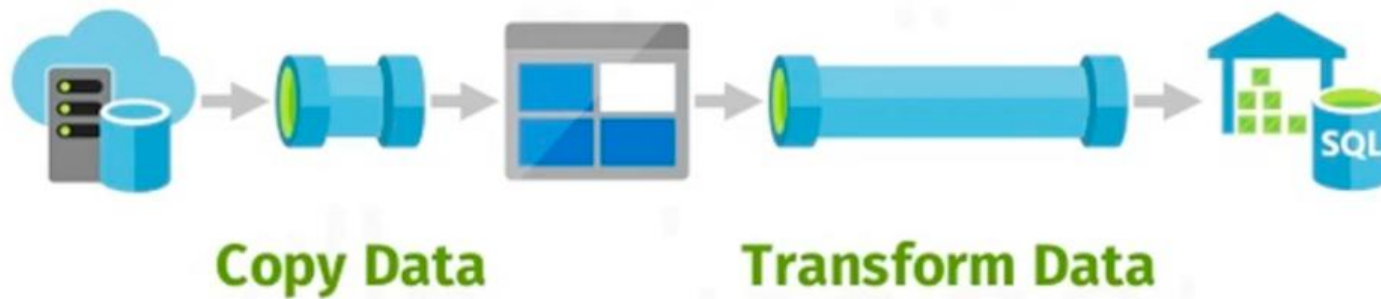


# What can you do in Azure Data Factory?



## Copy Data

More than 80 connectors to different services are available



## Transform Data

Using newly added Data Flow, now Data Factory is complete cloud based ETL tool.





Azure Data Factory

## Definition:

Azure Data Factory (ADF) is a hybrid data integration service that enables you to quickly and efficiently create automated data pipelines – without having to write any code!



**Azure Data Factory**

- Hybrid Data Integration Service
- Simplifies ETL at scale
- Enables modern data integration
- Drag and drop interface
- Over 80 connectors available
- Move, transform and save data
- Managed Service
- Create Data Driver workflows
- Orchestrate and automate data movement
- Transform and store data
- Operationalize the process
- ETL or ELT scenarios

# Data Factory on Azure Ecosystem

01

## Migration?

Data Factory excels in periodic data loads and transformation instead.



02

## Streaming?

ADF can orchestrate, but there are other dedicated services for streaming



03

## Transformations?

Data flows for simple ones, but you can use Databricks or HDInsight for more complex transforms



# SSIS vs Data Factory

## SSIS

More code-free transformations  
On Premises connectors (e.g excel)

## Data Factory

Much higher scalability  
Cloud and SaaS Connectors  
Event based Triggers  
Can use SSIS Packages



## Data Factory considerations

### Two versions

ADF V2 is the current and improved version

### Build options

PowerShell, .Net, Python, REST, ARM

### Highly integrated

DevOps, Key Vault, Monitor, Automation

### No data storage

Need to persist data by the end.

### Security standards

HTTP/TLS whenever possible





Deploying Clusters — Dask docu x

Dask: Status x

python 3.x - ImportError: cannot x

Introduction to Azure Data Facto x

+

https://docs.microsoft.com/en-us/azure/data-factory/introduction

Bookmarks Razorpay Dashboard BitPaper - Features Home TL-WR740N Privacy error (55) DSC Error in G... Manage excellatee... GitHub - Pierian-Da... Your Account Learning on Simpli... Reading list

Filter by title

Data Factory Documentation

Switch to version 1 documentation

Overview

Introduction to Data Factory

What's new in Azure Data Factory

Compare current version to version 1

Quickstarts

Create data factory - User interface (UI)

Create data factory - Copy data tool

Create data factory - Azure CLI

Create data factory - Azure PowerShell

Create data factory - .NET

Create data factory - Python

Create data factory - REST

Create data factory - ARM template

Create data flow

Tutorials

Samples

Concepts

How-to guides

Download PDF

for business intelligence (BI) applications to consume. Ultimately, through Azure Data Factory, raw data can be organized into meaningful data stores and data lakes for better business decisions.

## Code-Free ETL as a Service

INGEST

- Multi-cloud and on-prem hybrid copy data
- 90+ native connectors
- Serverless and auto-scale
- Use wizard for quick copy jobs

CONTROL FLOW

- Design code-free data pipelines
- Generate pipelines via SDK
- Utilize workflow constructs: loops, branches, conditional execution, variables, parameters, ...

DATA FLOW

- Code-free data transformations that execute in Spark
- Scale-out with Azure Integration Runtime
- Generate data flows via SDK
- Designers for data engineers and data analysts

SCHEDULE

- Build and maintain operational schedules for your data pipelines
- Wall clock, event-based, tumbling windows, chained

MONITOR

- View active executions and pipeline history
- Detail activity and data flow executions
- Establish alerts and notifications

## How does it work?

Data Factory contains a series of interconnected systems that provide a complete end-to-end platform for data engineers.

This visual guide provides a high-level overview of the Data Factory architecture:

WHAT IS AZURE DATA FACTORY?

DATA INTEGRATION

SCENARIO

HOW DATA FACTORY WORKS

WHY THE DATA FACTORY ARCHITECTURE WORKS

In this article

How does it work?

Top-level concepts

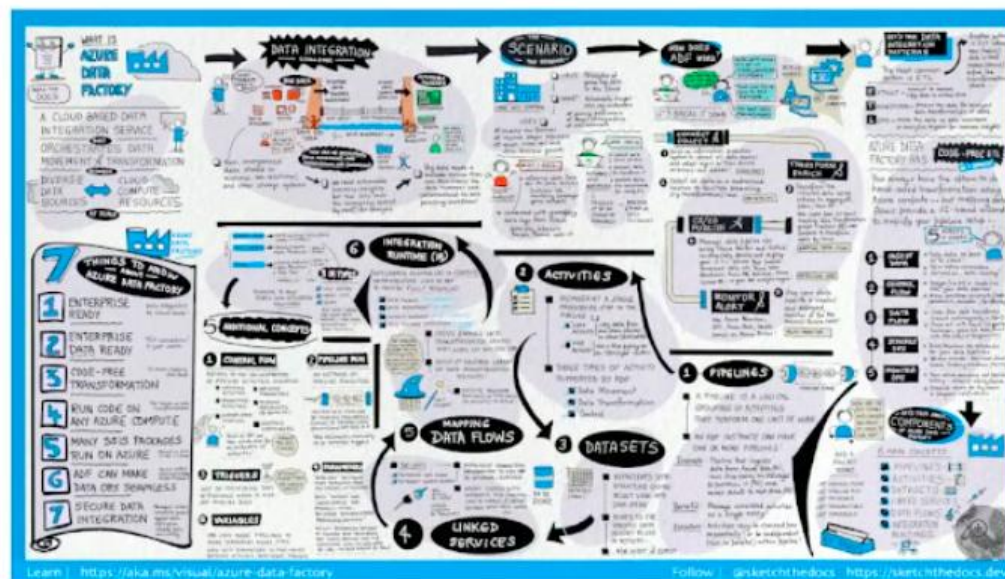
Next steps

Neuron

## How does it work?

Data Factory contains a series of interconnected systems that provide a complete end-to-end platform for data engineers.

This visual guide provides a high-level overview of the Data Factory architecture:



To see more detail, click the preceding image to zoom in, or browse to the high resolution image.

## Connect and collect

### In this article

[How does it work?](#)

[Top-level concepts](#)

[Next steps](#)

Filter by title

Data Factory Documentation

Switch to version 1 documentation

### Overview

Introduction to Data Factory

What's new in Azure Data Factory

Compare current version to version 1

### Quickstarts

Create data factory - User interface (UI)

Create data factory - Copy data tool

Create data factory - Azure CLI

Create data factory - Azure PowerShell

Create data factory - .NET

Create data factory - Python

Create data factory - REST

Create data factory - ARM template

Create data flow

### Tutorials

### Samples

### Concepts

### How-to guides

Download PDF

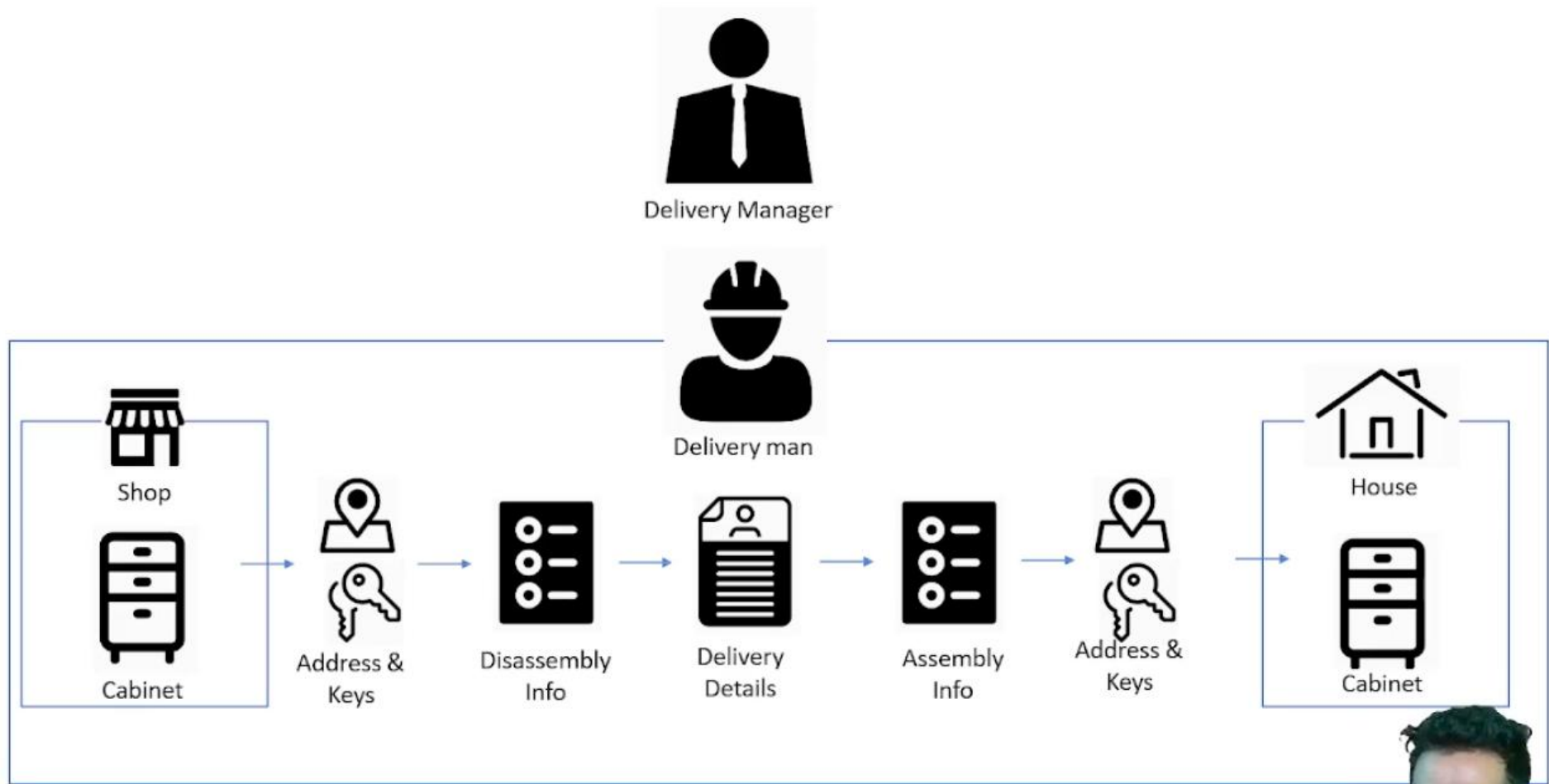


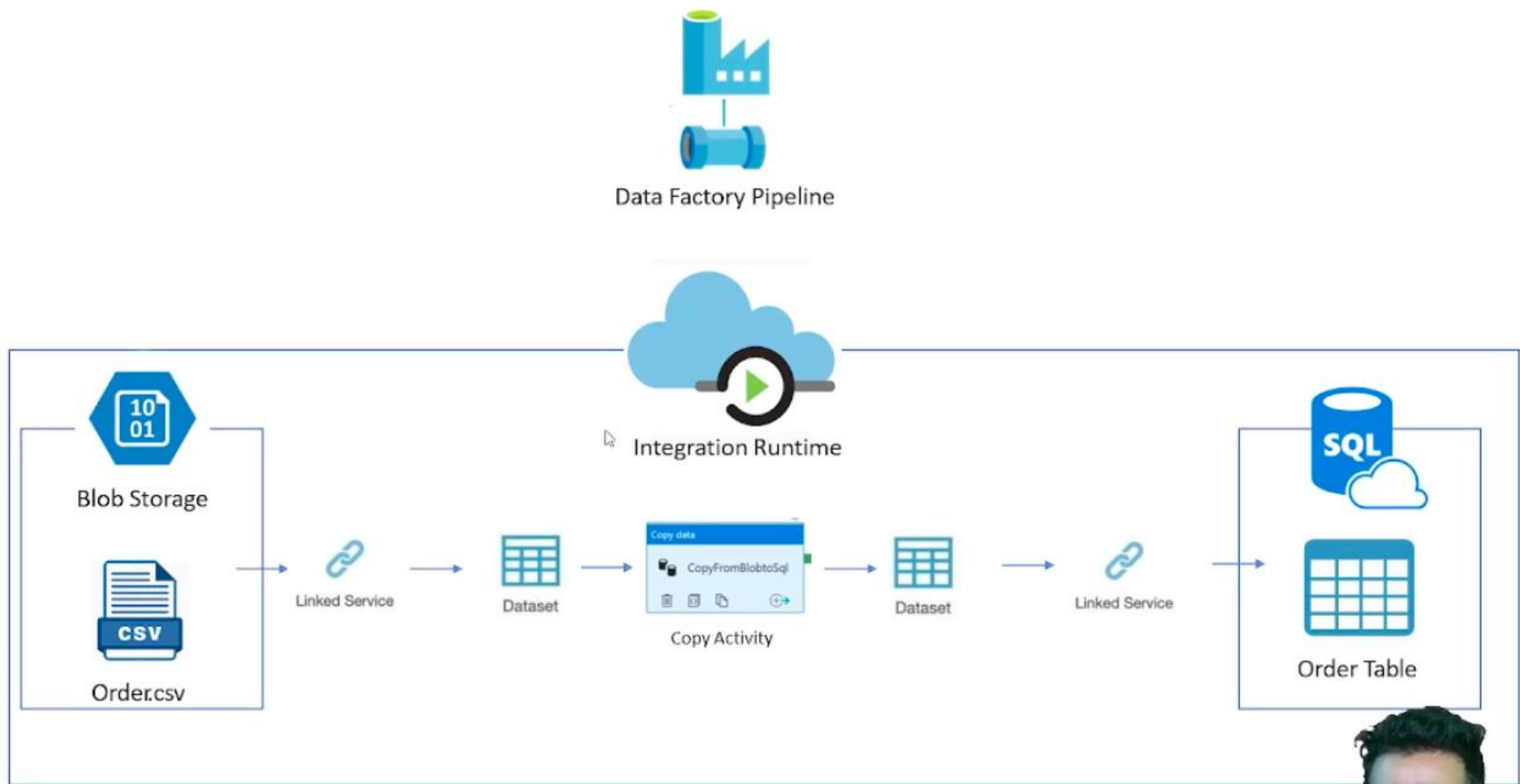


# Azure Data Factory Components

Slide 8 of 36







# Data Factory vs SSIS

## Azure Data Factory

Pipeline

Linked Service

Source

Sink

Activity

Data Flow

## SSIS

Package

Connection manager

Source

Destination

Control flow task

Data flow





Data Factory  
Pipeline

- Data Factories can contain one or more pipelines
- Logical group of Activities
- Manage Activities as a set
- One Pipeline can have one or more activities



## Azure Data Factory Activities

- Represents a processing step in the pipelines
- Actions to perform on data
  - Ingest data
  - Transform data
  - Store data
- Can be linked
  - Execute sequentially or
  - Run in parallel



## Activity types

01

### Data movement activities

Copy data amongst data stores located on-premises and in the cloud

Data stores – Blob storage, Cosmos DB, Amazon Redshift, Google BigQuery Hive, Maria DB...etc.



02

### Data transformation activities

Transform and enrich data

e.g. Hive, Pig, MapReduce, Spark or Databricks



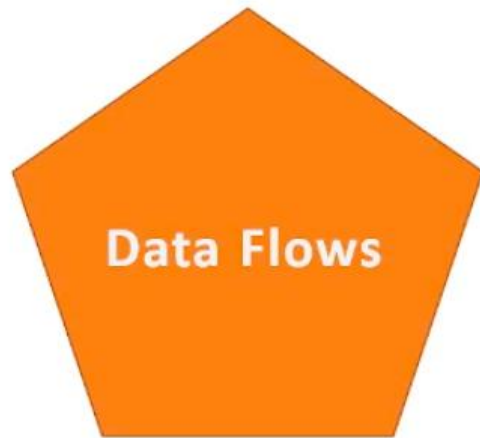
03

### Control activities

Control pipeline flow

e.g. ForEach, Web





- Data Flow is a new feature of Azure Data Factory (ADF) that allows you to develop graphical data transformation logic that can be executed as activities within ADF pipelines.
- Two types:
  - Mapping
  - Wrangling



# Integration Runtimes

- **Data Integration Capabilities**

- **Data Flow**

- **Data Movement**

- Format conversion, column mapping, serialization/deserialization etc.

- Provides the native compute to move data between cloud data stores in a secure, reliable, and high-performance manner.

- **Activity dispatch** (e.g. Databricks Notebook, HDInsight Hive, pig, spark activity, SP, ADL Analytics U-SQL activity)

- **SSIS Package execution**





## Dataset

- Simply point or reference the data
- Reference data used in anActivity
  - Files
  - Folders
  - Documents
  - Tables

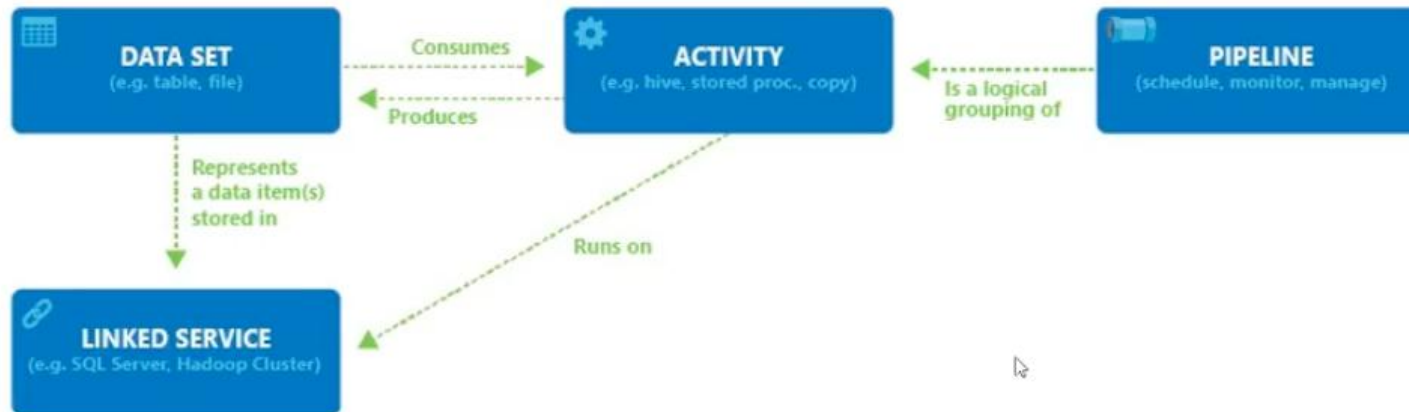




## Linked service

- Similar to connection string
- Represent the connection information to connect to external resources
  - Datastores like Azure SQL Server
  - Compute resource e.g. Spark Cluster

# ADF Components



# Integration Runtimes

- Provides fully managed, serverless compute infrastructure
  - You don't have to worry about infrastructure provision, software installation, patching, or capacity scaling.
  - Pay only for duration of actual use
- Bridges between the activity and linked service
  - Activity defines the action
  - Linked service define the location



# Integration Runtimes

Specify the infrastructure to run activities

## Azure Integration Runtime

Work on public networks

Responsible for data flows, data movements, and activity dispatches

## Self-hosted Integration Runtime

Work on public and private networks

Provide data movement and activity dispatch capabilities

Need to install on on-premises machine or a virtual machine inside private network

## SSIS Integration Runtime

Supports SSIS package execution

Works on public and private networks



# Integration Runtimes

IR type	Public network	Private network
Azure	Data Flow Data movement Activity dispatch	
Self-hosted	Data movement Activity dispatch	Data movement Activity dispatch
Azure-SSIS	SSIS package execution	SSIS package execution





# Integration Runtimes

- Default IR – AutoResolveIntegrationRuntime
- Create Azure IR
  - When you want to explicitly define the location of IR
  - Virtually group the activities executions on different IR for management purpose



ineuroninstance - Microsoft Azure

ineuroninstance - Azure Data Factory

ineuronblob - Microsoft Azure

adf.azure.com/en-us/authoring/pipeline/pipeline1?factory=%2Fsubscriptions%2F4642b277-3926-48dd-84fc-dffbc0cc6a92%2FresourceGroups%2Frg-df%2Fproviders%2FMicrosoft.DataFactory%2Ffactories%2Fineuroninstance

BookmarksRazorpay DashboardBitPaper - FeaturesHomeTL-WR740NPrivacy error(55) DSC Error in G...Manage excellatee...GitHub - Pierian Da...Your AccountLearning on Simpli...Reading list

Microsoft AzureineuroninstanceSearch

Home

Author

Monitor

Manage

Data Factory

Validate all

Publish all

Factory Resources

ssis p

Pipeline0

Dataset0

Data flows0

Power Query0

pipeline1

DelimitedText1

Activities

ssi

General

Execute SSIS package

Copy data

Copy d

Execute SSIS package

Execute SSIS packag

Parameters

Variable

New

New trigger

Name \*trigger1

Description

Type \*Schedule

Filter...

Schedule

Tumbling window

Storage events

Custom events

every 15 Minute(s)

Custom events

Specify an end date

Annotations

New

Start trigger

Start trigger on creation

OKCancel

Windows Taskbar

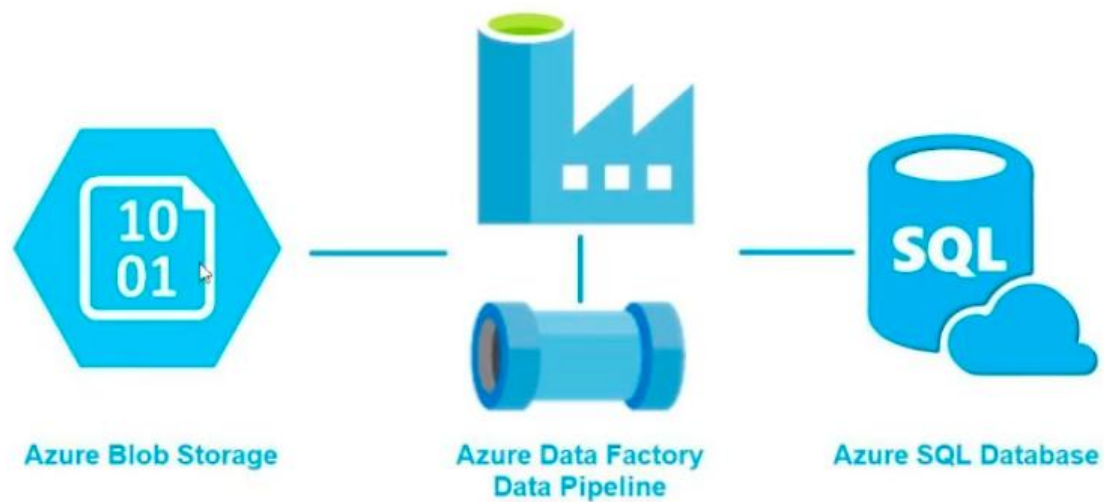
Neuron



- **Execute pipeline**
- **Many to many relationship b/w pipeline and trigger**
- **Three types of Trigger**
  - **Schedule Trigger** – Invoke pipeline on a wall-clock schedule
  - **Tumbling Window Trigger** – Operates on a periodic interval, also retain state
  - one-to-one relationship
  - Advance configuration options - Dependencies, delay, retry, concurrency
  - Properties - trigger().outputs.WindowStartTime/WindowEndTime
  - **Event-based Trigger** – trigger pipeline in response to an event
    - e.g. Arrival/deletion of file in Blob storage
    - Event trigger with Azure Event Grid Service
    - Properties – triggerBody().folderPath/fileName



## Demo: Copy Activity



# Data Flows



## Azure Data Factory Resources



Pipeline



Dataset



Linked Service



Integration Runtime



Data Flow

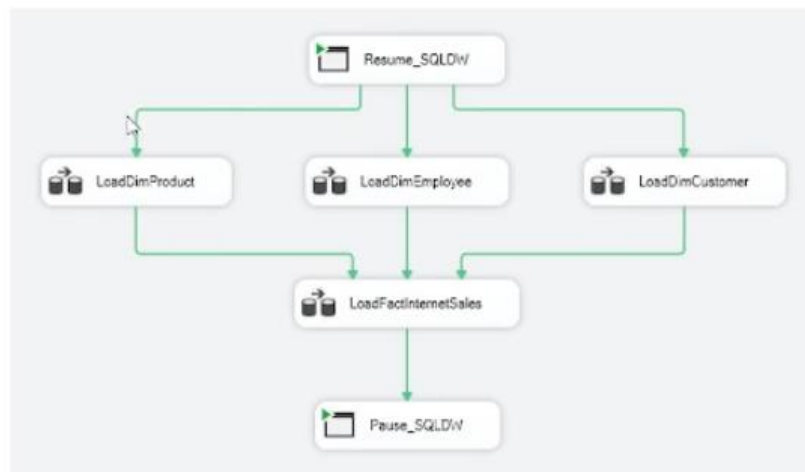
NEW

Allows you to develop graphical data transformation logic

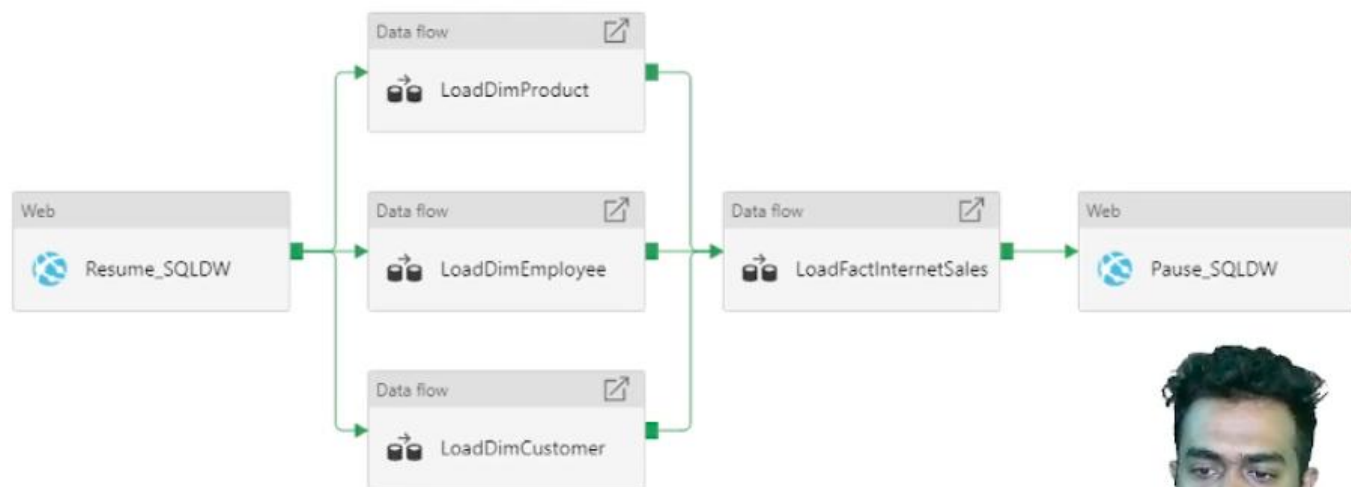




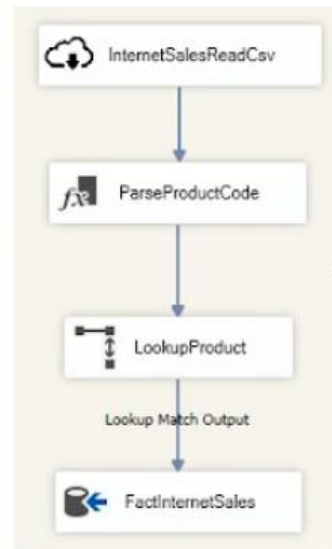
Example of the SSIS **Control Flow** tab for loading our data mart tables:



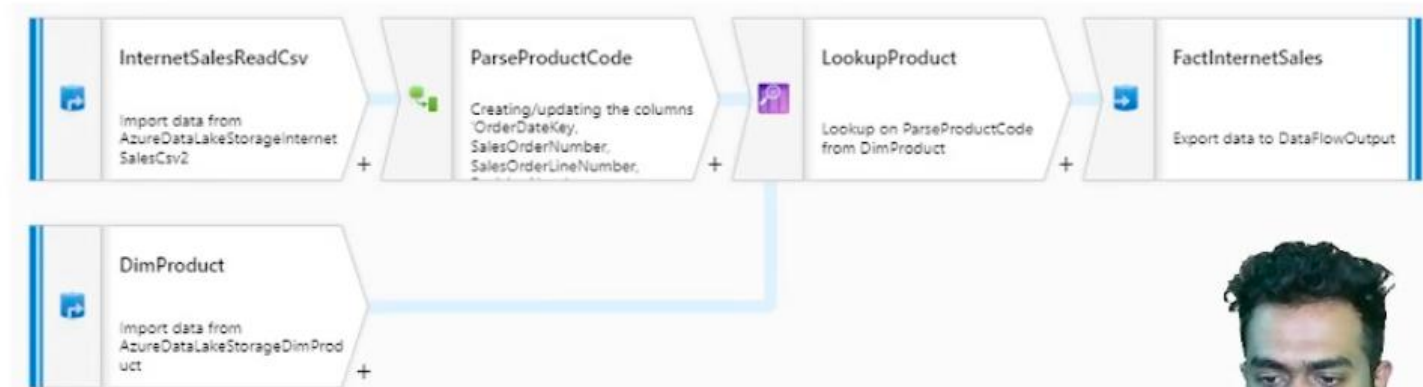
Example of the ADF Pipeline for loading our data mart tables:

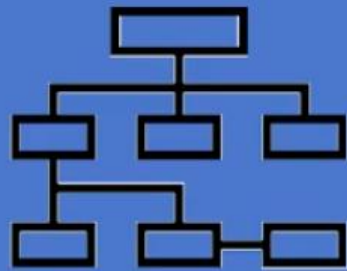


Example of SSIS **Data Flow** **tab** for loading the FactInternetSales table:



Example of ADF Mapping **Data Flows** for loading the FactInternetSales table:



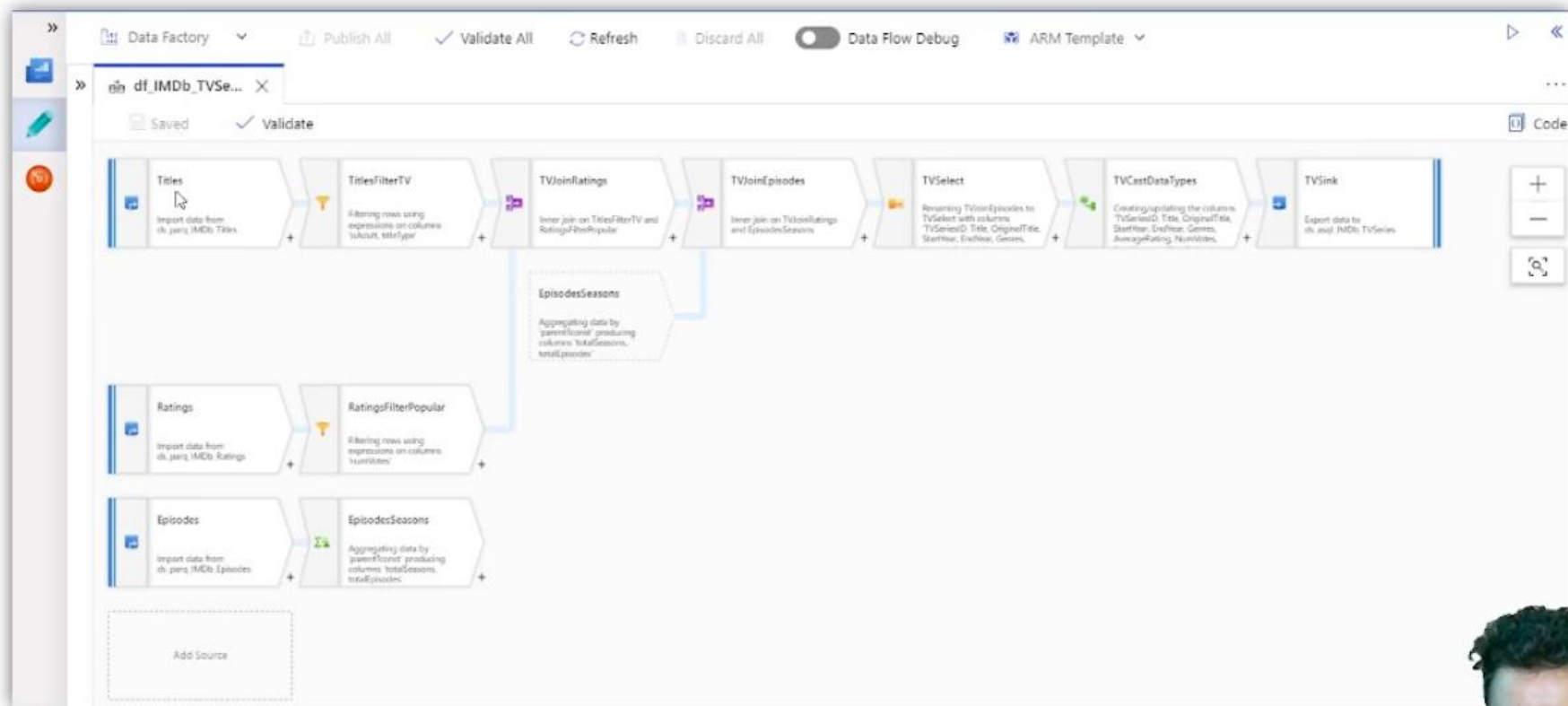


# Data Flow

**Mapping Data flow** – Transform Data  
(Known data and schema)

**Wrangling Data flow** – Prepare and explore  
data using power query (known or unknown  
datasets)

# Mapping Data Flows



## Data flows behind the scene



Behind the scene Data flow will execute on Azure Databricks using Spark



ADF internally handles all the code translation, spark optimization and execution of transformation