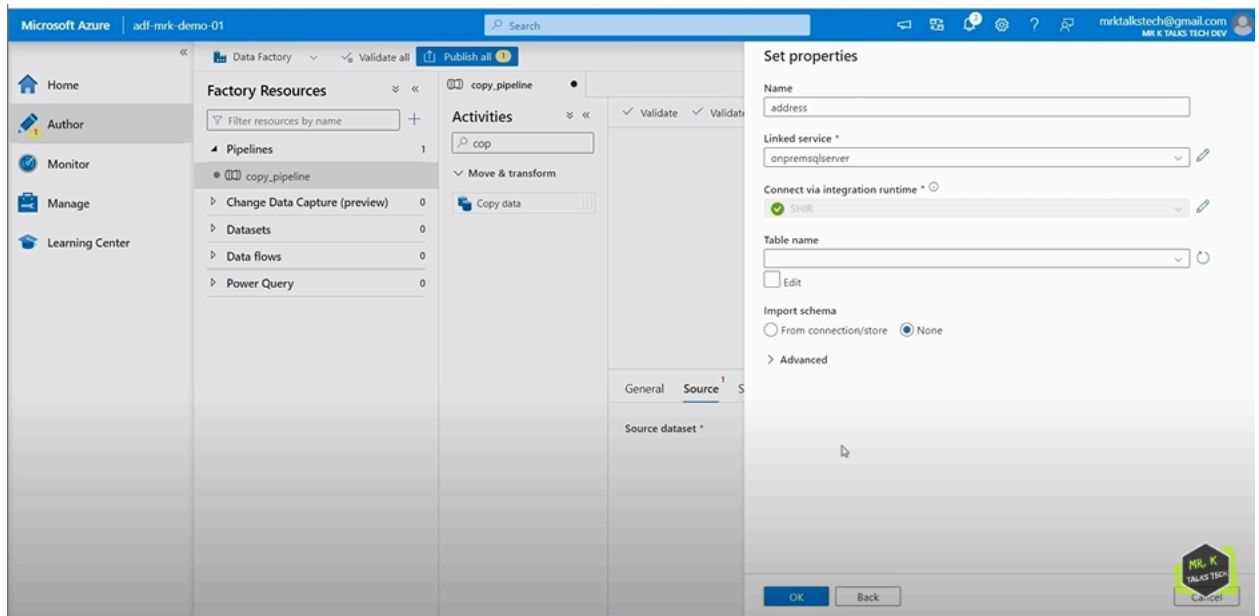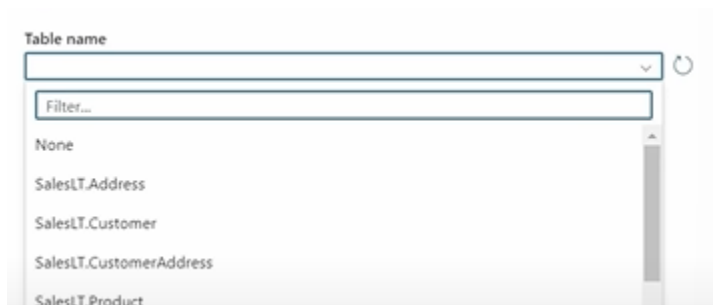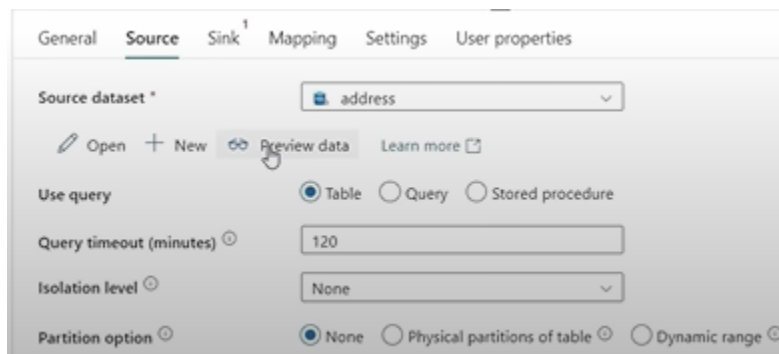Data Ingestion using ADF - part1

1. Now we have successfully integrated our on perm SQL server to the Azure Data factory using self host integrationruntime
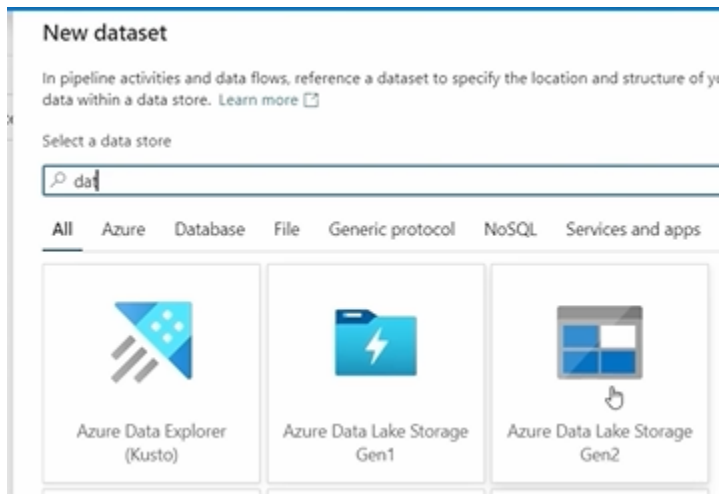


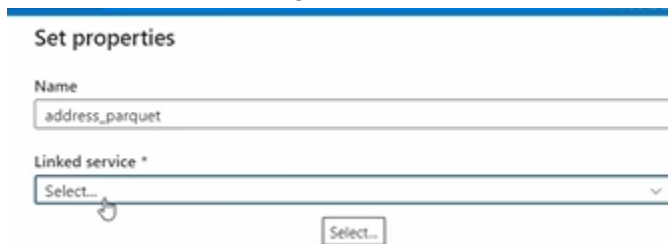2. Now we can choose the table names which are available in our onPerm



3. Now we'll just select address table for this demo..select it and then click OK
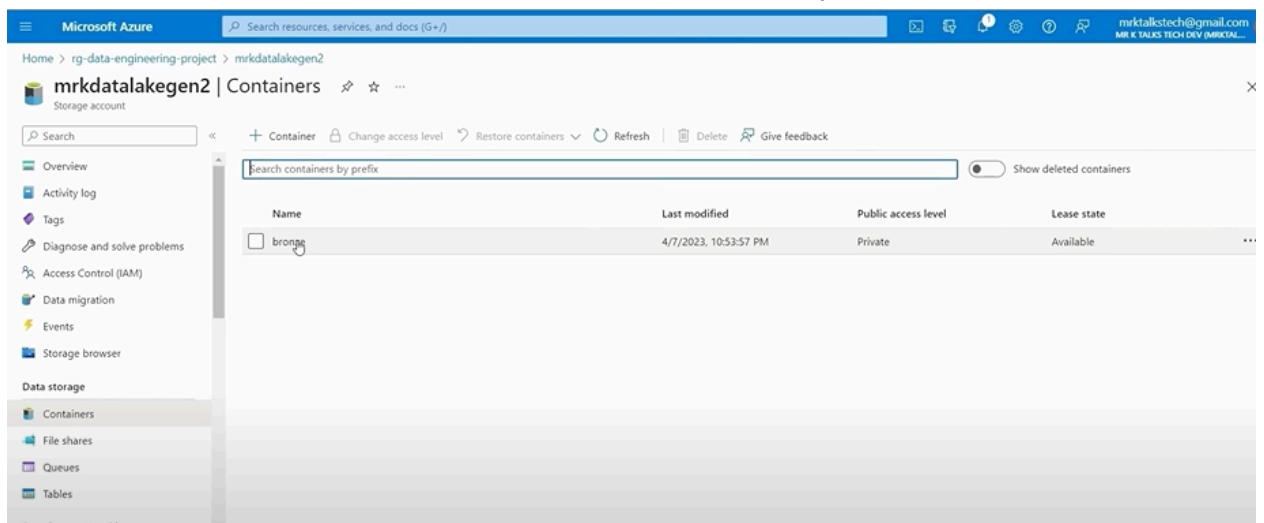4. And now we can preview our data too

5. Now lets configure the sink option…here our sink(dumping) would be ADL

**New dataset**

In pipeline activities and data flows, reference a dataset to specify the location and structure of y
data within a data store. Learn more

Select a data store

🔍 dat

| All | Azure | Database | File | Generic protocol | NoSQL | Services and apps |

Azure Data Explorer (Kusto)

Azure Data Lake Storage Gen1

Azure Data Lake Storage Gen2

6. Then next we have to give name and linked)service

**Set properties**

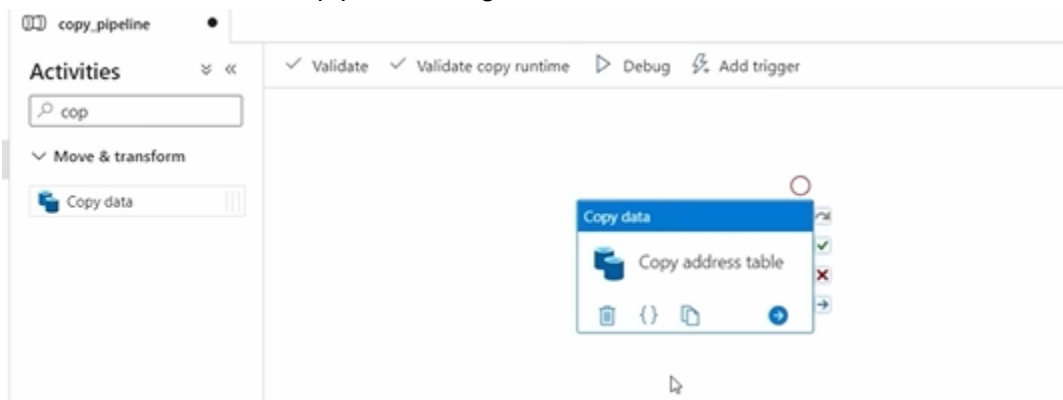Name

address_parquet

Linked service *

Select...

Select...

7. In linked_services…now our integration runtime would be AutoResolveIntegrationRuntime ..as we are moving data(within cloud) from ADF to ADL
8. We'll fill all the required things and click on test connection
9. Next we need to specify a location in ADL
10. So here we have created a container in ADL which is bronze layer



11. So we'll just give this as the file path…and click ok
12. Now after that..we have successfully completed source and sink
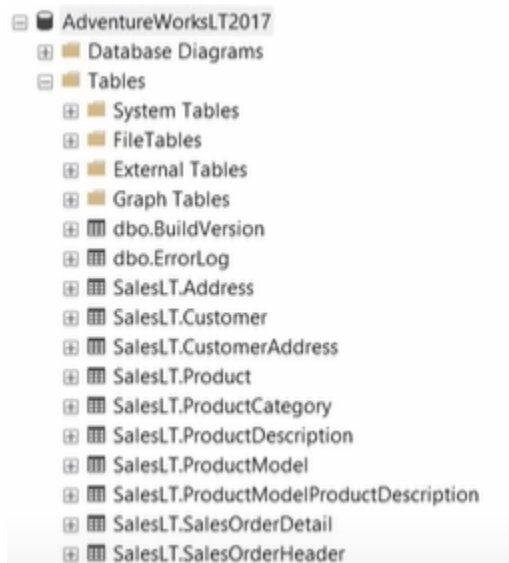
13. Next we'll run this demo pipeline using DEBUG



14. Now what this pipeline does is…copy the onPerm address table to ADL(bronze layer)


Data Ingestion using ADF - part2

1. Previously we have ingested only one table to show the demo
2. Now we ingest all the tables present in our OnPerm



3. We created a pipeline called copy all tables
4. Now we are actually interested in SalesLT schemas tables

5. SO we write a script to where we get schema name and table name

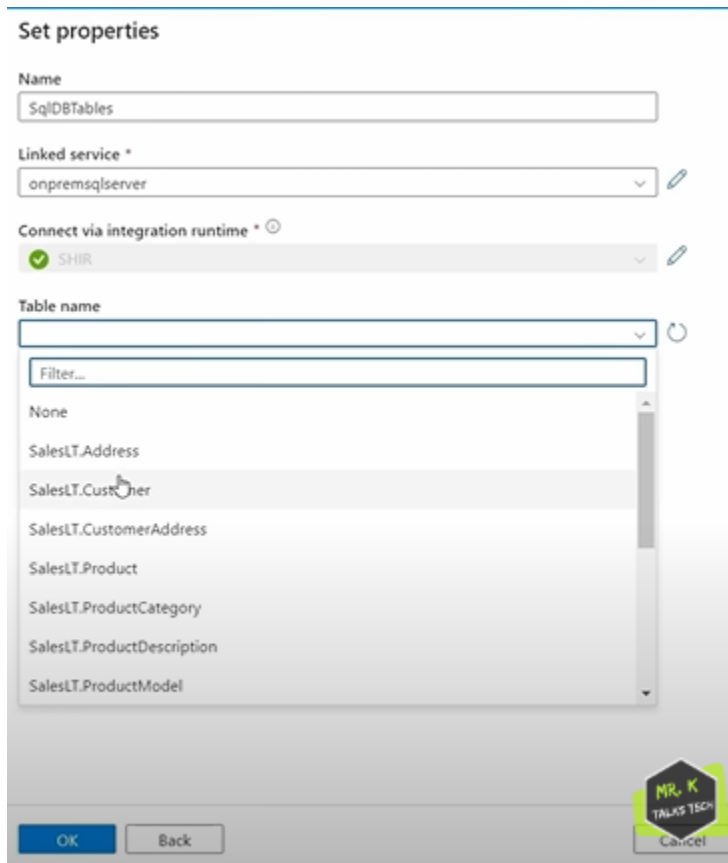| | SchemaName | TableName |
|---|---|---|
| 1 | SalesLT | Address |
| 2 | SalesLT | Customer |
| 3 | SalesLT | CustomerAddress |
| 4 | SalesLT | Product |
| 5 | SalesLT | ProductCategory |
| 6 | SalesLT | ProductDescription |
| 7 | SalesLT | ProductModel |
| 8 | SalesLT | ProductModelProductDescription |
| 9 | SalesLT | SalesOrderDetail |
| 10 | SalesLT | SalesOrderHeader |

```
SELECT
    s.name AS SchemaName,
    t.name AS TableName
FROM sys.tables t
INNER JOIN sys.schemas s
ON t.schema_id = s.schema_id
WHERE s.name ='SalesLT'
```

This gives list of all tables belongs to SalesLT schema

6. Now we'll copy this 10 tables to ADF
7. Next in our pipeline we'll add new activity which is LookUp activity



8. Next in the settings tab..we'll create a new Source dataset called SqlDBTables and choose our linked service which we created earlier and in the table name section..we'll
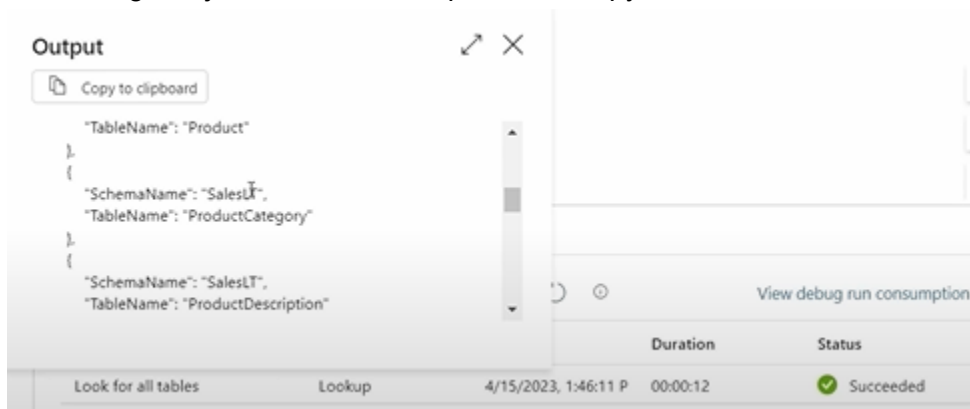
not choose any table and we click ok

**Set properties**

Name

SqlDBTables

Linked service *

onpremsqlserver

Connect via integration runtime * ⓘ

✓ SHIR

Table name

| Filter... |
| None |
| SalesLT.Address |
| SalesLT.Customer |
| SalesLT.CustomerAddress |
| SalesLT.Product |
| SalesLT.ProductCategory |
| SalesLT.ProductDescription |
| SalesLT.ProductModel |

OK   Back   Cancel

9. Now to retrieve the SalesLT schema's table in ADF..we use Query and give our query to

General  **Settings**  User properties

Source dataset    SqlDBTables

✎ Open  + New  👓 Preview data    Learn more ⧉

First row only        ✓

Use query       ○ Table  ● Query  ○ Stored procedure

Query *          INNER JOIN sys.schemas s    ✎ Edit
                ON t.schema_id = s.schema_id
                WHERE s.name ='SalesLT'

                Add dynamic content [Alt+Shift+D]

run   Query timeout (minutes) ⓘ   120

10. Now if we debug our copyalltables pipeline..then it gives us Input and output after it successfully debugged

11. Here using this json structured output ..we'll copy the tables from source



12. Next we'll add new activity called for each



13. Next we'll connect our lookup activity with for each



14. Now in the setting tabs of for each activity..we have to give items(dynamic content) and we give output of lookup activity

## Pipeline expression builder

Add dynamic content below using any combination of expressions, functions and system variables.

```
@activity('Look for all tables').output.value
```

Clear contents

**Activity outputs**   Parameters   System variables   Functions   Variables

🔍 Search

Look for all tables
Look for all tables activity output

Look for all tables
Look for all tables pipeline return value (preview)

Look for all tables count
Count of the rows

Look for all tables value array
Array of row data

MR. K
TALKS TECH

[ OK ]   [ Cancel ]

---

General   **Settings**   Activities (0)   User properties

Sequential          ☐

Batch count ⓘ       [                          ]

Items               [ @activity('Look for all tables').output... ]

15. Now inside these for each activity ..we'll add copy data activity



16. Next we have configure source and sink
17. First we set properties for source

18. Now in the query window of source..we'll add dynamic content



19. Our dynamic content is as follows


..here we'll be getting the output from lookup activity which is the outer loop's output


now we have configured the source successfully..next we'll configure the sink

20. FOr our sink we'll choose ADL with parquet format..and we set these properties



21. So inside the bronze layer..we'll need this type of folder structure



22. To add this type of structure..we go to sink dataset and open it

23. We'll add two parameters ..

 after that we can give

dynamic content values

24. We give this values..which are coming from lookup activity



25. Now in the connection tab..we'll give the folder structure using the parameters



26. Now we have completed our pipeline and we click on publish



27. We'll run our pipeline by using add trigger(trigger now)

28. Our pipeline is in progress



29. Now if we refresh our datalake..then we can see the data in the folders as we structured

30. Now we have officially completed the data ingestion process



Data Transformation using Databricks Part1

1. Now we'll use our databricks resource to make transformations on our data
2. So we need to create a compute cluster ..to work on notebooks
3. Click on create compute

4. We give default specifications for our cluster



5. Now we'll enable our  which gives us access to the ADL using our email credentials ..if we using the same email which is present in the shared blob data contributor



6. Next we'll click create cluster

7. Now lets create notebook



8. Now we'll create a notebook for mounting our ADL to the cluster



9. Now we'll change the source and the mount point

 next we'll execute this code

10. So what this means is ..we can access all the data in our bronze container ..using the mount point `dbutils.fs.ls("/mnt/bronze")`

11. To see the tables we use

```
1   dbutils.fs.ls("/mnt/bronze/SalesLT/")
```

```
Out[3]: [FileInfo(path='dbfs:/mnt/bronze/SalesLT/Address/', name='Address/', size=0, modificationTime=1681529111000),
 FileInfo(path='dbfs:/mnt/bronze/SalesLT/Customer/', name='Customer/', size=0, modificationTime=1681529113000),
 FileInfo(path='dbfs:/mnt/bronze/SalesLT/CustomerAddress/', name='CustomerAddress/', size=0, modificationTime=168152911700
0),
 FileInfo(path='dbfs:/mnt/bronze/SalesLT/Product/', name='Product/', size=0, modificationTime=1681529112000),
 FileInfo(path='dbfs:/mnt/bronze/SalesLT/ProductCategory/', name='ProductCategory/', size=0, modificationTime=168152911600
0),
 FileInfo(path='dbfs:/mnt/bronze/SalesLT/ProductDescription/', name='ProductDescription/', size=0, modificationTime=16815
```

12. Now we'll create mounts for silver and for gold too
13. *Important point : Since we have used credential passthru..it is not always mandatory to mount the container..instead..we can just use full datalake path..to access the data inside a container
14. Now we'll get the data from the bronze layer and do some transformations then load to silver