# Extract and process data while maintaining security, monitoring and standards on the Azure Cloud platform.

Problem Statement:  Data engineers usually receive  requirements for data processing & transformation from different sources, Azure offers a wide range of services & tools to help you achieve this.
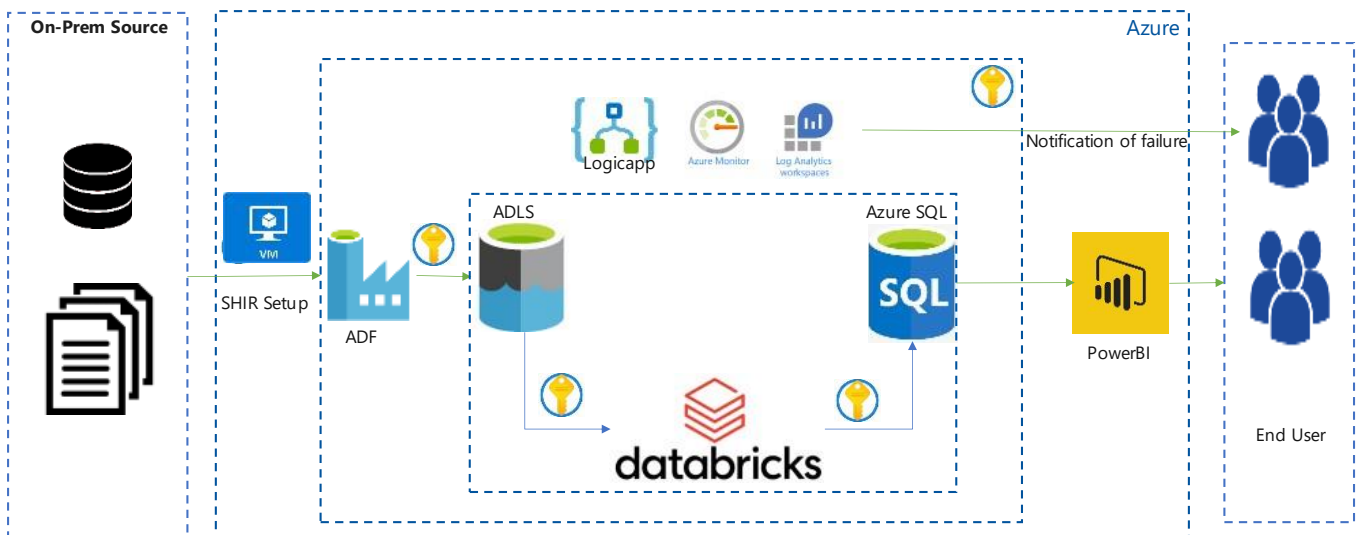
We received a requirement to ingest files from on-premises to azure cloud platform & perform data quality checks while maintaining security, monitoring & data lake transformation layers before loading data into data warehouse.

## Services use for Transformation

1. **Azure account:** For using cloud platform
2. **Azure Data lake : -** For storing data & maintaining transformation layer
3. **Azure DataFactory:-** For data ingestion from onprem to azure & data processing
4. **Azure Databricks:** For data transformation & delta tables
5. **Azure Synapse:** For DWH & data modelling purpose
6. **Azure Keyvualt:** For maintaining credentials & secrets at centralize location.
7. **Azure monitor** : For monitor Azure resources
8. **Azure log analytics:** For querying logs
9. **Logic Apps:** For sending email on failure & success
10. **Azure VM :** For setup SHIR

## Architecture Diagram:

# Logical architecture

## Technical requirements:

Set up Azure environment with listed services

## 1.Azure data lake setup & maintain 3 different layers in ADLS

- RAW: Store raw data from source

- Refined: store cleansed data after processing

- Processed: store transformed data



| | Name |
|---|---|
| ☐ | $logs |
| ☐ | proccesssed |
| ☐ | raw |
| ☐ | refined |

## 2. Set up VM for self-hosted IR to connect on prem server to pick up file

| Name ↑↓ | Type ↑↓ | Sub-type ↑↓ | Status ↑↓ |
|---|---|---|---|
| ☁ AutoResolveIntegrationR... | Azure | Public | ✓ Running |
| 🖥 integrationRuntime1 | Self-Hosted | --- | ✓ Running |

## 3. Set Up Azure Data Factory



The following checks have been performed in ADF for data consistency:

- Check that the file is available in the path. If it's not, there should be a timeout after 1 minute:



- Check whether the file size is greater than 20b or not. If the file size is greater than 20b then it needs to be processed:

We used 2 get metadata to check first list of files & then stored filename in variable to check size of each file.

**Get metadata output passed to foreach loop activity to run loop for every file**



**Inside for each loop used variable to store filename & using 2nd getmetadata get file size.**
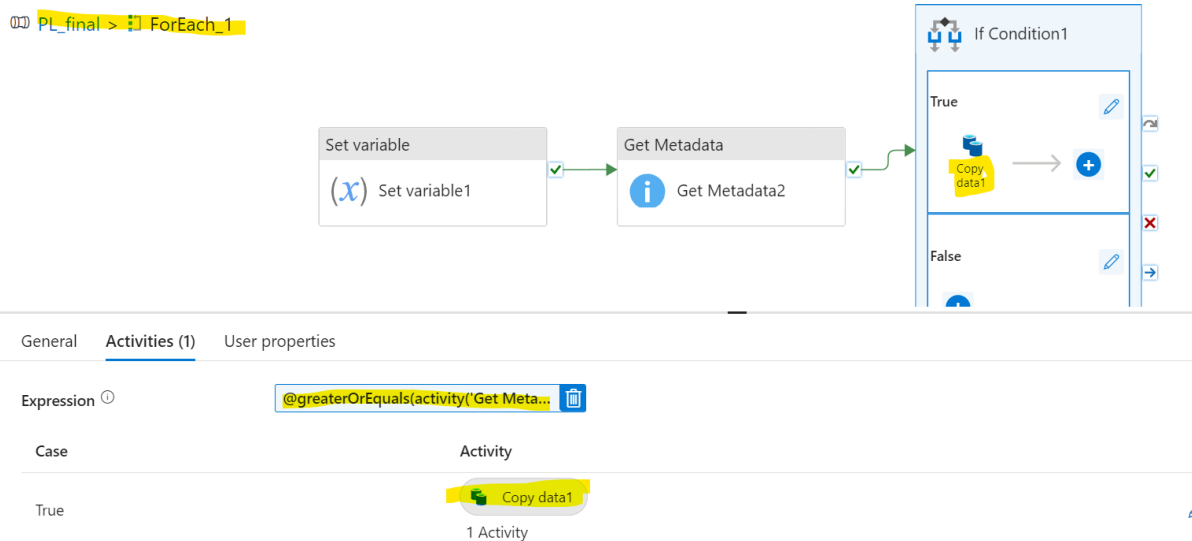
Passed size condition in ifcondition activity to pick files which are greater then 20b



```
@greaterOrEquals(activity('Get Metadata2').output.size,21)
```

PL_final > ForEach_1

Set variable — Set variable1

Get Metadata — Get Metadata2

If Condition1
True — Copy data1
False

| General | Activities (1) | User properties |

Expression ⓘ   @greaterOrEquals(activity('Get Meta...

| Case | Activity |
| --- | --- |
| True | Copy data1 |
| | 1 Activity |

- Setup email configuration on failure & success of pipeline to get notify

Create master pipeline to call logic app to get notified on pipeline failure



Web — on_Success_email

Execute Pipeline — Execute Pipeline1 PL_final

Web — on_Failure_mail

4.     Configured Logicapp to trigger email & passed Logicapp URL in be activity:

**5. Setup keyvault to store credentials**



Once all the checks have been performed, the DataFactory pipeline is run to maintain the transformation layer in ADLS.

Setup Databricks :

In Databricks couple of transformation applied , I am attaching some of them only as per confidentiality,

- If date is NULL or blank, give default date as '2020-11-28'. Format of date column should be YYYY-MM-DD.
- Remove the entries which has URL field value as 'ERROR'.
- Transform the values of column country with their acronyms. For eg: Austria would be replaced by 'AUST', Belgium by 'BELG' etc.

```python
#read the parquet file from adls
source_file="abfss://refined@chaitanyacapastonelake.dfs.core.windows.net"
df_source=spark.read\
    .format('parquet')\
    .option('inferschema',True)\
    .load(source_file)

adf=df_source
```

Cmd 3

```python
# Apply transformations

target_df = adf.withColumn("date", when(col("date").isNull(), "2020-11-28").otherwise(col("date")))
target_df = target_df.filter(col("url") != "ERROR")
```

```
24    #run the loop
25    for rows in table1.find_all('tr'):
26      column=rows.find_all('td')
27      if len(column)>=2:
28        country_name=column[0].text.strip()
29        acronym=column[1].text.strip()
30        country_acronym[country_name]=acronym
31
32    #covert the column 'country' of transformed table to upper-case
33    target_adf=target_df.withColumn("country",upper(col("country")))
34
35    # Create a DataFrame from the acronym data
36    acronym_df = spark.createDataFrame(country_acronym.items(), ["country", "acronym"])
37
38    # Join the original DataFrame with the acronym DataFrame to replace values
39    adf1 = target_adf.join(acronym_df, "country", "left").select("*")
40
41    # Interchange the positions of two columns (e.g., swap "Age" and "Country") and then drop country table
42    result_df = adf1.withColumnRenamed("country", "temp").withColumnRenamed("acronym", "country").withColumnRenamed
      ("temp", "acronym").drop("acronym")
43
```

Complete code & ARM template added in github.

Once data transformation done loading data in sql to build pbi dashboard on so can be consumed be end user.

## Challenges

Integrate all Azure services together in a secure way with the help of VNET & Keyvault, so data can be processed from on prem to Azure & loaded successfully into database for consumption.

During the process, the criticality of data is the most important thing to protect so that the whole flow works smoothly as per the architecture diagram.

## Benefits

In all projects as data engineering will play a key role to perfoem various data operation & maintain security so its help of this high level flow can implement security best practices can be followed.