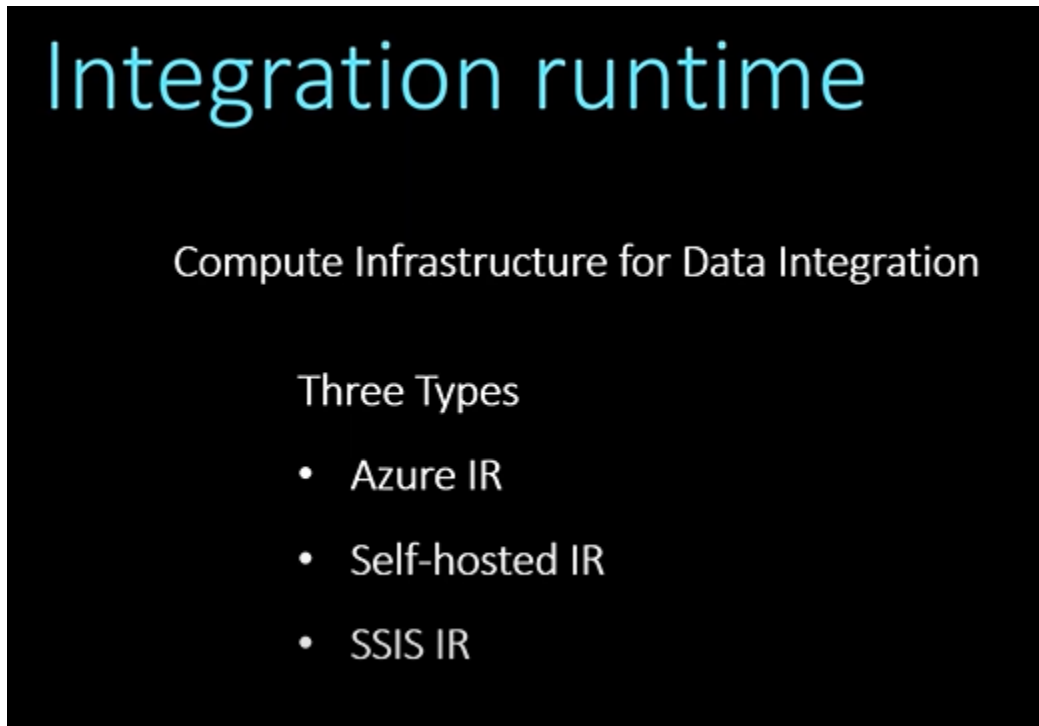Integration Runtime, Linked Service and Datasets in ADF

1. SO what is integration runtime? It is a compute infrastructure for data integration.. We know that ADF is mainly used for data integration from diff sources…so it might need some compute power to do this… this compute power can be obtained from integration runtime
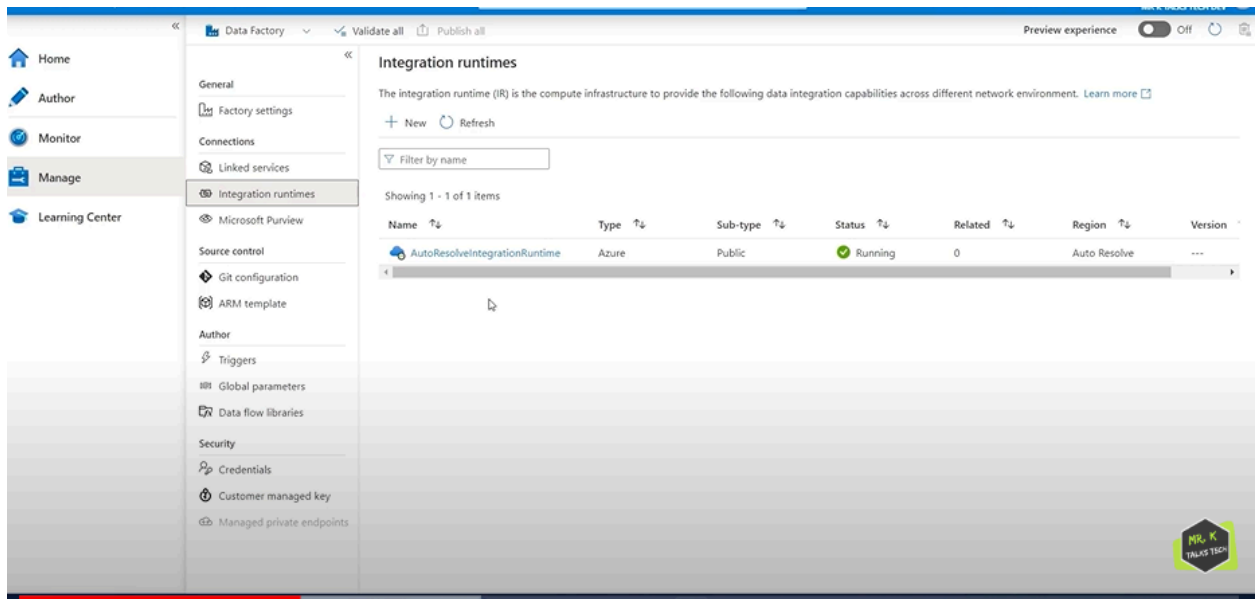2. There are three types



3. Azure IR is mainly used to connect to cloud based data sources…like ADF to ADL etc
4. Self hosted IR is mainly used to connect onPerm SQL DB…here we need to install a package on onperm System…then azure can connect with this with the help of package
5. SSIS IR…

**Scenario:** Imagine you have a well-established SSIS package that cleans and transforms data from various sources before loading it into a data warehouse on your on-premises SQL Server. This SSIS package is critical to your data pipeline.

**Challenge:** Migrating your data infrastructure to Azure, you don't want to rewrite the entire SSIS package from scratch. Ideally, you want to continue leveraging your existing SSIS workflows within your Azure Data Factory pipelines.

**SSIS IR to the rescue:** Here's where SSIS IR comes in. It creates a managed cluster of virtual machines in Azure with the SSIS engine pre-installed. You can then deploy your SSIS package to the SSIS IR and execute it seamlessly within your Azure Data Factory workflows. This allows you to migrate your existing SSIS logic to the cloud without significant recoding.

6. By default here we'll be having autoResovleIntegration runtime in our ADF



7. To create our own Integration runtime ..click in new
8. Next topic is Linked Service

9.

# Linked Service

- It is much like connection strings, which define the connection information needed for the ADF to connect to the data source
- More than 85 in-built linked service connectors are available inside ADF
- You need an Integration runtime to create a linked service connection

10. Datasets

1. **Data Migration from On-premises to Azure:**
- **Scenario:** You're migrating your data warehouse from an on-premises SQL Server database to Azure Synapse Analytics.
- **Linked Services:**
  - Create a Linked Service for the on-premises SQL Server, specifying server name, database name, username, and password (potentially using Self-Hosted Integration Runtime for private network access).
  - Create a Linked Service for Azure Synapse Analytics, defining connection details like workspace name and access key.

11.

# Datasets

- It is the structure/ format of the data
- You need to have a linked service connection to create a Dataset

12. Lets see how we can create datasets in azure

13. To create a dataset ..go to author tab



first we need to specify data source…next we will choose the formar

- **Scenario:** You want to copy data from a CSV file stored in Azure Blob Storage to a table in Azure SQL Database.
- **Linked Services:**
  1. Create a Linked Service for **Azure Blob Storage**. Provide details like your storage account name, access key, and the specific container holding the CSV file.
- **Dataset:**
  1. Create a Dataset for the **CSV file**. Specify the following:
     - The Linked Service you created for Azure Blob Storage (connecting ADF to the data source).
     - The path to the CSV file within the Blob Storage container.
     - The schema of the CSV data, including column names and data types (e.g., string, integer, date). This tells ADF how to interpret the data in the file.
- **ADF Pipeline:**
  1. In your data movement activity (e.g., copy activity), reference the Dataset you created. This instructs ADF on how to access the CSV file (using the linked service) and how to interpret its structure based on the defined schema.

14. Lets understand this with a better example

15. Consider we have source and destination



16. To perform this we'll use ADF to connect the source and destination



17. Here first our ADF will get data from the source
18. And first our ADF must connect to the source..to make that connection…we create Azure IR…and then next we create linked servie which is data lake gen2 …then next we create datasets..and here we need to specify the data format of Source

19. Next we have to connect ADF with destination ..and then we can execute copy activity in ADF



20. Now in copy activity we mention source and destination…and we execute this pipeline

21. If source is onperm SQL DB



Copy data using ADF copy data tool

1. Here what we'll do is



copy data from one ADL and paste it in 2nd ADL

2. Here in our resource group..we have 2 data lakes



3. And in the source data lake we have sample csv file



4. And in dest datalake we dont have any files
5. Lets go to ADF

6. Now click on copy data tool in ADF



7. In ADF the most commonly used activity is copy data..to transfer the data from one loc to another loc



here we'll be choosing built-in copy task and run the pipeline once

8. Now next we have to specify the data source

and in the connection..we have to specify Azure IR

**New connection**

Azure Data Lake Storage Gen2   Learn more ☐

Name *

source_datalake

Description

Connect via integration runtime * ⓘ

AutoResolveIntegrationRuntime                              ⌄

Authentication type

Account key                                                            ⌄

Account selection method ⓘ

⦿ From Azure subscription   ◯ Enter manually

Azure subscription ⓘ

Select all                                                              ⌄

Storage account name *

                                                                          ⌄   ↻

Test connection ⓘ

⦿ To linked service   ◯ To file path

Annotations

+ New

[ Create ]   [ Cancel ]                          🖉 Test connection

9. In the authentication type we choose account key(here every ADL has a key



)..now our ADF will connect to  ADL using these keys



10. Next we need to specify this
11. Next we test our connection and create this connection
12. Here we have create Azure IR, and linked service…next we have to create a dataset

13. For that first step would be selecting the csv file

## Source data store

Specify the source data store for the copy task. You can use an existing data store connection or specify a new data store.

Source type: Azure Data Lake Storage Gen2

Connection *: source_datalake   ✎ Edit   + New connection

**File or folder**

If the identity you use to access the data store only has permission to subdirectory instead of the entire account, specify the path to browse.

[                                    ]   📁 Browse

**Options**

☐ Binary copy ⓘ

Recursively ⓘ
☑

☐ Enable partition discovery ⓘ

Max concurrent connections ⓘ
[                                    ]

**Filter by last modified**

Start time (UTC)                End time (UTC)
[                    ]           [                    ] ⓘ

click on browse ..go to source and click on our file

## Browse

Select a file or folder.

Root folder > source

📄 SampleCSVFile_11kb.csv

**File or folder**

If the identity you use to access the data store only has permission to subdirectory instead of the entire account, specify the path to browse.

[ source/SampleCSVFile_11kb.csv ]   📁 Browse

14. Next we need to specify the file format of source

**Copy Data tool**

- ✓ Properties
- ② Source
  - ● Dataset
  - ● Configuration
- ③ Destination
- ④ Settings
- ⑤ Review and finish

**File format settings**

File format ⓘ

| DelimitedText | ∨ |   [ Detect text format ]   [ 👓 Preview data ]

Column delimiter

| Comma (,) | ∨ |

☐ Edit

Row delimiter

| Default (\r,\n, or \r\n) | ∨ |

☐ Edit

☐ First row as header ⓘ

〉 Advanced

Compression type

| None | ∨ |

Additional columns ⓘ

+ New

[ ‹ Previous ]  [ Next › ]

## 15. Next we setup the destination



## 16. Next we'll specify the folder path ..to where our data gets dumper

17. Next we configure file format ..which is csv
18. Now we give a name to our task

Copy Data tool

Properties ✓
Source ✓
Destination ✓
Settings ④
⑤ Review and finish

Settings

Enter name and description for the copy data task, more options for data movement

Task name *          copy_data_pipeline

Task description

Data consistency verification ⓘ    ☐
Fault tolerance ⓘ
Enable logging ⓘ    ☐
Enable staging ⓘ    ☐
❯ Advanced

19. Summary of what we have created

Copy Data tool

Properties ✓
Source ✓
Destination ✓
Settings ✓
Review and finish ⑤
Review ●
○ Deployment

Summary

You are running pipeline to copy data from Azure Data Lake Storage Gen2 to Azure Data Lake Storage Gen2.

Azure Data Lake Storage Gen2  ⟶  Azure Data Lake Storage Gen2

Properties                                                          ✎ Edit
Task name          copy_data_pipeline
Task description

Source                                                             ✎ Edit
Connection name        source_datalake
Dataset name           SourceDataset_qqg
Column delimiter       ,
Escape character       \
Quote char             "
First row as header    false
File name              SampleCSVFile_11kb.csv

❮ Previous    Next ❯

20. Now if we deploy our copy_data_pipeline..then we can our file in dest
21. Next we'll learn how we can do this copy data from scratch


Create a copy data pipeline from scratch

1. Here we'll create a datapipeline from scratch

2. We'll be using the same use case



3. We have sample file in our source ADL container



4. We'll copy this file and dump it in dest container in dest DL



5. Here we already have one copy data pipeline in our ADF..which has been created by copy data activity

6. Lets create this pipeline from scratch
7. Now we'll use AZURE IR to create a runtime



so here we use the default one
8. Next we need to create a linked service



here we choose ADL service..next

9. Now test the connection bw ADL and ADF using test connection..then next click on create

10. Next we create another linkedservice for dest ADL as well

our connection is successful for ADF and dest ADL

11. Next we need to create the datasets



Here we'll choose ADL as our data source…and select the format

12. Next we set properties to our data set and select linked service and file path





Here we have created dataset for our source          ...
Next we'll do the same for dest

13. To save all our work ..click on publish all



14. Now we have everything to create our pipeline
15. NExt to create a pipeline..go to author tab and click create pipeline

16. Now we'll drag the copy activity in our pipeline



17. Next here we have to config the source,sink

18. Click on source and select our source dataset



19. Next sink
20. Next click on publish to save our work
21. Now we'll use DEBUG mode and run our pipeline



22. Our pipeline is currently running

23. Now we can see our files in the dest Data lake



Triggers in DF and setting up Scheduled Trigger in ADF



1.

**Types**

- Schedule Trigger
- Storage Events Trigger
- Tumbling Window Trigger
- Custom Triggers

2. Types

3. Here we'll learn the Schedule trigger



**Schedule Trigger**

This trigger runs a pipeline on a specific schedule, such as hourly, daily, weekly, or monthly.

4.



Schedule triggers can be created and configured using the ADF portal

5. Lets see how we can schedule a trigger

6. Previously we have created a pipeline which copies the data from source DL to Dest DL



7. Now we'll add a file to our source DL and test our pipeline



here we have uploaded our file

8. And in our pipeline..we'll change the source path ..to the new file in the ADL

9.  Here  Trigger now is same as debug
10. Click on new/edit to add triggers..click bew trigger
11. Next we give properties for our trigger

 we can also specify an end
date to stop our trigger

## 12. Create the trigger and publish the changes

**Publish all**

You are about to publish all pending changes to the live environment. Learn more

**Pending changes (2)**

| NAME | CHANGE | EXISTING |
|------|--------|----------|
| ⌄ Datasets | | |
| ⊞ source_csv | (Edited) | source_csv |
| ⌄ Triggers | | |
| ⚡ scheduled_trigger | (New) | - |

Publish    Cancel

## 13. Here we can our pipeline has succeeded

14. Here we can see our pipeline ran successfully and copies the file



15. This pipeline will be run every two minutes