# Azure Data Factory v2

## Data Integration in the Cloud

Paul Andrew | Adatis

10/03/2018

**Microsoft®** MVP
Most Valuable
Professional

@MrPaulAndrew

Microsoft Partner

Gold Data Analytics
Gold Data Platform
Gold Cloud Platform

Microsoft

adatis

# All Talk Content on GitHub

https://github.com/mrpaulandrew

**SQLSaturdays**

Content and slides from various SQL Saturday events.

● PLpgSQL    PASS SQLSATURDAY

{Location}-{Month}-{Year}

# Agenda

**Data Factory Recap**

Concepts
Components

**ADFv2**

Features Update
The Integration
Runtime

**Data Factory Extensibility**

SSIS, Functions,
Custom Activities

**Conclusions**

Design Patterns
ETL/ELT in Azure

# Agenda

**Data Factory Recap**

Concepts
Components

**ADFv2**

Features Update
The Integration Runtime
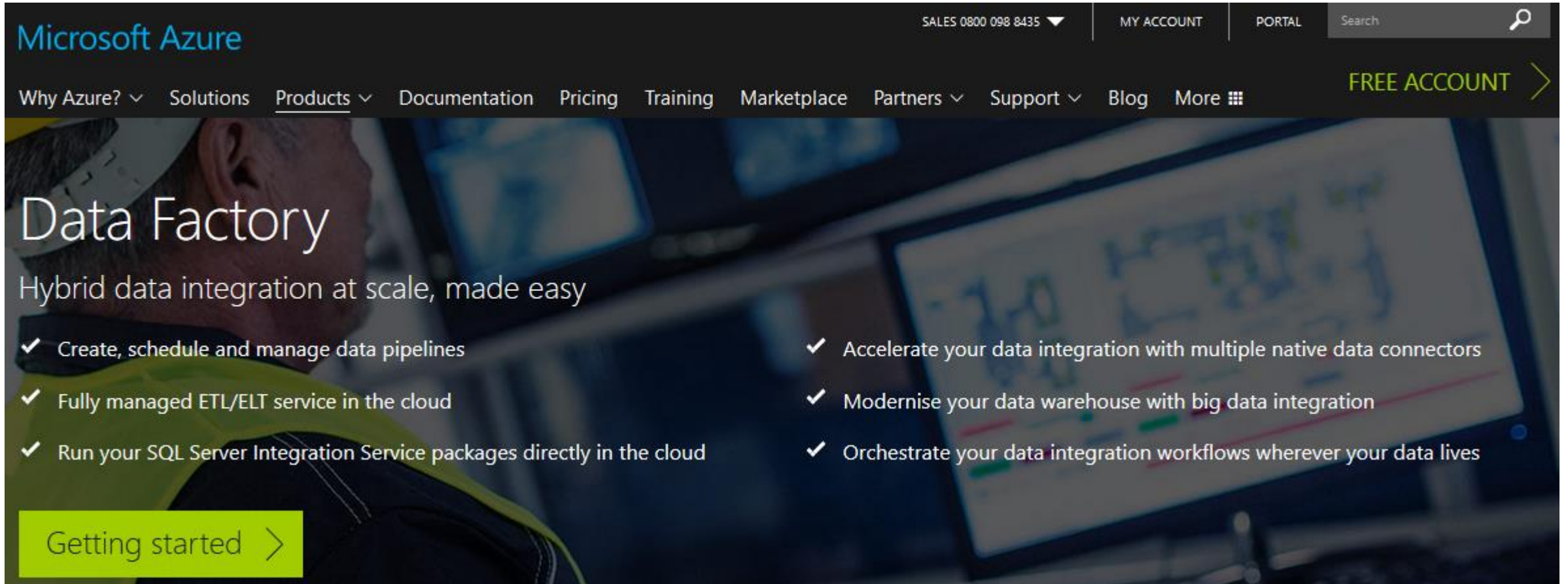
**Data Factory Extensibility**

SSIS, Functions, Custom Activities

**Conclusions**

Design Patterns
ETL/ELT in Azure

# What is Azure Data Factory?

https://azure.microsoft.com/en-gb/services/data-factory/

# What is Azure Data Factory?

Copy Data

Transform

**SQL**

Performs copy activities using it's own scalable compute

**Azure Data Factory**

JSON

Triggers other components for heavy lifting & intensive transformation

# What does Azure Data Factory do?



Sources

Copy Data

Transform Activities

Destinations

# Data Factory Components

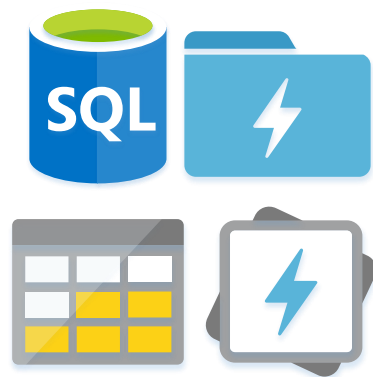Copy Data

Transform

**1** **Linked Services** – How do I connect?
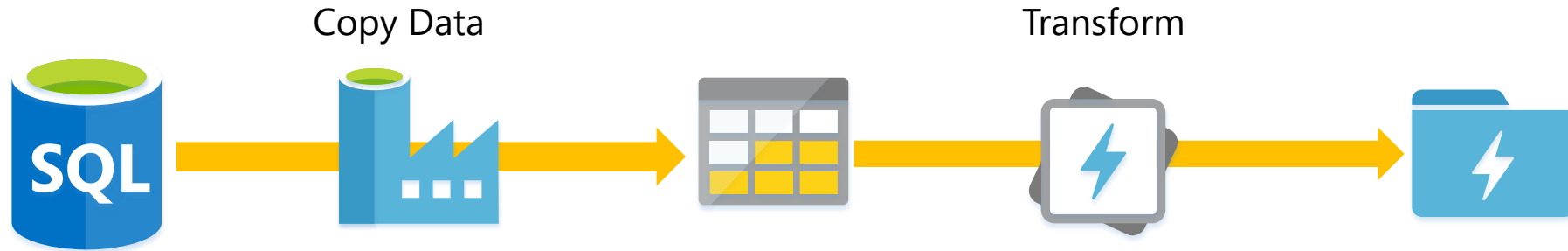
Like the SSIS Connection Manager!

SQLDBLinkedService

ConnectionString: *Server=MyServer;Database=myDataBase*
*UserName: "MrPaulAndrew"*
*Password: ***************

# Data Factory Components

Copy Data

Transform

SQL

1  **Linked Services**

2  **Data Sets** – What slices/partitions does my data have?

SQL  dbo.DimCustomer

/RAW/Orders/2018/01/01/Orders.csv

# Data Factory Components



Copy Data

Transform

1. **Linked Services**

2. **Data Sets**

3. **Activities** – What do we want to happen?
With what conditions?

**U-SQL Activity**

Script: *wasb//:myscripts/ProcessOrders.usql*
AUs: *5 units*
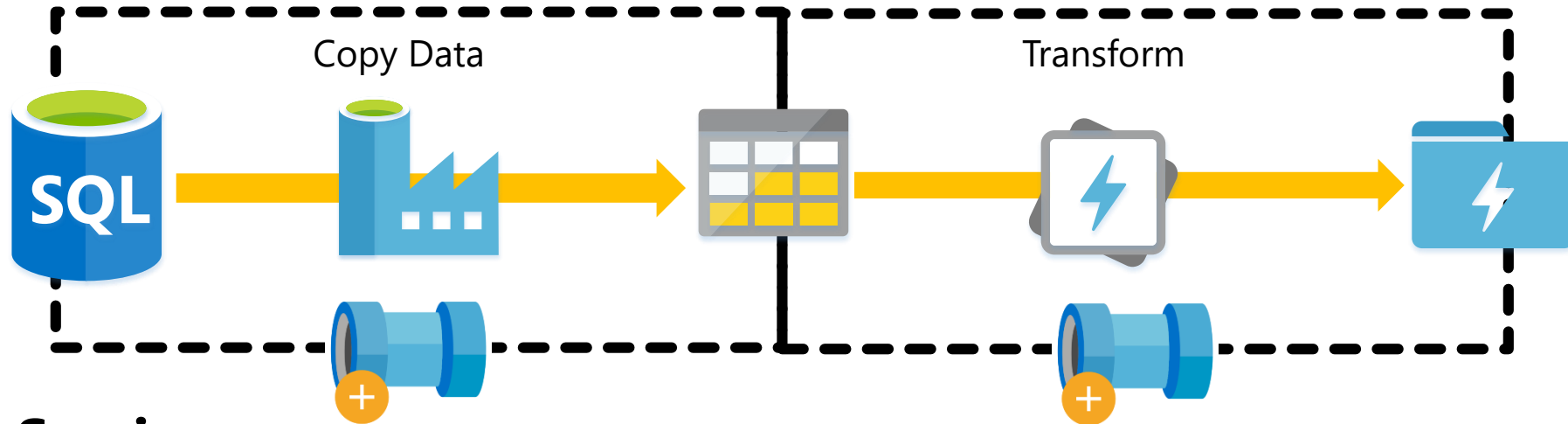Priority: *1000*
Parameters: *@Output = "RAW/Orders/..."*

# Data Factory Components



Copy Data        Transform

**1**   **Linked Services**

**2**   **Data Sets**

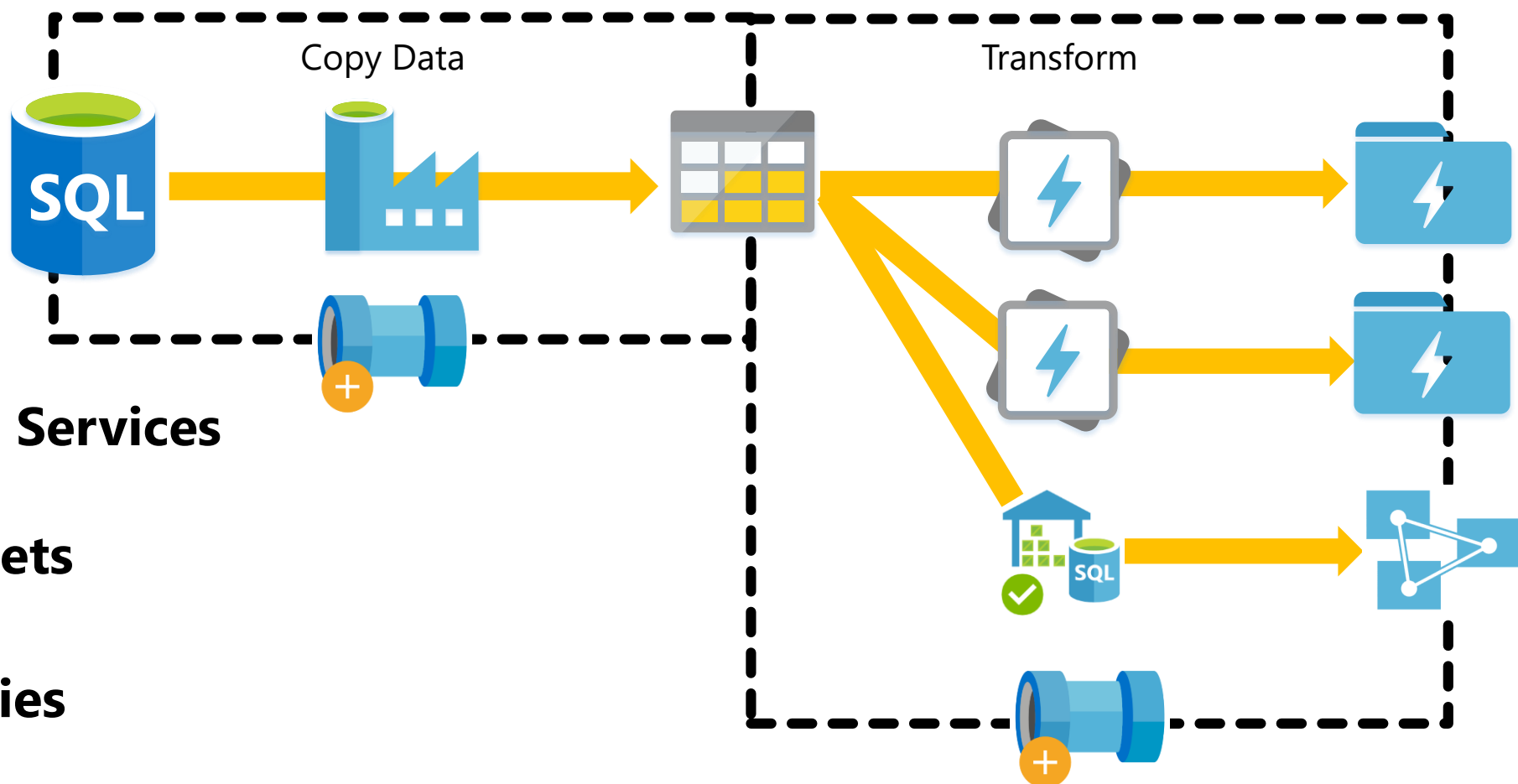**3**   **Activities**

**4**   **Pipelines** – What groups of work do I want to do?

# Data Factory Components



Copy Data

Transform

1 **Linked Services**

2 **Data Sets**

3 **Activities**

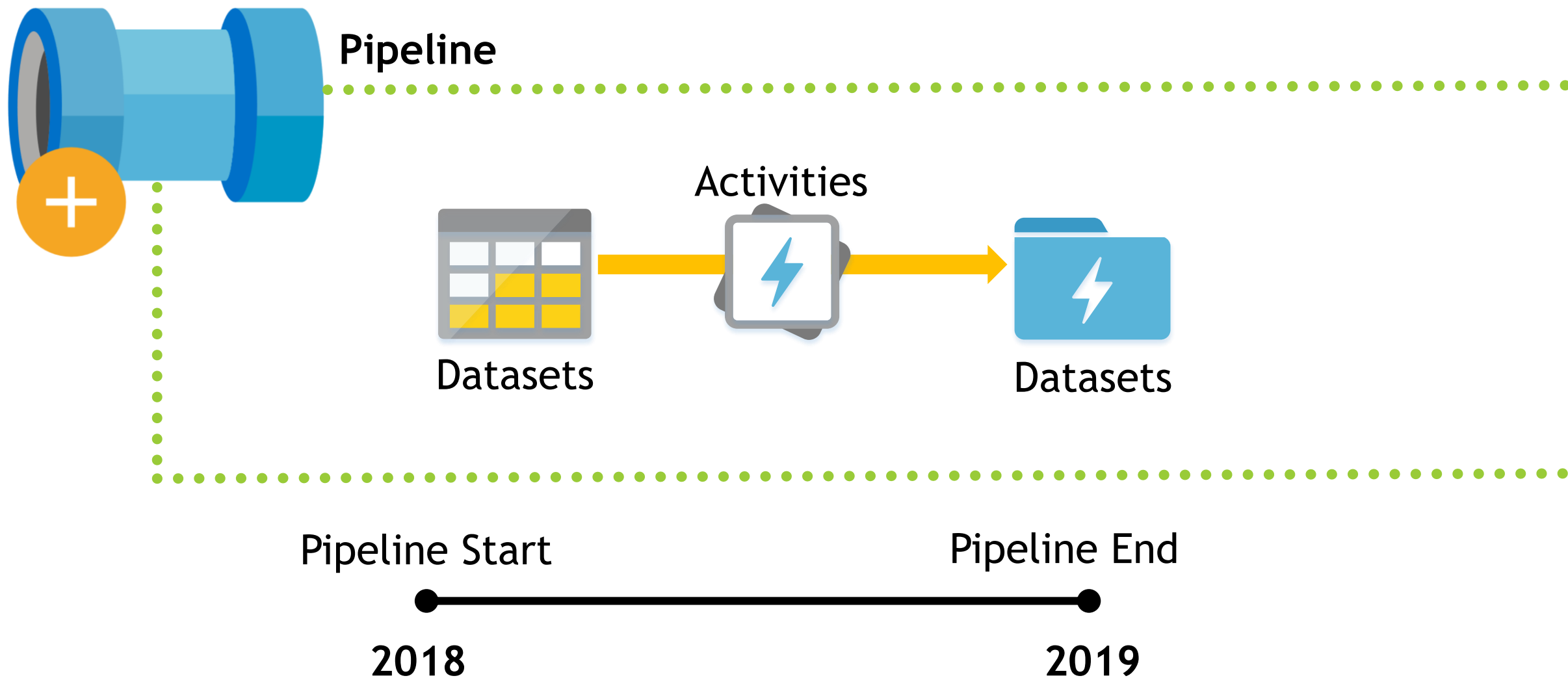4 **Pipelines** – What groups of work do I want to do?

# Azure Data Factory Concepts

**Time Slices** – triggering an activity execution.

**Pipeline**

Activities

Datasets

Datasets

Pipeline Start

Pipeline End

**2018**

**2019**

# Azure Data Factory Concepts Continued

**Time Slices** – triggering an activity execution.



**Pipeline**

Activities

Datasets

Datasets

Interval: *Month*
Frequency: *1*

2018

2019

# Azure Data Factory Concepts Continued

**Time Slices** – triggering an activity execution.



**Pipeline**

Activities

Datasets

Datasets

Loading

2018

2019

Time Slice Problems…

# Agenda

| Data Factory Recap | ADFv2 | Data Factory Extensibility | Conclusions |
|---|---|---|---|
| Concepts Components | Features Update The Integration Runtime | SSIS, Functions, Custom Activities | Design Patterns ETL/ELT in Azure |

# Data Factory Issues & Limitations

- **Time Slices** – Complex, difficult to change & provision
- **Pricing** – High and low frequency
- **Control Flow** – Not conditional. Pass or fail
- **Developer Tools** – VS or Portal JSON templates
- **Hard Coded Pipelines** – No dynamic values
- **C# Coding** – Often required for anything complex
- **Monitoring** – Times slice bound, focused on datasets
- **Connectivity** – Limited to Microsoft supported linked services

# So What's Changed?



| | | |
|---|---|---|
| ✓ | Data Movement (Copy) | ✓ |
| ✓ | Activities | ✓ |
| ✓ | Pipelines | ✓ |
| ✓ | Datasets | ✓ |
| ✓ | Time Slices/Tumbling Windows | ✓ |
| ✗ | Event Triggers | ✓ |
| ✗ | Recurring Schedules | ✓ |
| ✗ | Parameters | ✓ |
| ✗ | Expressions | ✓ |
| ✗ | Conditional Logic | ✓ |
| ✗ | Use of SSIS Packages | ✓ |
| ✗ | Graphic Developer Canvas | ✓ |
| ✗ | Drilldown Monitoring | ✓ |

V1

V2
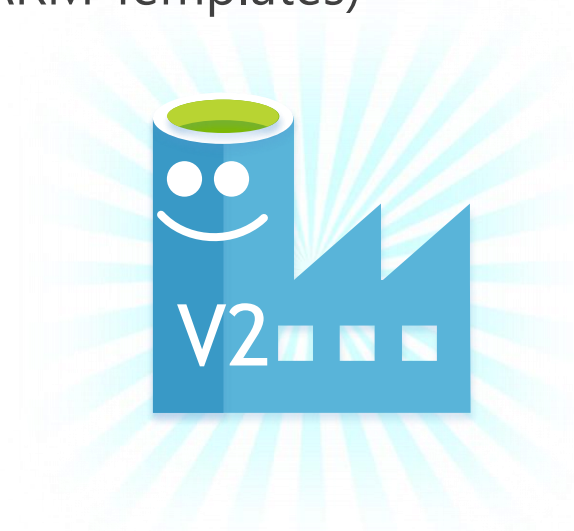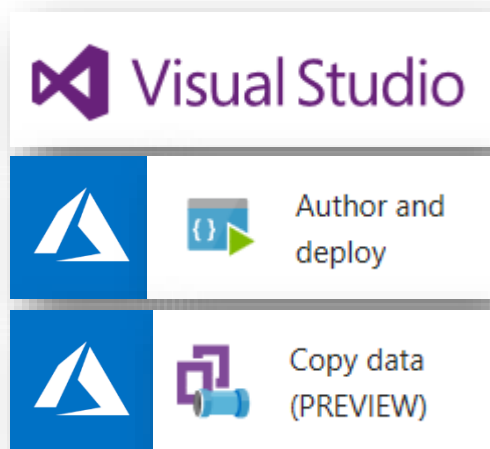
# Developer Tools

- JSON Templates
- Data Factory Wizard
- Reverse Engineer From Azure
- Deployment Wizard

- No JSON Templates
- No Deployment Wizard
- Data Factory Wizard
- Reverse Engineer From Azure (via ARM Templates)



Visual Studio Code editor window showing:
```
ADFPipeline.json - Untitled (Workspace) - Visual Studio Code
File  Edit  Selection  View  Go  Debug  Tasks  Help
{} ADFPipeline.json
1    {
2        "PipelineName": "Why???"
3    }
```

Visual Studio — Author and deploy — Copy data (PREVIEW)

V1

Visual Studio (crossed out) — Author & Monitor

V2

Microsoft .NET

PowerShell

Only Visual Studio 2015
http://bit.ly/2tsyD90
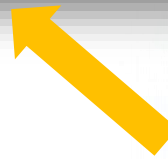
# Monitoring



**V1**

**V2**

ADFv2DemoFactory01 | Monitor **Pipeline Runs** ⌄
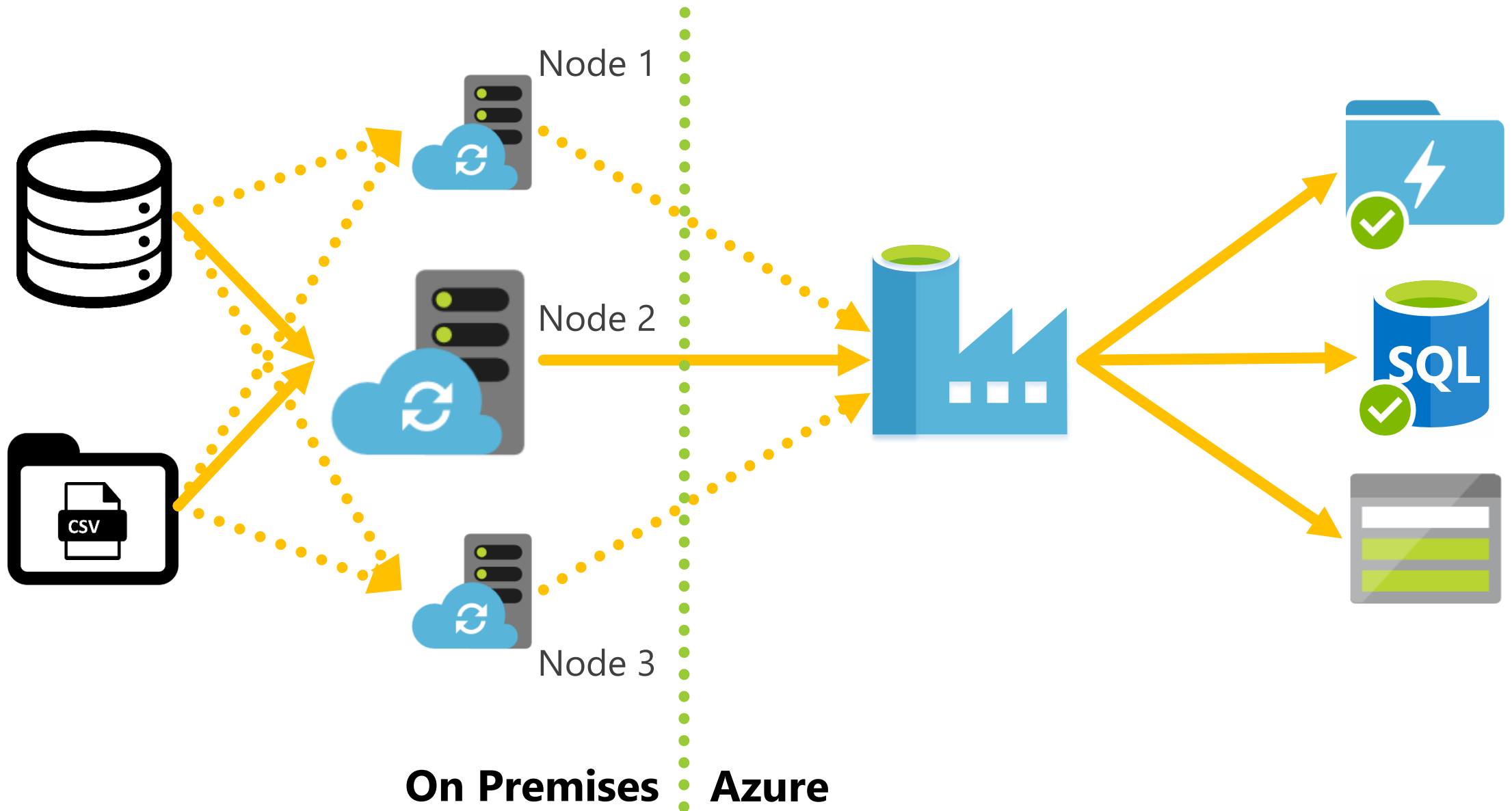
⟳ Refresh

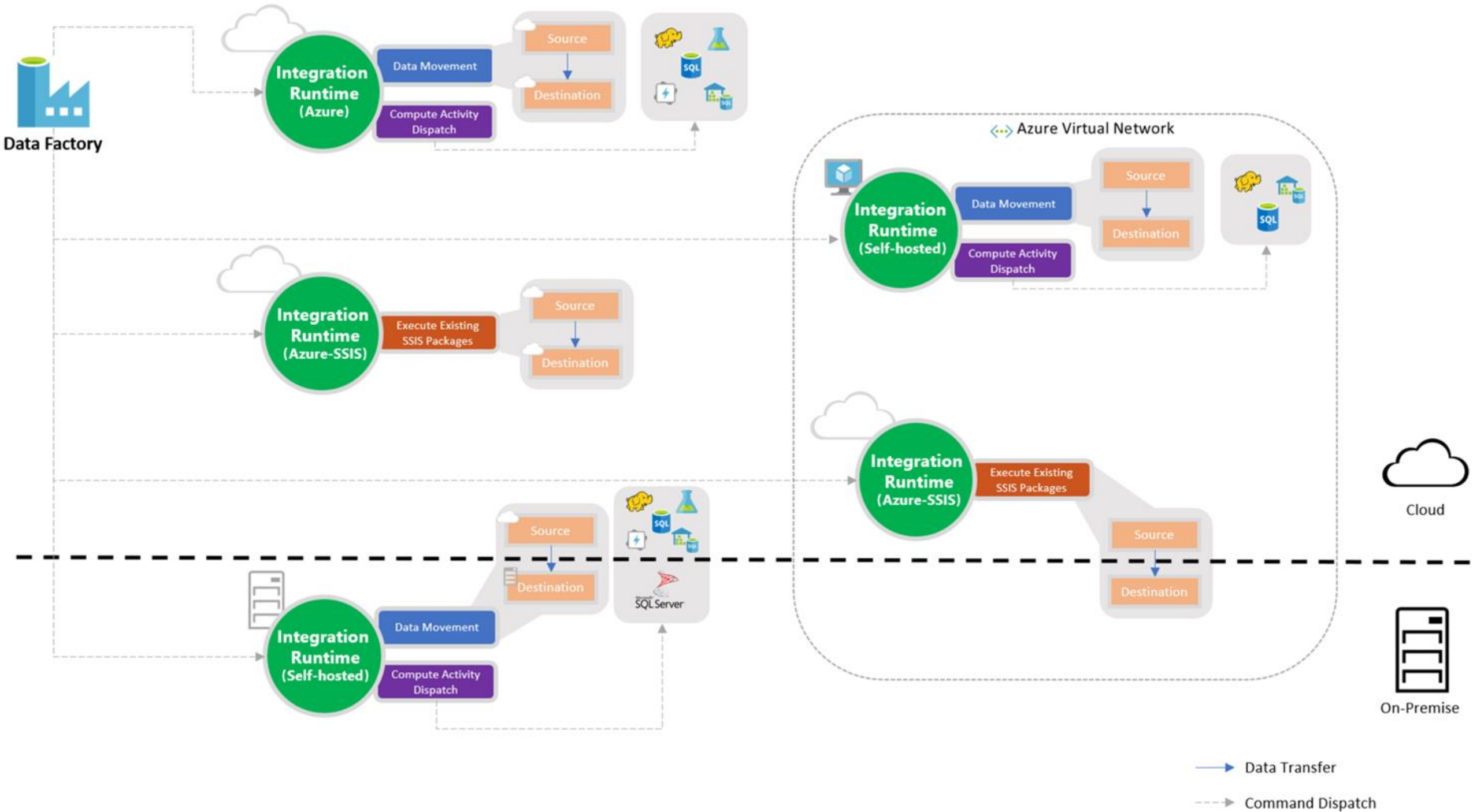📅 **Last 24 Hours** 01/14/2018 1:35 PM - 01/15/2018 1:35 PM ⌄          🌐 **Time Zone** (UTC+00:00) London ⌄

**All**   Succeeded   In Progress   Failed   Cancelled

| Pipeline Name ▽ | Actions | Run Start ⇅ | Duration | Triggered By | Status | Parameters | Error |
|---|---|---|---|---|---|---|---|
| RunSSISPackage | ▶ | 01/15/2018, 1:36:42 PM | 00:00:11 | Manual trigger | ✔ Succeeded | | |

Data factory

| RESOURCE EXPLORER ✕ | ⏱ PaulsFunFactoryV1 |

▶ ⏸ ⬛   Start time (UTC): 01/08/2018 01:32 pm   End time (UTC): 01/16/2018 0...

◢ **Data Factories**
  ◢ PaulsFunFactoryV1
    ◢ Pipelines
      FileCleaning
      UploadFileToADLStore
    ◢ Datasets
      FakeOrdersClean
      FakeOrdersLanding
      FakeOrdersSourceFile
    ◢ Linked services
      BatchCompute
      BlobStore
      DataLakeStore
      LaptopGateway
      USQLEngine
    ◢ Gateways
      PaulsLappy
      jhjlhkj

FakeOrdersSourceFile          **Active**
FILESHARE          **UploadFileToADLStore**
PIPELINE
**1** activities          FakeOrdersLanding
AZURE DATA LAKE...

**FileCleaning**
PIPELINE
**1** activities          FakeOrdersClean
FREQ: MONTH
INTVL: 1     AZURE DATA LAKE...

Last edited by Paul Andrew 4 ...

**ACTIVITY WINDOWS**

⟳ ⟳ 📋 ▽   No filter applied.                    Last refreshed a few seconds ...

| Pipeline ▽ | Activity ▽ | Window Star... | Window End ▽ | Status ▽ | Type ▽ | Last Attempt... | Last Attempt... | Duration | Retry At... |
|---|---|---|---|---|---|---|---|---|---|

There are currently no activity windows to display.

# The Integration Runtime *(AKA The Data Management Gateway)*

Node 1

Node 2

Node 3

**On Premises** **Azure**

# Agenda

**Data Factory Recap**

Concepts
Components

**ADFv2**

Features Update
The Integration
Runtime

**Data Factory Extensibility**

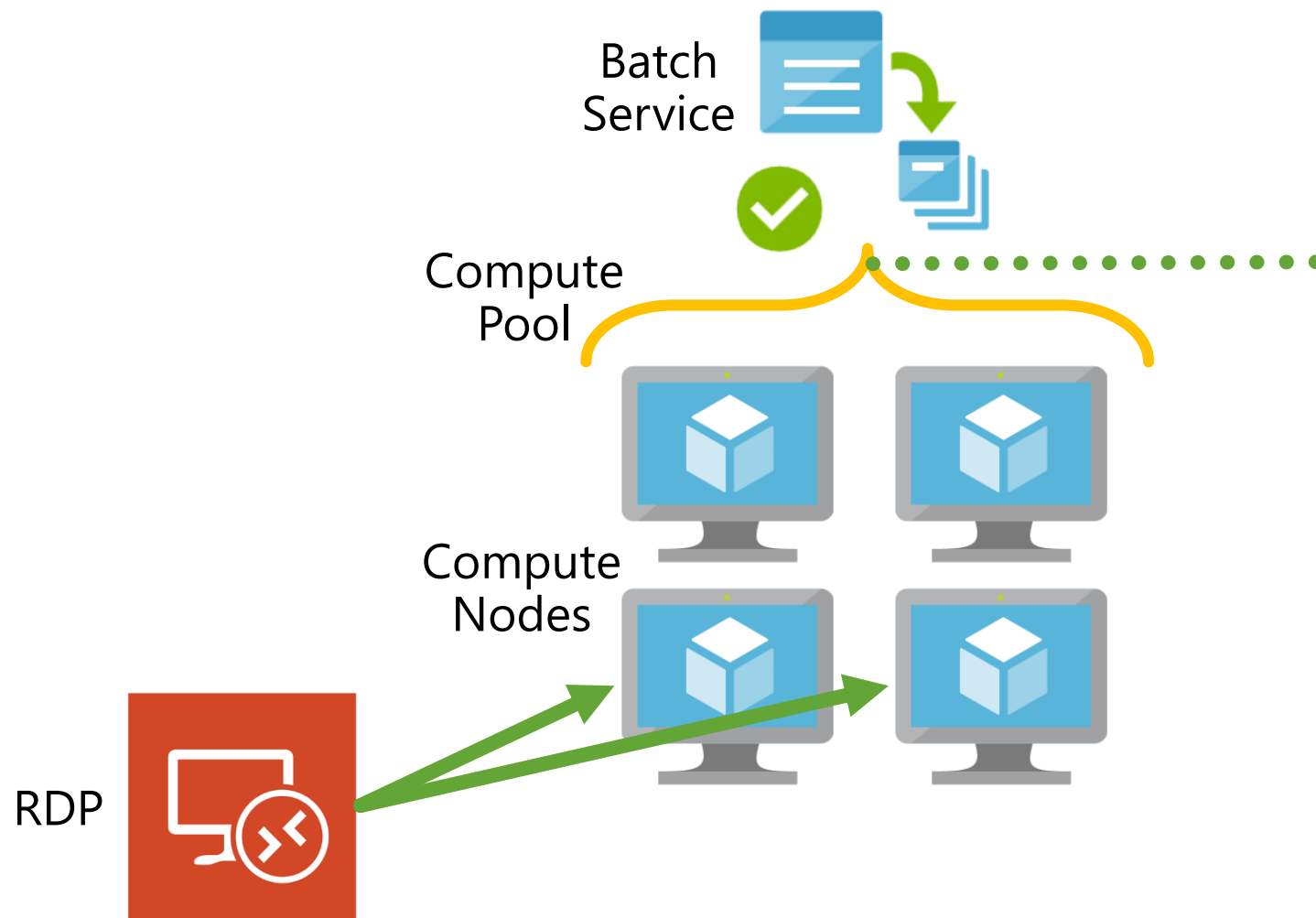SSIS, Functions,
Custom Activities

**Conclusions**

Design Patterns
ETL/ELT in Azure

# ADF Extensibility

**1** **Custom Activities** – C# classes called by the Azure Batch Service
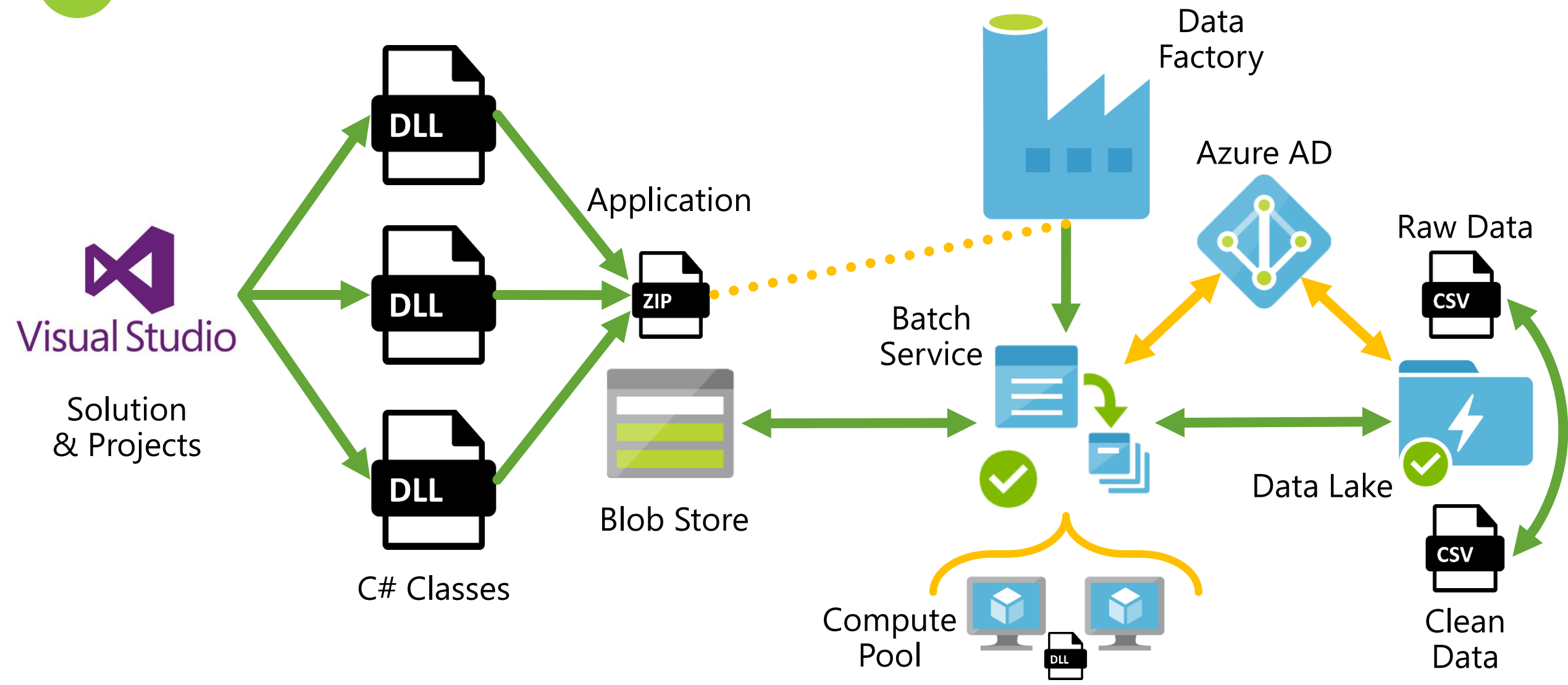
Batch Service

Compute Pool

Compute Nodes

RDP

VM node size set per compute pool:

| A1 Standard ★ | A2 Standard ★ | A3 Standard ★ |
|---|---|---|
| 1 Cores | 2 Cores | 4 Cores |
| 1.8 GB | 3.5 GB | 7 GB |
| 1 TB OS disk size | 1 TB OS disk size | 1 TB OS disk size |
| 70 GB Resource disk size | 135 GB Resource disk size | 285 GB Resource disk size |
| 2 Max data disk | 4 Max data disk | 8 Max data disk |
| Unable to display pricing | Unable to display pricing | Unable to display pricing |

- 1 compute node = 1 virtual machine.

- 1 job per compute node.

- Max of 4 tasks per node.

- OS on D drive, not C.

- Special environment variables.

# ADF Extensibility Continued

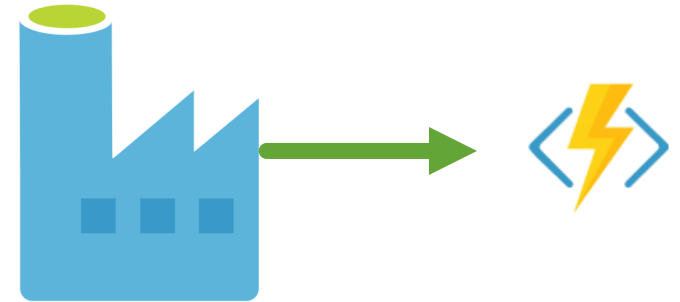**1** **Custom Activities** – C# classes called by the Azure Batch Service



Data Factory

Azure AD

Raw Data

Application

ZIP

DLL

DLL

DLL

Visual Studio

Solution & Projects

C# Classes

Blob Store

Batch Service

Compute Pool

Data Lake

Clean Data

# ADF Extensibility Continued

**1** **Custom Activities**

**2** **Rest API Calls** – Web Activities Calling Azure Functions
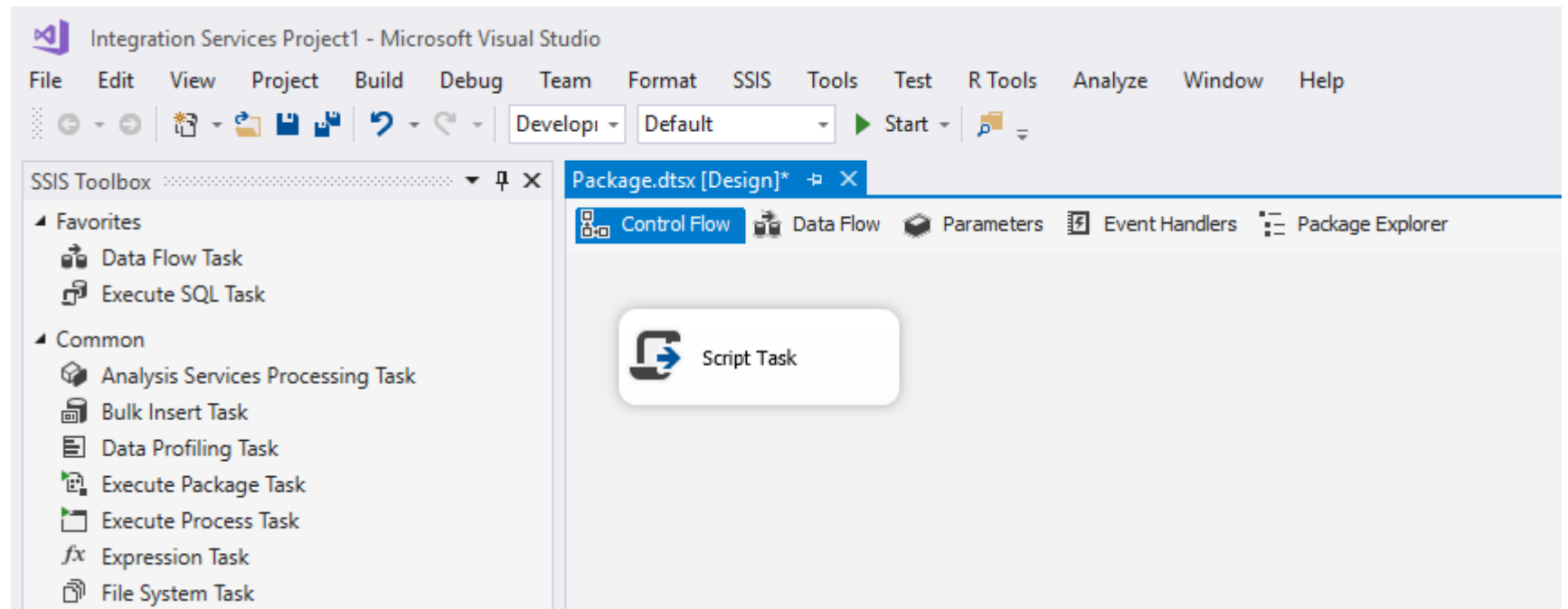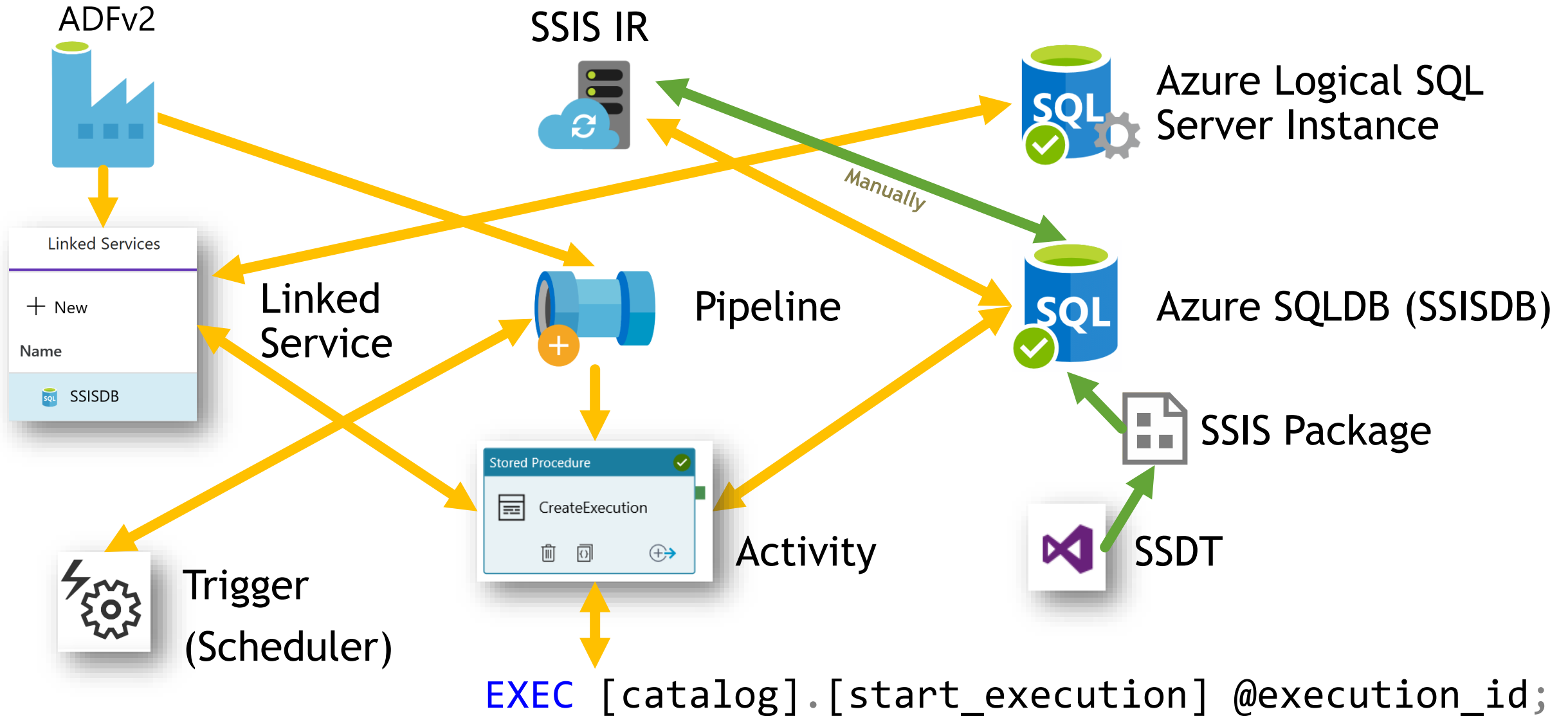
# ADF Extensibility Continued

**1** **Custom Activities**

**2** **Rest API Calls**

**3** **SSIS** – Packages with Control Flows and Data Flows

# What do we need to schedule an SSIS Package in Azure?



ADFv2

SSIS IR

Azure Logical SQL Server Instance

Linked Services

+ New

Name

SSISDB

Linked Service

Pipeline

Manually

Azure SQLDB (SSISDB)

Stored Procedure ✓

CreateExecution

Activity

SSIS Package

Trigger (Scheduler)

SSDT

```
EXEC [catalog].[start_execution] @execution_id;
```

# Demo: ADFV2

# Agenda

**Data Factory Recap**

Concepts
Components

**ADFv2**

Features Update
The Integration Runtime

**Data Factory Extensibility**

SSIS, Functions, Custom Activities

**Conclusions**

Design Patterns
ETL/ELT in Azure
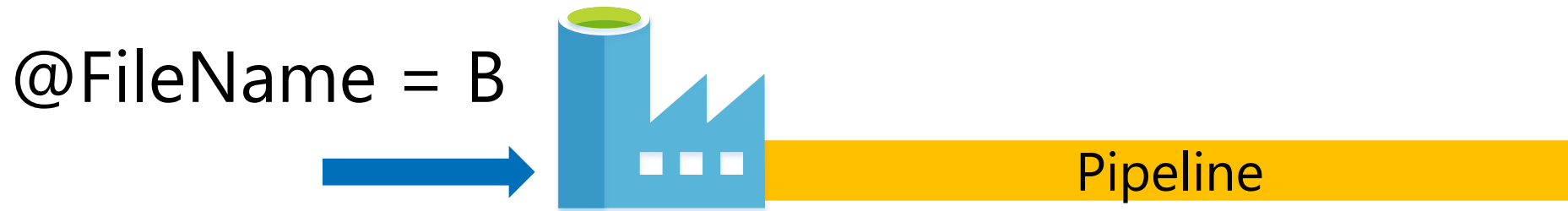
**Maybe, limited use.**

**Yes, definitely.**

# Dynamic Pipelines using Parameters

@FileName = B
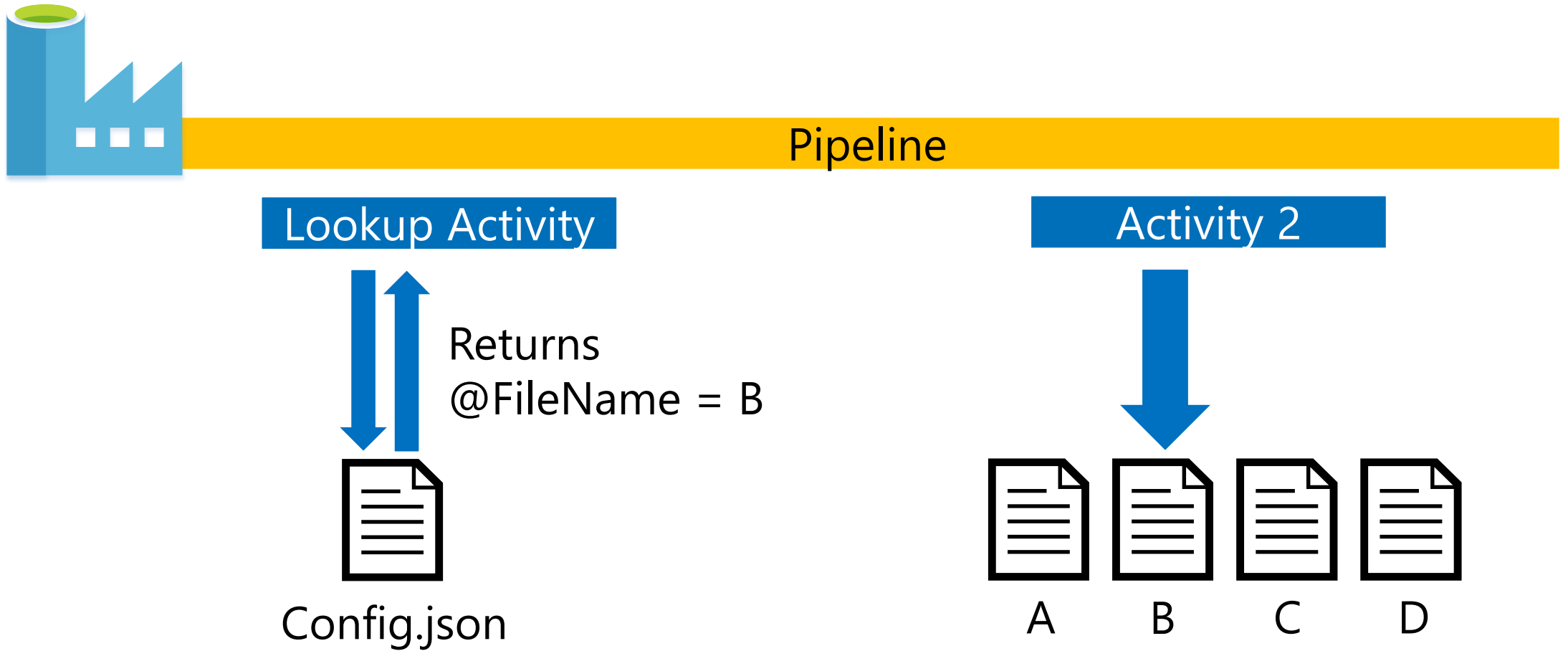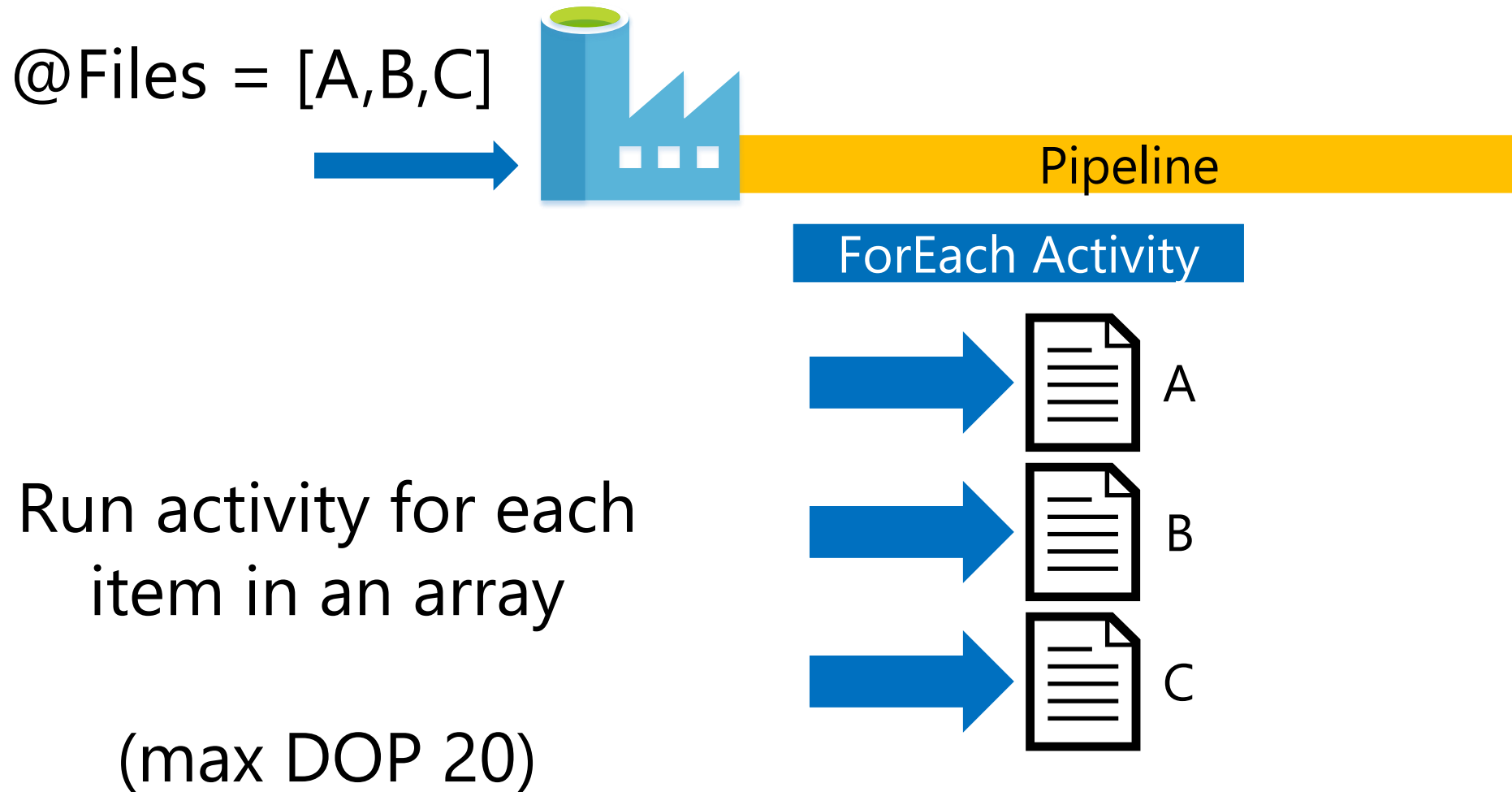
Pipeline

Activity 1

Dynamically change
parameters based
on inputs

A    B    C    D

V2

# Dynamic Pipelines using Lookup Activity

# ForEach Pipelines

@Files = [A,B,C]

Pipeline

ForEach Activity

A

B

C

Run activity for each item in an array

(max DOP 20)

# Execute Pipeline Activity

# Pipeline Triggers



Manual Trigger

Time Slices

Wall Clock

Pipeline 1

# SSIS Integration Runtimes



Pipeline

Stored Proc

SSIS Catalog

SSIS IR

The SSIS IR with Azure V-Net Access

# Event-Triggered Loading (The New ETL?)

**1** File added to Network Folder

**3** Logic App triggers Data Factory event. Provides filename as parameter

Pipeline 1

**2** Logic App watches folder, triggering at the next polling interval

**4** Pipeline runs, processing file as soon as it is available. Only one ADF Pipeline required

# Reusable Pipelines

# Batch Queuing



**2** Pipeline Triggered

Pipeline A

Lookup Activity

ForEach Activity

**3**

Lookup reads list of files into parameter

**1** List of files to process written to json doc

Pipeline B

Pipeline B

Pipeline B

**4** Child pipeline triggered for each file in lookup parameter. Up to 20 concurrent executions

# Thanks for Listening

## Paul Andrew

@MrPaulAndrew



**Blog:**    http://mrpaulandrew.com
**Email:**   paul@mrpaulandrew.com