## 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.

Azure Data Factory (ADF) provides a range of **activities** that facilitate data ingestion and transformation processes. These activities are used in pipelines to create workflows for data movement, transformation, and integration. Below, I explain the different types of activities available in Azure Data Factory with examples, use cases, and their limitations.

---

Types of Activities in Azure Data Factory

## 1. Data Movement Activities

These activities move data between data sources.

### Copy Activity

- **Description**: Transfers data from a source to a destination. It supports a wide range of formats (e.g., CSV, JSON, Parquet, Avro).
- **Why Use**:
    - To move data between on-premises, cloud storage, databases, or REST APIs.
    - Efficient for **batch processing**.
- **When to Use**:
    - Ingesting raw data from multiple sources into a data lake.
    - Transferring data between different storage accounts or databases.
- **How to Use**:
    1. Create a pipeline in Azure Data Factory.
    2. Add a Copy Activity and configure the **source** and **sink** datasets.
    3. Map columns between the source and destination if required.
    4. Use **linked services** to connect to data sources (e.g., Blob Storage, SQL Server).
- **Example**:
    - Move CSV files from an on-premises SFTP server to Azure Blob Storage.
- **Limitations**:
    - No complex transformations (only direct column mapping or file format conversion).
    - Latency may increase with large file sizes.

## 2. Data Transformation Activities

These activities allow data manipulation.

**Data Flow Activity**

- **Description**: Enables **visual data transformations** using the **Mapping Data Flow** feature.
- **Why Use**:
    - Perform ETL (Extract, Transform, Load) operations like filtering, aggregations, joins, and derived columns.
- **When to Use**:
    - Converting unstructured data into structured formats.
    - Aggregating sales data by region and date.
- **How to Use**:
    1. Add a Data Flow Activity to your pipeline.
    2. Create a Data Flow and design a visual transformation.
    3. Link source data (e.g., JSON or Parquet) and define transformations (e.g., remove duplicates, join tables).
    4. Define the output (sink) dataset.
- **Example**:
    - Aggregate daily sales data from multiple regions into a single table for analytics.
- **Limitations**:
    - Requires an **Azure Integration Runtime**.
    - Higher cost compared to using SQL-based transformations.

### 3. Control Flow Activities

These activities control pipeline execution flow.

**ForEach Activity**

- **Description**: Iterates over an array or a list of items (e.g., files, table rows) and executes child activities for each item.
- **Why Use**:
    - To handle multiple input files or process data from multiple tables.
- **When to Use**:
    - Processing daily files (e.g., 100 CSV files in a folder).
    - Looping through a list of databases.
- **How to Use**:
    1. Add a ForEach Activity to your pipeline.
    2. Define the input (array or list) using dynamic expressions.
    3. Add child activities (e.g., Copy Activity) inside the ForEach loop.
- **Example**:
    - Process multiple CSV files uploaded to Azure Blob Storage daily.
- **Limitations**:
    - Nested ForEach loops are not allowed.
    - Limited to sequential or batch execution.

---

**If Condition Activity**

- **Description**: Executes activities based on a condition (true/false).
- **Why Use**:
    - To implement decision-making logic in pipelines.
- **When to Use**:
    - Performing different tasks based on file size, type, or timestamp.
- **How to Use**:
    1. Add an If Condition Activity to the pipeline.
    2. Define the condition using dynamic expressions.
    3. Add child activities under the **True** or **False** branches.
- **Example**:
    - Process large files differently than small files.
- **Limitations**:
    - Complex conditions can be hard to debug.

### Wait Activity

- **Description**: Pauses the pipeline execution for a specific duration.
- **Why Use**:
  - To introduce delays (e.g., waiting for file generation).
- **When to Use**:
  - Ensuring upstream systems complete their processes before proceeding.
- **How to Use**:
  1. Add a Wait Activity to the pipeline.
  2. Set the wait time in seconds.
- **Example**:
  - Wait for 5 minutes before processing files uploaded to a folder.
- **Limitations**:
  - Inefficient for long delays.

### 4. Data Integration Activities

These activities involve system-level integration.

### Lookup Activity

- **Description**: Retrieves data (e.g., single value or table rows) from a data source.
- **Why Use**:
  - To fetch metadata or configuration settings.
- **When to Use**:
  - Fetching table names or file paths dynamically.
- **How to Use**:
  1. Add a Lookup Activity to the pipeline.
  2. Configure the source dataset and query.
- **Example**:
  - Retrieve file paths from a SQL table for processing.
- **Limitations**:
  - Limited output size (64 KB).

### Web Activity

- **Description**: Calls REST APIs to perform custom tasks (e.g., trigger Logic Apps or call external services).
- **Why Use**:
    - To integrate with third-party services.
- **When to Use**:
    - Fetching data from external APIs or triggering notifications.
- **How to Use**:
    1. Add a Web Activity to the pipeline.
    2. Specify the API endpoint, HTTP method, and headers.
    3. Use dynamic expressions to pass parameters.
- **Example**:
    - Trigger a Logic App to notify stakeholders about file uploads.
- **Limitations**:
    - Requires the API to be accessible publicly or via private endpoints.

---

**Example Scenario**

Use Case: Processing Multiple File Formats from Multiple Sources

1. **Scenario**:
    - **Source 1**: On-premises FTP server generates daily CSV files.
    - **Source 2**: REST API provides JSON files in real-time.
    - **Output**: Data is transformed into a unified Parquet format and stored in Azure Data Lake.
2. **Pipeline Design**:
    - **Copy Activity**: Move CSV files from the FTP server to Azure Blob Storage.
    - **Web Activity**: Call the REST API to fetch JSON files and store them in Blob Storage.
    - **Data Flow Activity**: Merge CSV and JSON data into a unified Parquet format.
    - **ForEach Activity**: Iterate through all uploaded files and process them.
    - **If Condition Activity**: Check the file size and route it for special processing if it exceeds a threshold.
3. **Output**:
    - All processed data is stored in Azure Data Lake as Parquet files, ready for analytics.