

# Step-by-Step Implementation in Azure Data Factory (ADF)

## Scenario:

We have four data sources:

1. An on-premises SQL Server producing CSV files.
2. An Azure Blob Storage containing JSON files.
3. An AWS S3 bucket with Avro files.
4. An Azure Data Lake Storage (ADLS) with Parquet files.

## Objective:

Ingest data from these sources into a central Azure Data Lake Storage (ADLS) and perform transformations.

## Step 1: Create Linked Services

**Linked Services** are connections to data sources. Each data source will have its own linked service.

1. **On-premises SQL Server (CSV files):**
  - Use Self-hosted Integration Runtime to connect to on-premises SQL Server.
  - Define the connection string and credentials.
2. **Azure Blob Storage (JSON files):**
  - Create a linked service to connect to Azure Blob Storage.
  - Provide the storage account name and access key.
3. **AWS S3 Bucket (Avro files):**
  - Create a linked service to connect to AWS S3.
  - Provide the access key and secret key for the AWS account.
4. **Azure Data Lake Storage (Parquet files):**
  - Create a linked service to connect to ADLS.
  - Provide the storage account name and access key.

## Step 2: Create Datasets

**Datasets** represent the data structures within the data stores, typically pointing to the data you want to use in your activities.

1. **CSV Dataset:**
  - Define the schema (columns) and file path in the on-premises file system.
2. **JSON Dataset:**
  - Define the schema and file path in Azure Blob Storage.
3. **Avro Dataset:**
  - Define the schema and file path in the AWS S3 bucket.
4. **Parquet Dataset:**
  - Define the schema and file path in Azure Data Lake Storage.

## Step 3: Create Pipelines

**Pipelines** are the workflows that orchestrate data movement and transformation activities.

1. **Copy Activity:**
  - For each data source, create a copy activity to move data from the source to the destination (ADLS).
2. **Transformation Activities:**
  - Use Data Flow activities if transformations are required (e.g., data cleaning, aggregations).

## Summary

1. **Create Linked Services** to connect to each data source (on-premises SQL Server, Azure Blob Storage, AWS S3, and ADLS).
2. **Create Datasets** to define the schema and format of the data for each file type (CSV, JSON, Avro, and Parquet).
3. **Create Pipelines** with copy activities to orchestrate the data movement from the sources to the central ADLS.

## Types of Activities in ADF

### 1. Copy Activity:

- Used to copy data from a source to a destination.
- Supports various file formats and data stores.
- **Purpose:**
- The Copy Activity is designed to copy data from a source to a destination. It supports a wide range of data stores, including databases, file systems, and cloud storage services.
- **When to Use:**
- Use the Copy Activity when you need to move or copy data from one location to another, such as from an on-premises SQL Server to Azure Data Lake Storage (ADLS), or from an AWS S3 bucket to Azure Blob Storage.
- **Limitations:**
- While the Copy Activity is powerful, it may have limitations in terms of handling complex transformations. For more advanced transformations, you may need to use Data Flow Activity or other transformation tools.
- Performance can be affected by the size of the data and the network bandwidth between the source and the destination.

### 2. Data Flow Activity:

- Used for data transformation using mapping data flows.
- Allows complex transformations like data cleansing, aggregation, and joins.
- The Data Flow Activity allows you to create data transformation logic using a graphical interface. It is built on Apache Spark and is designed for big data processing.
- **When to Use:**
- Use the Data Flow Activity when you need to perform complex data transformations, such as data cleansing, aggregations, joins, and data format conversions.
- **Limitations:**
- Data Flow Activity can be resource-intensive and may require significant compute resources, especially for large datasets.
- It may have a learning curve for users who are not familiar with data transformation concepts and the graphical interface.

### 3. **Lookup Activity:**

- Retrieves a dataset from any of the supported data stores.
- Used to return data for use in subsequent activities.
- **Purpose:**
- The Lookup Activity retrieves data from a specified dataset. It is often used to return configuration or control data that is used in subsequent activities.
- **When to Use:**
- Use the Lookup Activity when you need to retrieve a single value or a small dataset that will be used to control the flow of the pipeline.
- **Limitations:**
- The Lookup Activity is not designed for retrieving large datasets. It is best suited for small datasets or single values.
- Performance can be impacted if the query retrieves a large number of rows.
- 

### 4. **ForEach Activity:**

- Iterates over a collection and executes specified activities for each item.
- **Purpose:**
- The ForEach Activity iterates over a collection and executes specified activities for each item in the collection.
- **When to Use:**
- Use the ForEach Activity when you need to perform repetitive tasks for each item in a collection, such as processing multiple files in a directory or iterating over a list of database records.
- **Limitations:**
- The ForEach Activity can lead to performance issues if the collection contains a large number of items, as it will execute the nested activities for each item sequentially or in parallel.
- Parallel execution is limited by the number of available resources and concurrency settings.

## 5. **Execute Pipeline Activity:**

- Invokes another pipeline from the current pipeline.
- **Purpose:**
- The Execute Pipeline Activity allows you to invoke another pipeline from within the current pipeline. This is useful for modularizing and reusing pipeline logic.
- **When to Use:**
- Use the Execute Pipeline Activity when you need to call a reusable pipeline or when you want to break down a complex workflow into smaller, manageable pipelines.
- **Limitations:**
- The Execute Pipeline Activity adds an additional layer of complexity to the workflow, which may make debugging and monitoring more challenging.
- There may be limitations on the depth of nested pipeline executions.