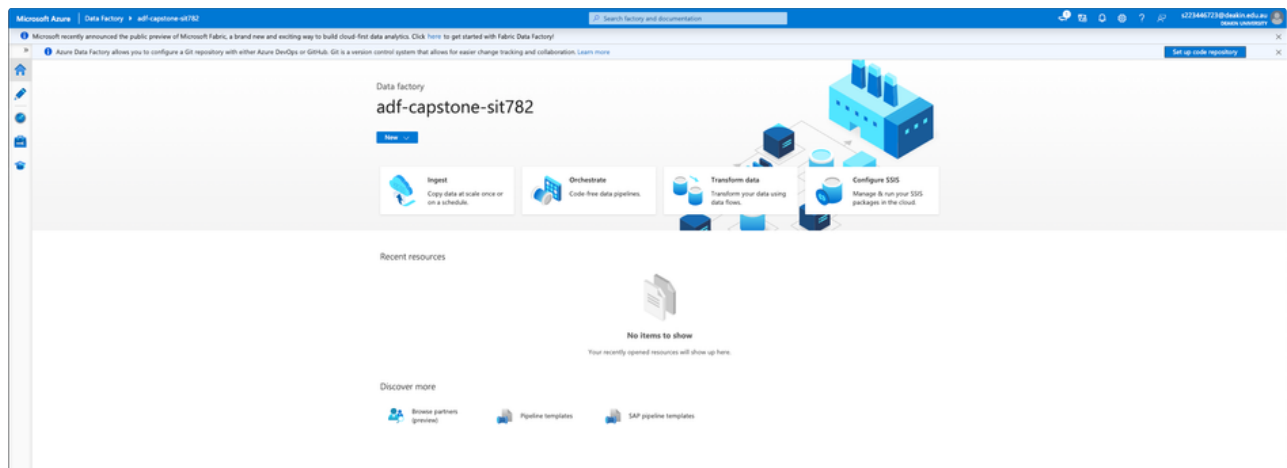


4. Once in the page, you will be able to see all the Data Factories which has been provisioned. Select the ADF which was created in the previous step.
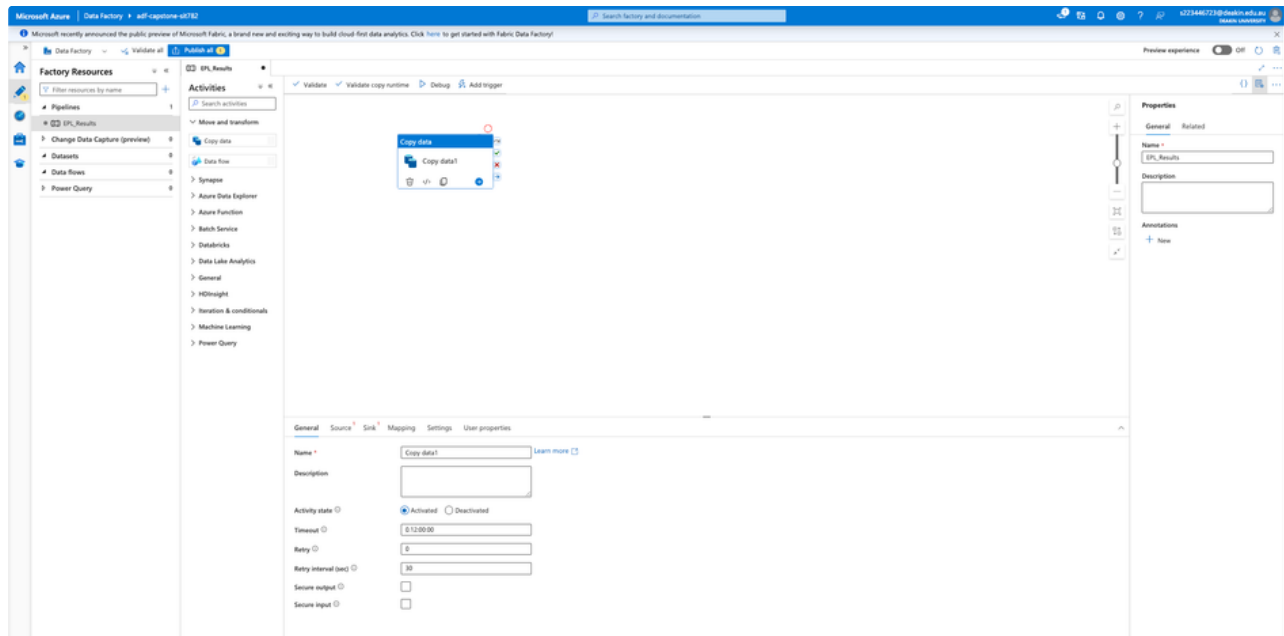
5. You will land in a page similar to what is shown below. This is the home page for your Data Factory. From here you can provision access, check logs, set up networking, monitor alerts and check resource health status. Click on "Launch Studio" next.
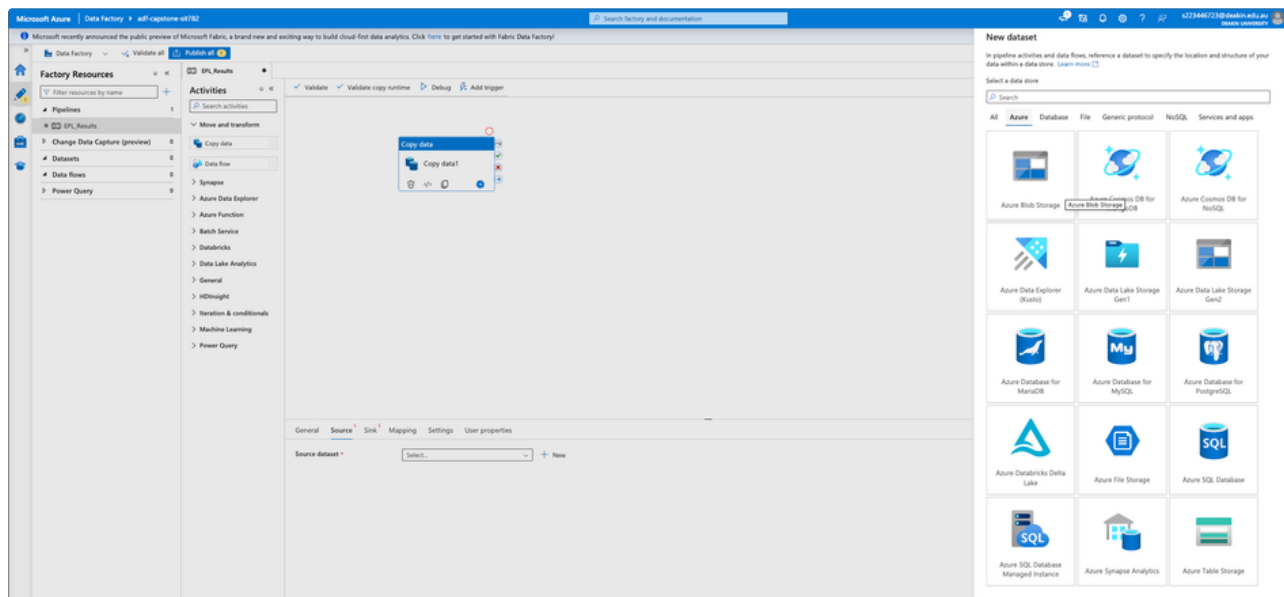


6. Once you click on "Launch Studio", it will open a separate page as shown below. This is your homepage for Azure Data Factory instance which you have created earlier. Note that, for each separate Data Factory which has been provisioned, it will open a separate page. From this page, we will be able to perform all data engineering and data warehouse activities. For this example, we will ingest the EPL results file which was uploaded in the Azure Blob Storage earlier. Click on "Orchestrate" on the main page next.
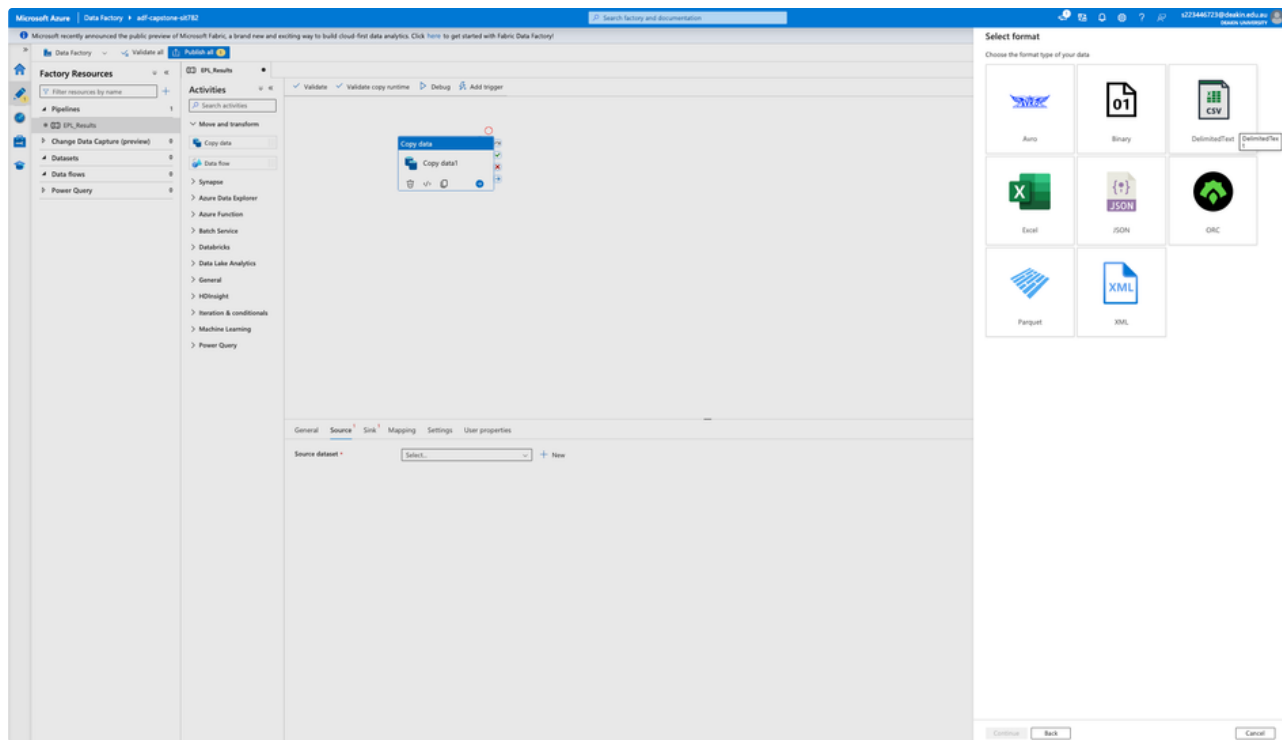
7. Once you click on "Orchestrate", it will open up a blank pipeline. Name the pipeline as "EPL_Results" on the right hand pane. Next click on "Move and transform" to expand it and drag and drop "Copy data" to the main window. You can rename the "Copy data" stage in the "Name" section below. I will keep this as-is for now.
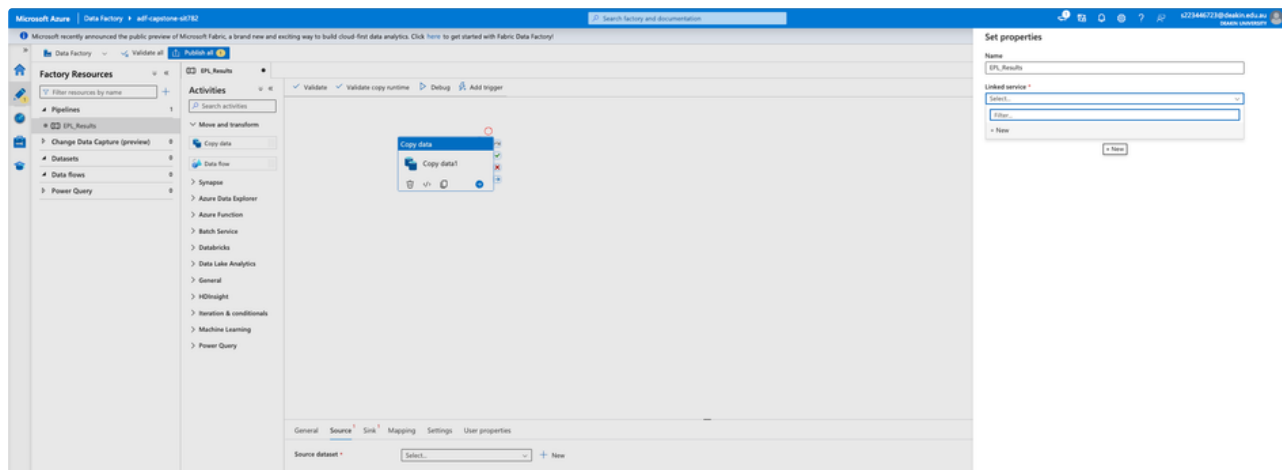


8. Now we will need to define the Source and Sink for this pipeline. Our source is the EPL Results file present in the Blob Storage and our Sink is the table which we created earlier ( 📄 Azure SQL Database - how to provision, access and development guide ). Click on "Source" and then click on "New". This will open the right hand pane as shown below.
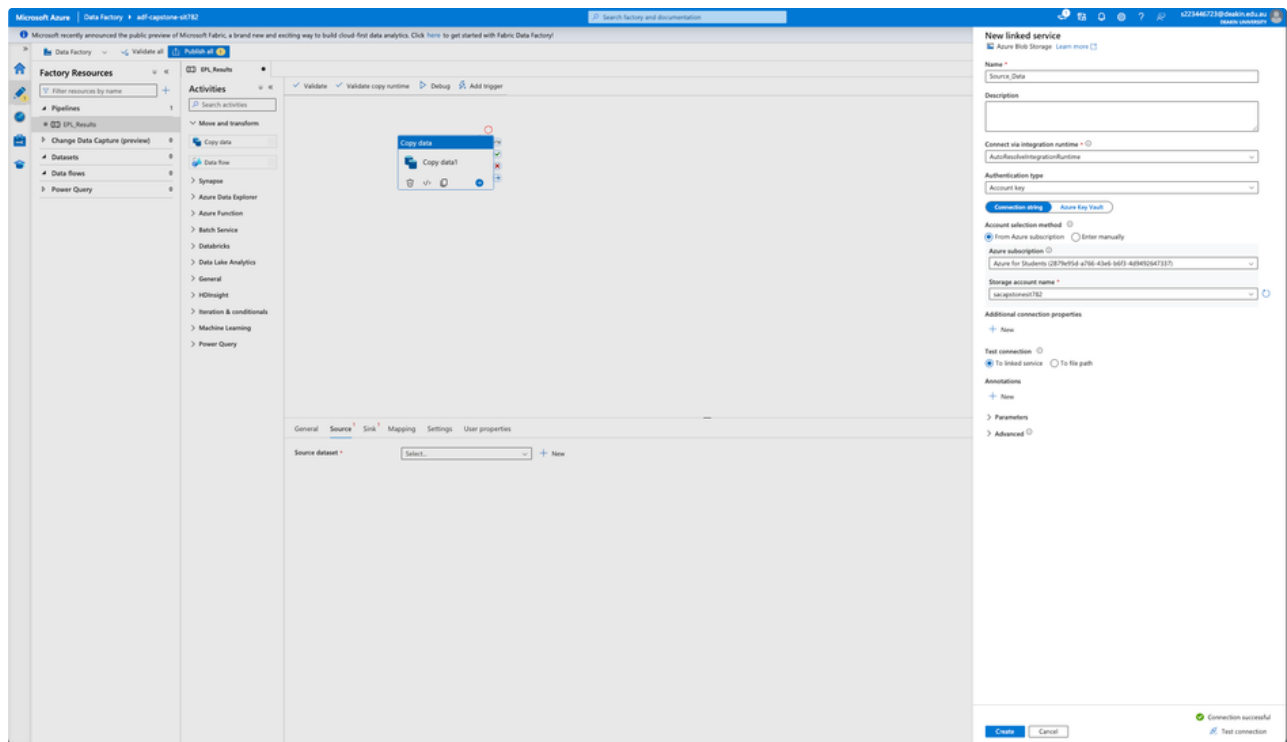


9. Click on "Azure" on the right hand pane and then click on "DelimitedText", since the file uploaded to Blob Storage is a csv file. If this is any other file, choose the appropriate format from the list below.
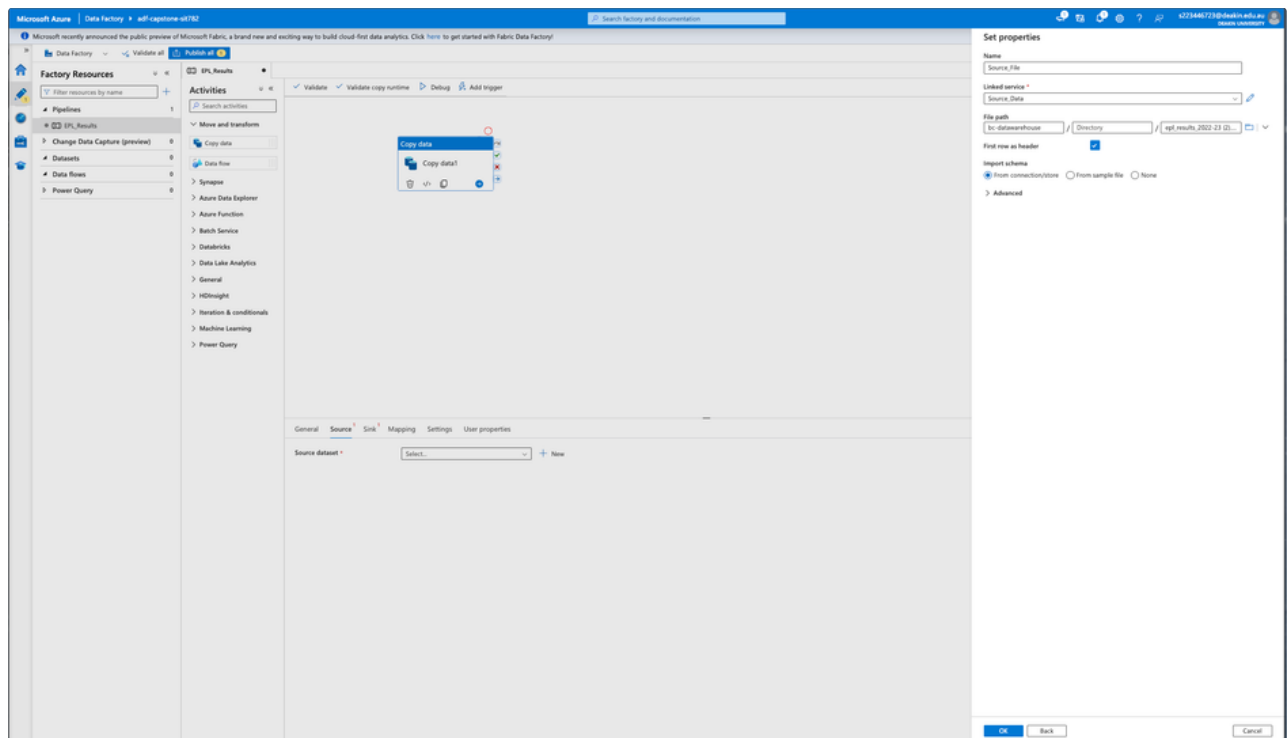
10. Now, name the file. For this POC, I have named it as EPL_Results. Now we need to define the "Linked service" for this Source, which is the connection to Blob Storage. Click on "New" under "Linked service", since this is the first time we will be creating this Linked service.



11. Name the Linked service. Under "Integration runtime", select "AutoResolveIntegrationRuntime". The other option is Self Hosted Integration Runtime, but that would mean more infrastructure setup. Next, select the Azure subscription and the storage account name. This is the same name which was defined when Azure Blob Storage was provisioned. Click on "Test connection" to test your connectivity and once it is successful, click on "Create".

12. Next we need to define the path to the file. Click on the Folder icon in "File path" and select the file. Under "Import schema", select "From connection/store" to import the file schema into the pipeline.



13. Once everything has been correctly set up by following the steps above, you will see the following under "Source" tab. Here you can see that there are multiple options. Since we have a single fixed file, we have chosen "File path in dataset". If this involves multiple files, then choose "List of files".

General   **Source**   Sink[1]   Mapping   Settings   User properties

Source dataset *
📄 Source_File   ⌄    ✏ Open    + New    👓 Preview data    Learn more ☐

File path type
● File path in dataset    ○ Prefix    ○ Wildcard file path    ○ List of files ⓘ

Filter by last modified ⓘ
Start time (UTC)    End time (UTC)

Recursively ⓘ    ✓

Enable partitions discovery ⓘ    ☐

Max concurrent connections ⓘ

Skip line count

Additional columns ⓘ    + New

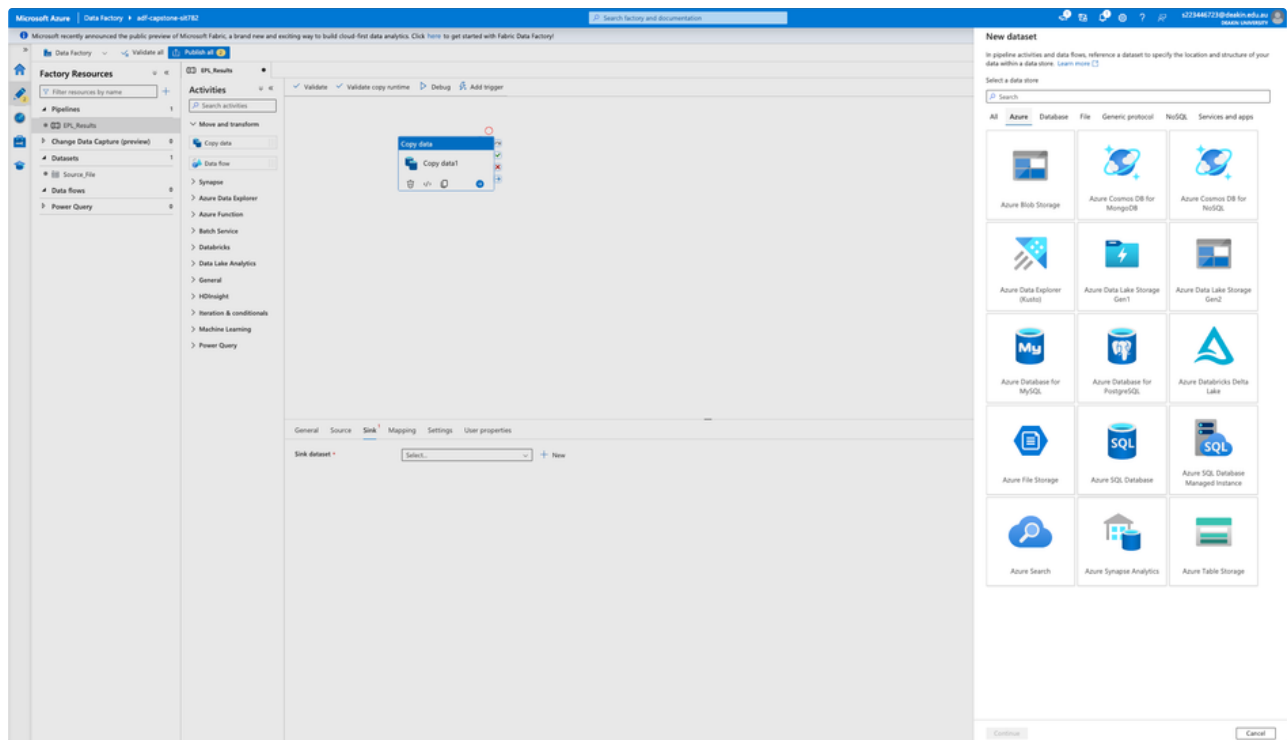14. Next, check that the data in the "Source" is coming up correctly. This can be done by selecting "Preview data" on the "Source tab. Check the data in the window which opens up. Ideally, the data should look good, but if there are unicode character issues, then check the source file for any un-translated characters.
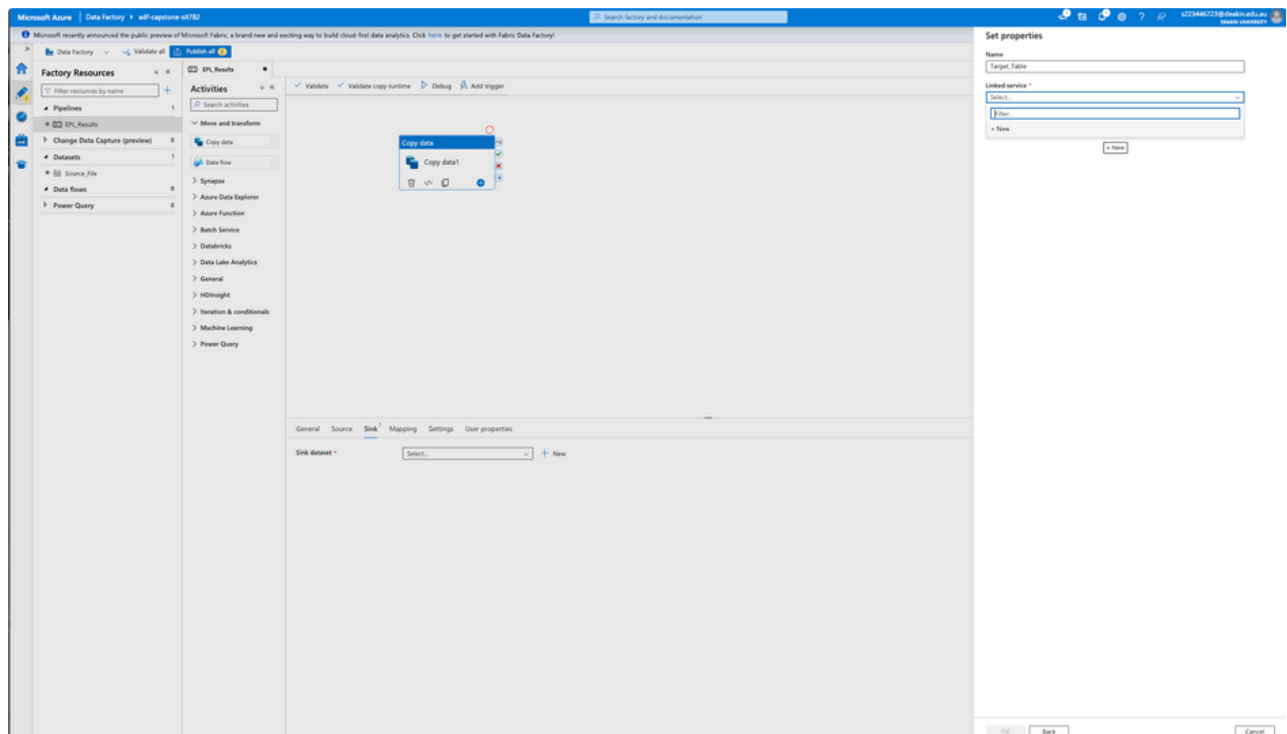


15. Now, we will define the "Sink". Click to "Sink" tab and then click on "New". This will open up a familiar page similar to the below screenshot -
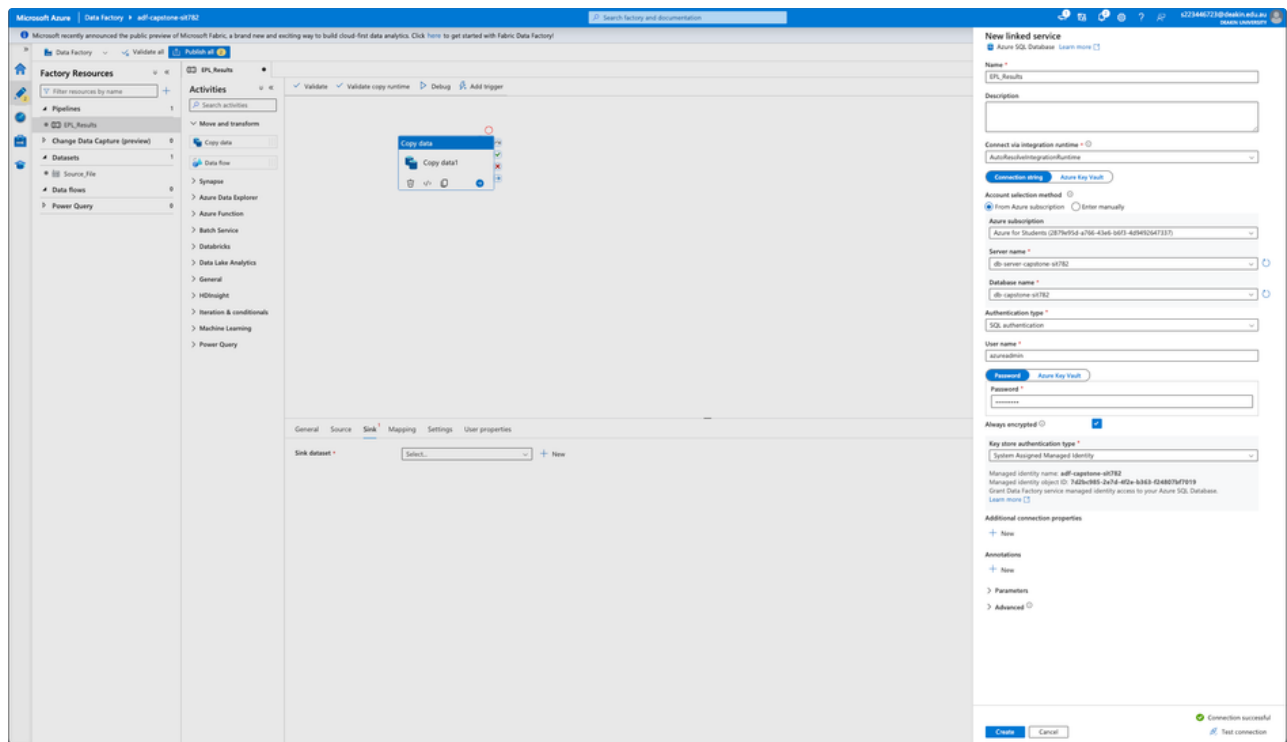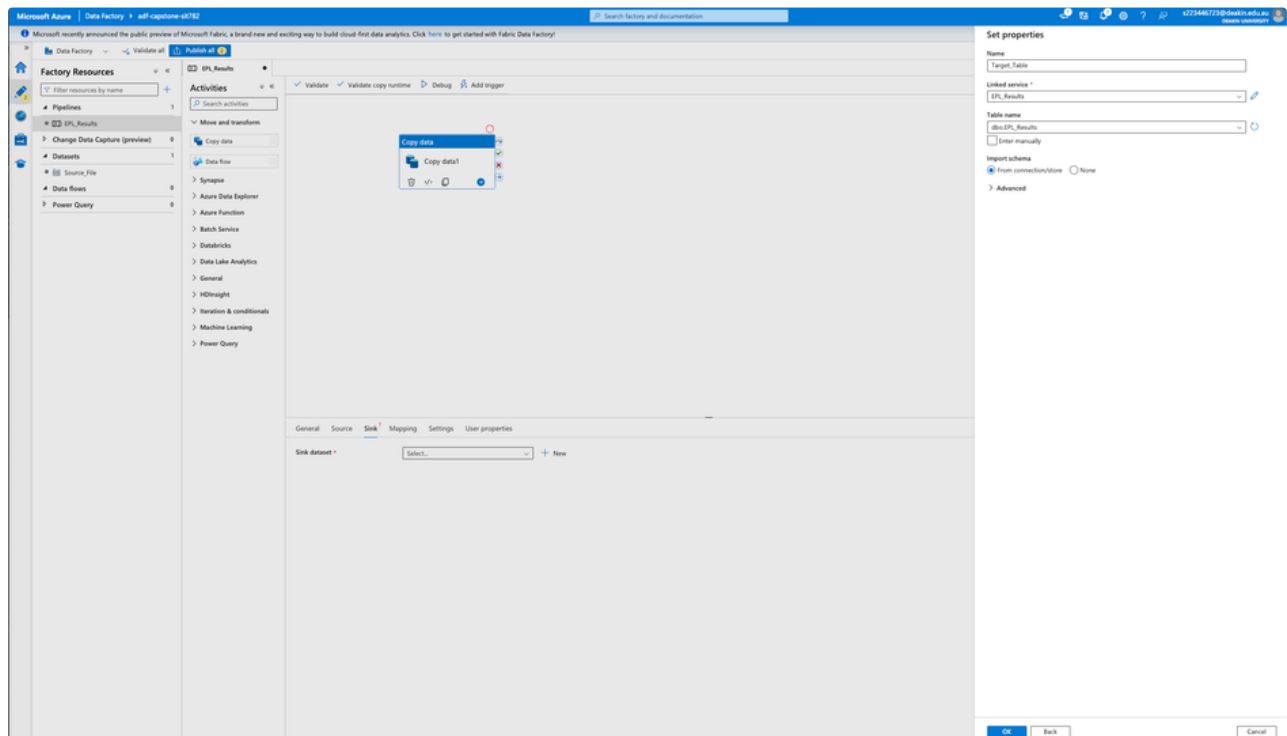
16. Click on "Azure" and then select "Azure SQL Database". We will next configure the connection to the SQL database which we created earlier ( 📄 Azure SQL Database - how to provision, access and development guide ). Name the Sink/target as appropriate and select "New" under "Linked services". I will set up a new Linked service for the SQL Database, since this is the first time I am creating the Linked service. If the Linked service already exists, then it will appear in the drop-down list.
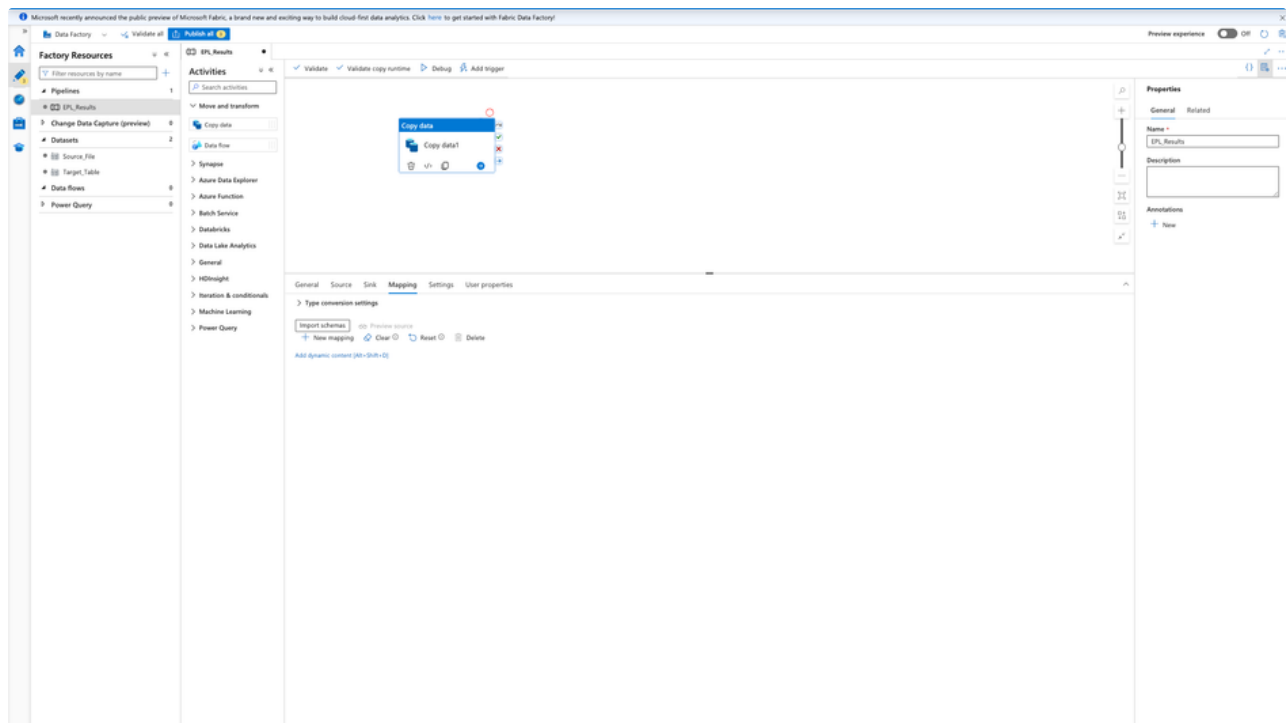


17. Name the linked service next and select the subscription from the drop-down list. Next, select the database server name and the database. Select the "Authentication type" as "SQL authentication" and fill in the details of the user id which was created when the database was provisioned. Input the credentials for the user id and select the check box next to "Always encrypted". This will ensure that the data is always encrypted. Check the connection created by clicking on "Test connection" below. Once it is successful, click on "Create".
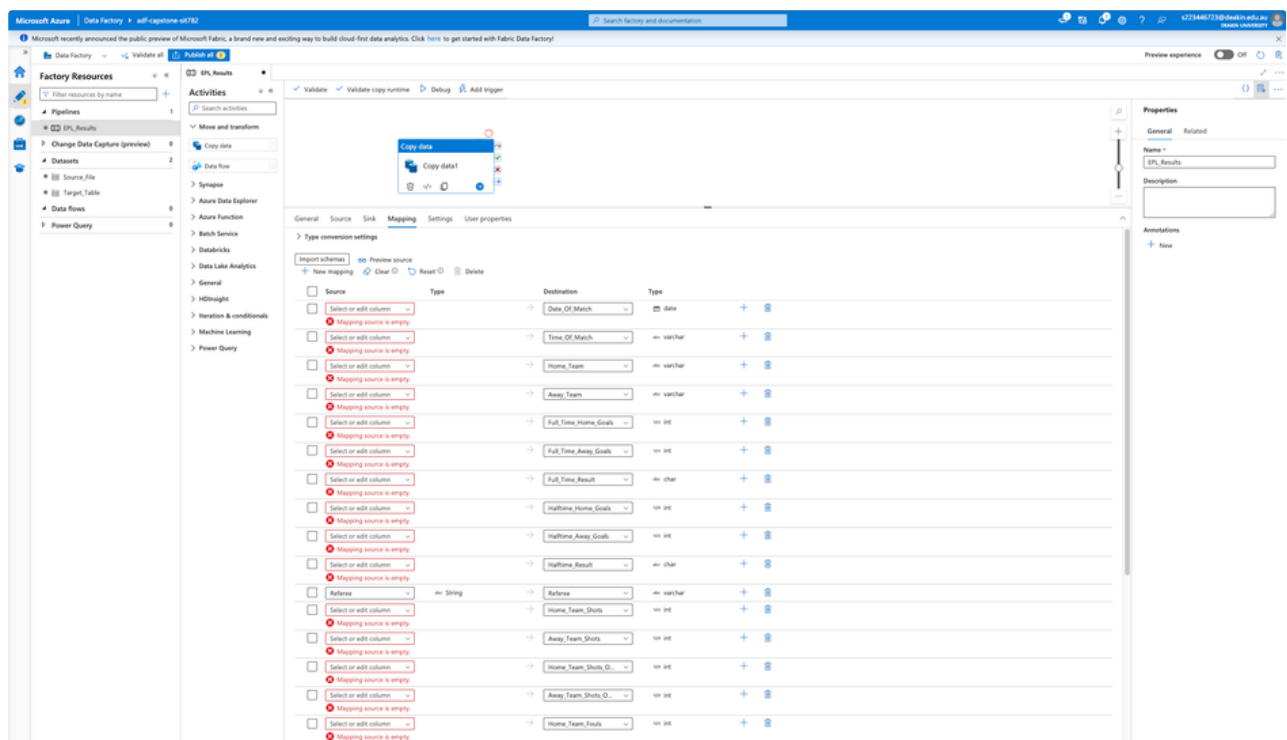
18. Next, we need to select the table where we want to load the data. Under "Table name", select "dbo.EPL_Results" and under "Import schema", select "From connection/store" to get the schema of the target table. This is critical for the next steps. Click on "OK" to create the Sink.



19. Now that we have created the Linked services for Source and Sink, and have also imported the source and target schemas, we will next need to map them. Select "Mapping" tab and click on "Import schemas".

20. Once you click on "Import schemas", both source and target schemas will appear as shown below. We will need to manually map each source column to target columns for the mapping to be complete. If the name of source and target column is similar, it will auto-map, as can be seen from the "Referee" column below. However, it is necessary to check each and every column to ensure proper mapping has been done.



21. Once you select all the mappings for source and sink, you will see that the error sign above the "Mapping" tab will disappear. Now, we will first save the pipeline, by clicking on "Publish all" on the top left. Clicking on that starts the validation process and publishing the pipeline saves your work.

22. Click on "Debug" next, this will trigger the pipeline. Once the pipeline completes, you will see something similar to this -



23. Now, check the table and confirm if the data has been loaded by logging into Azure Data Studio, as configured before. Selecting all the records from this small table will enable you to see if every field has been loaded properly.

If you have followed all the steps outlined above, you should be able to successfully provision a new Azure Data Factory, create Linked services and load data from Azure Blob Storage to Azure SQL Database. This can be used as template for any such data engineering/data warehouse activities in your assigned project.