# Complex Orchestration

## With Dynamic Data Factory Pipelines

Paul Andrew | Principal Consultant & Solution Architect

@MrPaulAndrew    In/MrPaulAndrew

https://github.com/mrpaulandrew

**CommunityEvents**

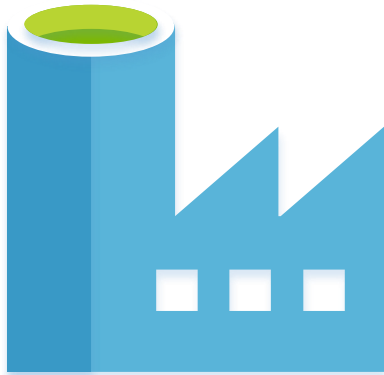Demo code, content and slides from various community events.

● C++

{Event/Location}-{Month}-{Year}
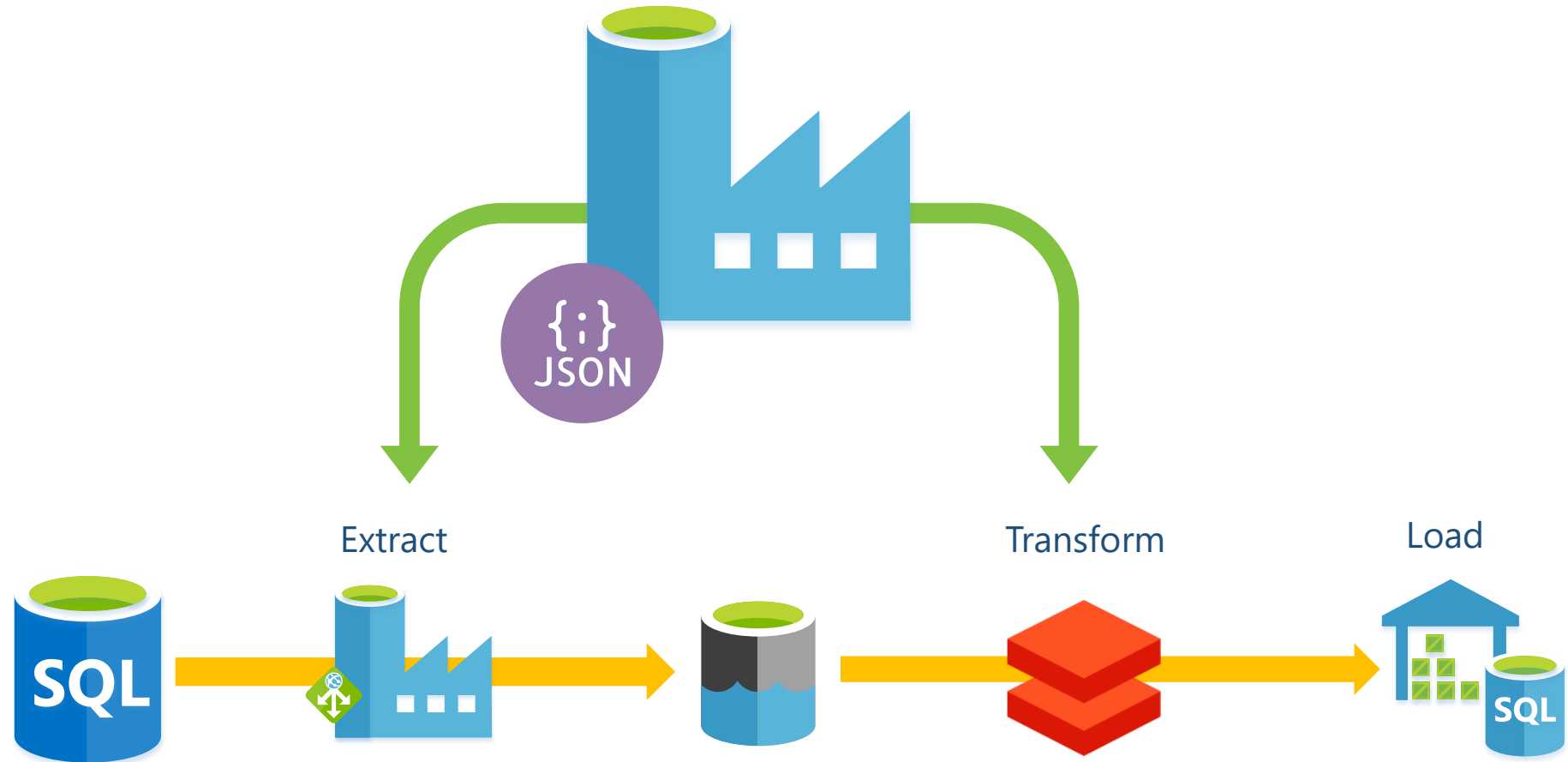
# Session Agenda

- Data Factory – A Quick Overview

- Dynamic Pipelines

- Extending Data Factory
  - Web Activities
  - Custom Activities

- True Scale Out Execution
  - SSIS Integration Runtime

- Data Factory – In Production
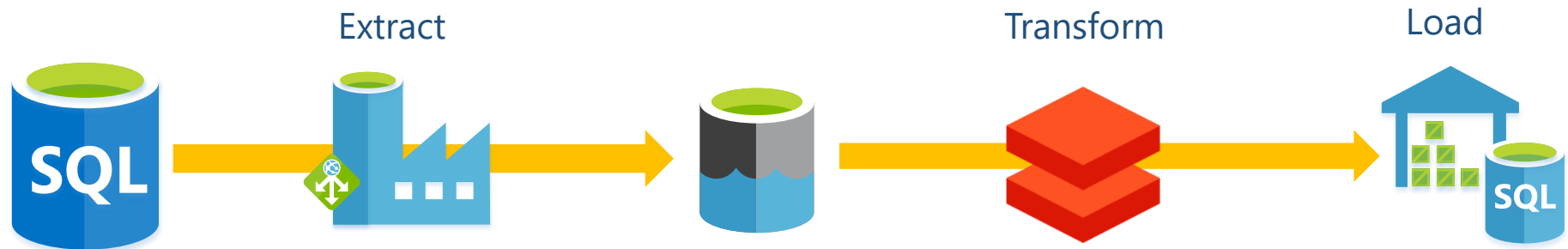  - Bootstrapping
  - DevOps

# What is Azure Data Factory?



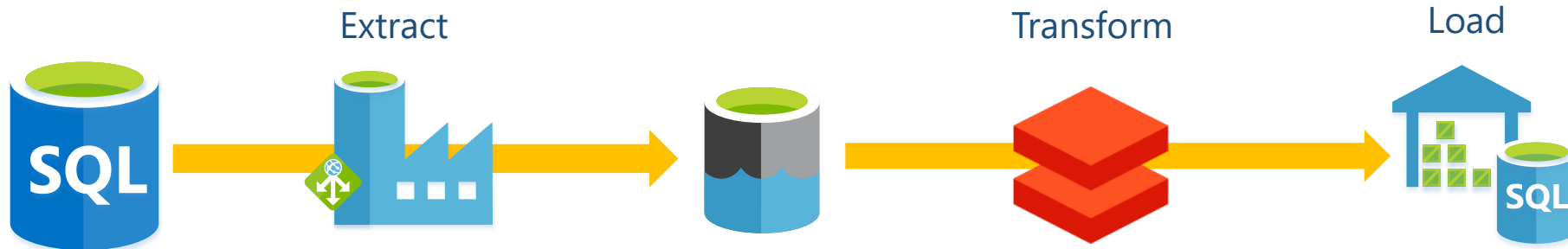Extract          Transform          Load

# What is Azure Data Factory?

Extract

Transform

Load

# Data Factory Components

Extract        Transform        Load

1 **Linked Services** ✓
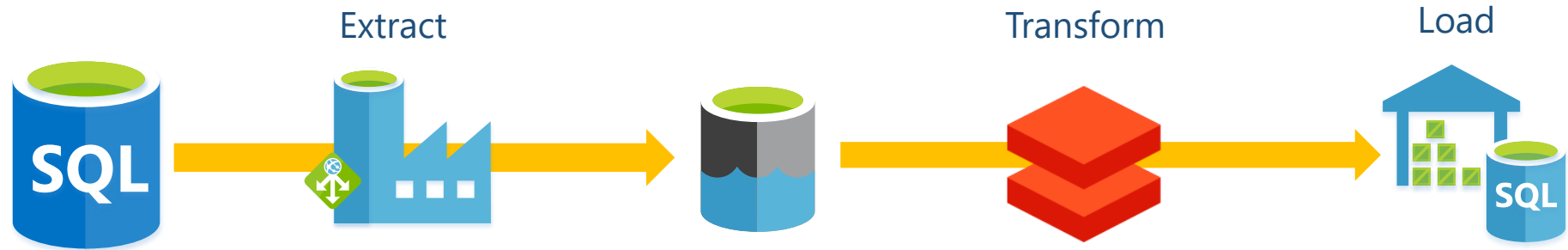2 **Data Sets** ✓
3 **Activities** ✓
4 **Pipelines** ✓
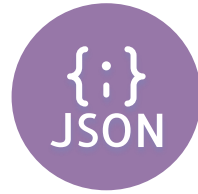5 **Triggers** ✗

{ ; } JSON

```json
{
    "name": "GenericSQLDB",
    "type": "Microsoft.DataFactory/factories/linkedservices",
    "properties": {
        "parameters": {
            "ServerInstance": {
                "type": "String"
            },
            "DatabaseName": {
                "type": "String"
            },
            "SQLUser": {
                "type": "String"
            },
            "SQLPassword": {
                "type": "String"
            }
        },
        "type": "AzureSqlDatabase",
        "typeProperties": {
            "connectionString": "Integrated Security=False;Encrypt=True;ConnectionTimeout=30;
            Data Source=@{linkedService().ServerInstance};
            InitialCatalog=@{linkedService().DatabaseName};
            UserID=@{linkedService().SQLUser};
            Password=@{linkedService().SQLPassword}"
        }
    }
}
```

# Data Factory Components

Extract

Transform

Load

**SQL**

**SQL**

| 1 | **Linked Services** ✓ |
| 2 | **Data Sets** ✓ |
| 3 | **Activities** ✓ |
| 4 | **Pipelines** ✓ |
| 5 | **Triggers** ✗ |

{;}
JSON

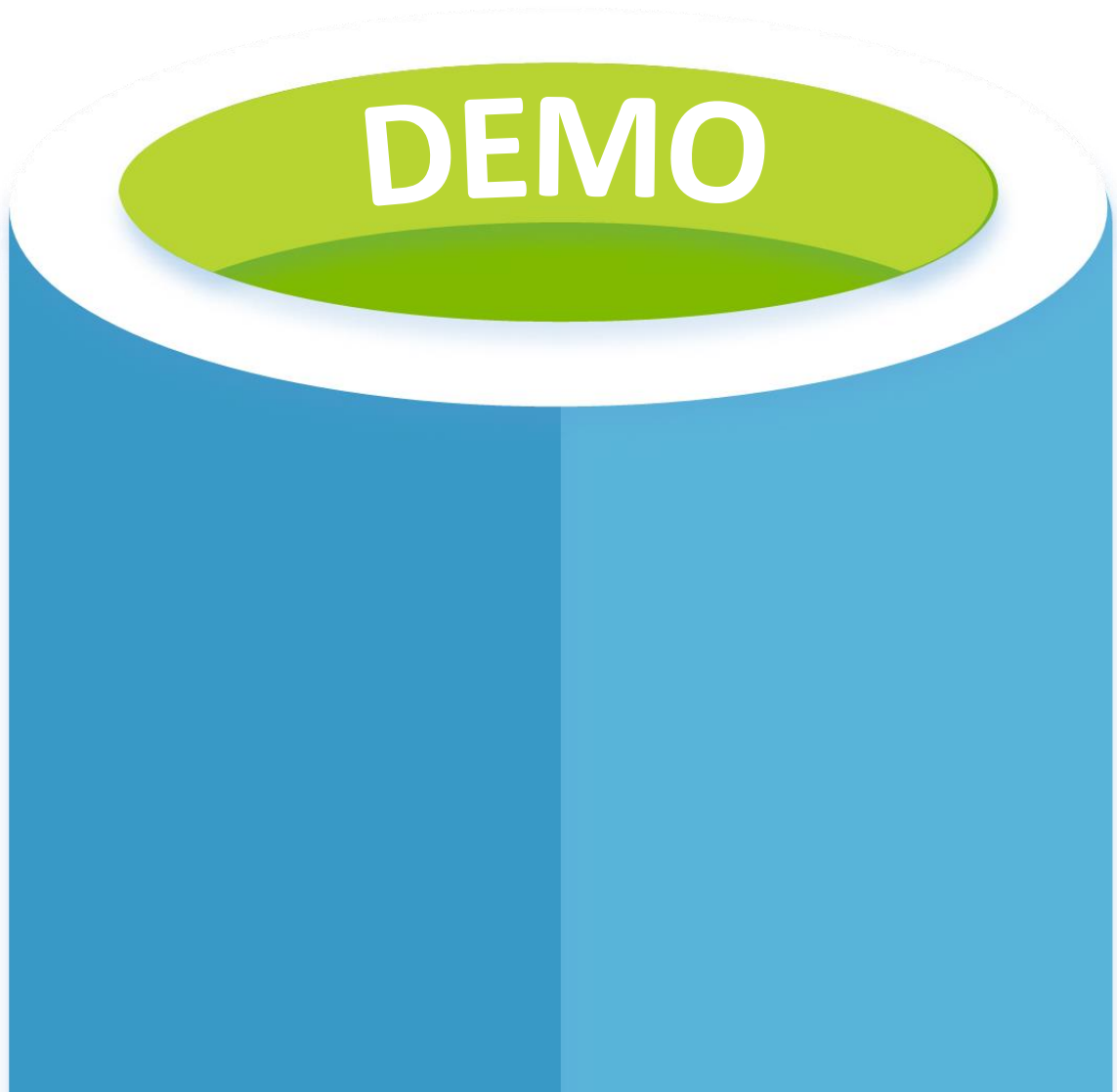**Expression Builder**

@{........} ⬅ Parameters
System Variables

- Collection
- Conversation
- Date
- Logical
- Math
- String

Add dynamic content [Alt+P]

# Dynamic Data Factory Pipelines

# Demo Architecture 1

# Demo Architecture 2



Source
Logical Instance

Target
Logical Instance

Source DB

Metadata

Data Factory

Target DB

Set Variable ✓
{x} Set Source Password

Set Variable
{x} Set Target Password

Lookup ✓
GetTableList

ForEach ✓
Copy All Tables

Stored Procedure
Truncate Target Table

Copy Data
Copy Table

{;} JSON

1x Linked Service
1x Dataset

# Web Activities vs Web Hook Activities

# Web Hook vs Web Activity

| Web | Web Hook |
|---|---|
| PUT<br>POST<br>GET<br>DELETE | POST |
| 1 Minute Timeout | Configurable Timeout |
| Retry Capabilities | No Retry |
| Linked Services<br>Datasets | No Artifact Support |
| One Way Call | Call Back URL |

# Web Hook vs Web Activity

| **Asynchronous** | | **Synchronous** |
|:---:|:---:|:---:|
| Web | | Web Hook |
| PUT POST GET DELETE | | POST |
| 1 Minute Timeout | | Configurable Timeout |
| Retry Capabilities | | No Retry |
| Linked Services Datasets | | No Artifact Support |
| One Way Call | | Call Back URL |

https://mrpaulandrew.com/2019/06/18/azure-data-factory-web-hook-vs-web-activity/

# Web Hook vs Web Activity

Web Hook

# Web Hook vs Web Activity

Scale Up SQLDB → Process Database → Scale Down SQLDB

# Web Hook vs Web Activity

Scale Up
SQLDB

Process
Database

Scale Down
SQLDB

```powershell
#Get web hook call body information:
if ($WebhookData) {
    $parameters = (ConvertFrom-Json -InputObject $WebhookData.RequestBody)
    if ($parameters.callBackUri) {$callBackUri = $parameters.callBackUri}
}

$SizeInBytes = 250 * 1024 * 1024 * 1024

#Set SQLDB tier:
Set-AzureRmSqlDatabase `
    -Edition "Standard" `
    -RequestedServiceObjectiveName "S9" `
    -MaxSizeBytes $SizeInBytes | out-null

#Call back to ADF:
if ($callBackUri) {
    Invoke-WebRequest -Uri $callBackUri -Method POST
}
```

# Web Hook vs Web Activity

Scale Up SQLDB

ADF Web Hook Activity Body

```json
1  {
2      "WebhookName":"SQLDBScaler",
3      "RequestBody":{
4          "count": 1,
5          "value": [
6              {
7                  "EnvironmentName": "AzureCloud",
8                  "ResourceGroupName": "ResourceGroup01",
9                  "ServerName": "SQLInstance01",
10                 "DatabaseName": "Database01",
11                 "TargetEdition": "Standard",
12                 "TargetTier": "S3"
13             }
14         ],
15         "effectiveIntegrationRuntime": "PaulsFunFactoryIR01 (North Europe)",
16         "callBackUri": "https://PMNortheurope.svc.datafactory.azure.com/workflow/callback/##RUNID##?callbackUrl=##TOKEN##"
17         "RequestHeader":
18             {
19                 "Connection":"Keep-Alive",
20                 "Expect":"100-continue",
21                 "Host":"s9events.azure-automation.net",
22                 "x-ms-request-id":""
23             }
24     }
```

```powershell
#Get web hook cal
if ($WebhookData)
    $parameters
        if ($parameters.callBackUri)  {$callBackUri = $parameters.callBackUri}
}

$SizeInBytes = 250 * 1024 * 1024 * 1024

#Set SQLDB tier:
Set-AzureRmSqlDatabase `
    -Edition "Standard" `
    -RequestedServiceObjectiveName "S9" `
    -MaxSizeBytes $SizeInBytes | out-null

#Call back to ADF:
if ($callBackUri) {
    Invoke-WebRequest -Uri $callBackUri -Method POST
}
```
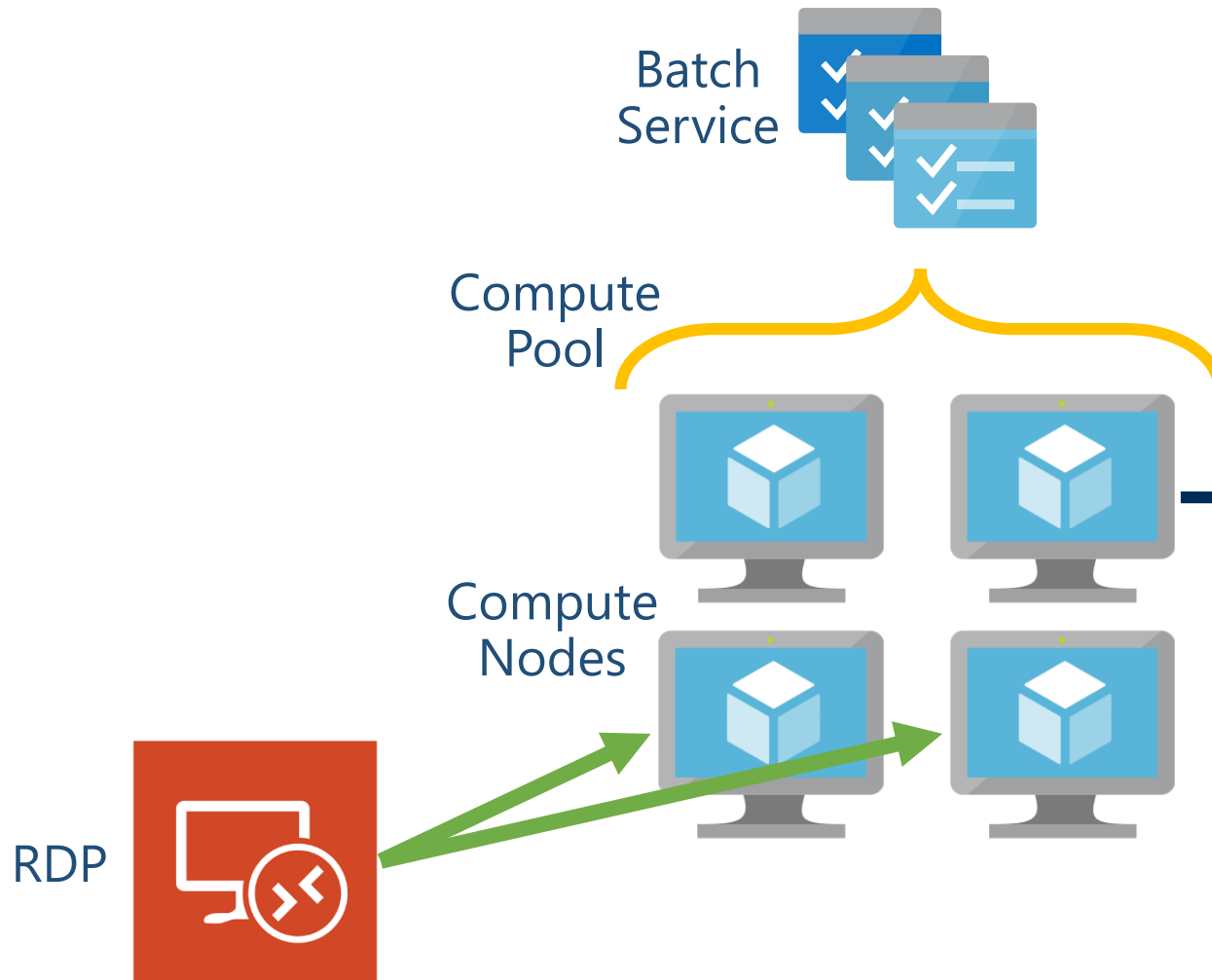
# Custom Activities

# Azure Batch Service

Scale out compute delivered using PaaS technology with IaaS underneath.

Batch Service

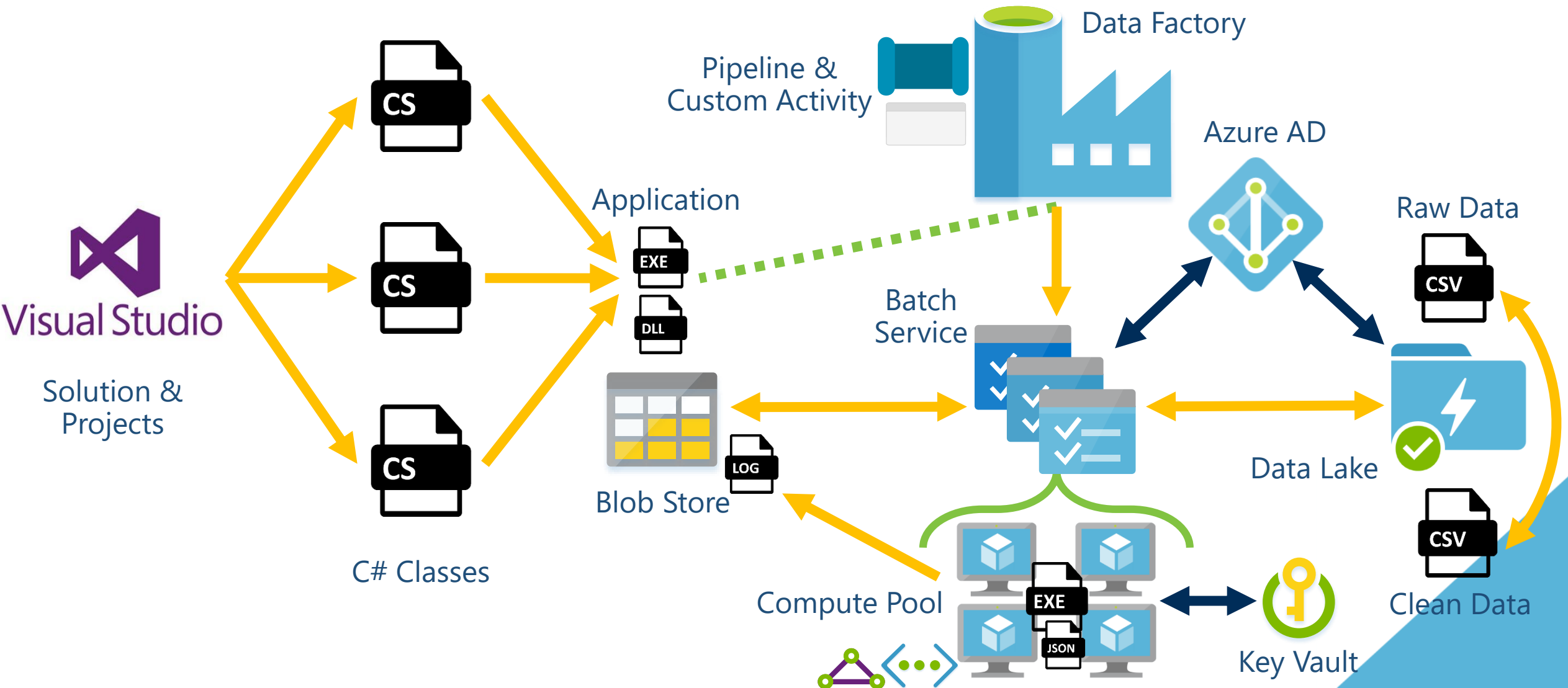Compute Pool

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.

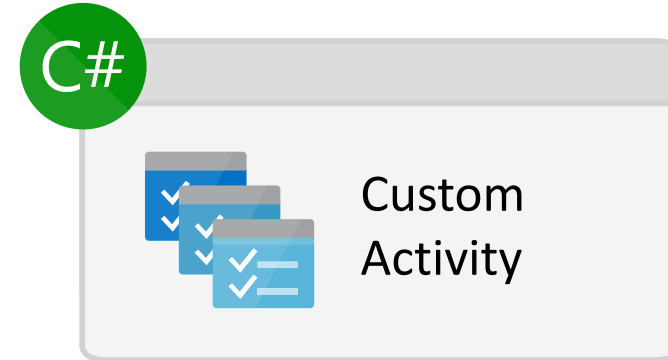▶ Special environment variables.

# Building a Custom Activity

A .Net Console App Executed Using Azure Batch Service.

# Extensibility Conclusions

**C#**

Do Stuff

10 minutes execution

**C#**

Custom Activity

Auto scale out compute &
Scale up per compute node

**REST**

Web

Asynchronous

**REST**

Web Hook

Synchronous
(Call back)

Data Transformation in Azure with SSIS

* MAXDOP 80

SSIS Package

8x Packages per Node

10x Nodes

SSIS Integration Runtime

Express Route

VNet

Data Factory

Logical or Managed Instance

SSISDB

Visual Studio – SQL Server Data Tools (SSDT)

Sales Order Header

Sales Order Details

Merge Join

Aggregate

Order Line Count Table

# The SSIS IR vs Hosted IR with Express Route



Hosted IR(s)
**4x nodes**

Copy Activity

VNet

Express Route

SSIS IR

Storage

80x parallel copies

CSV  XLS  TXT

**On Premises**  **Azure**

# The SSIS IR Start/Stop



VNet

SSIS IR

Express Route

Storage

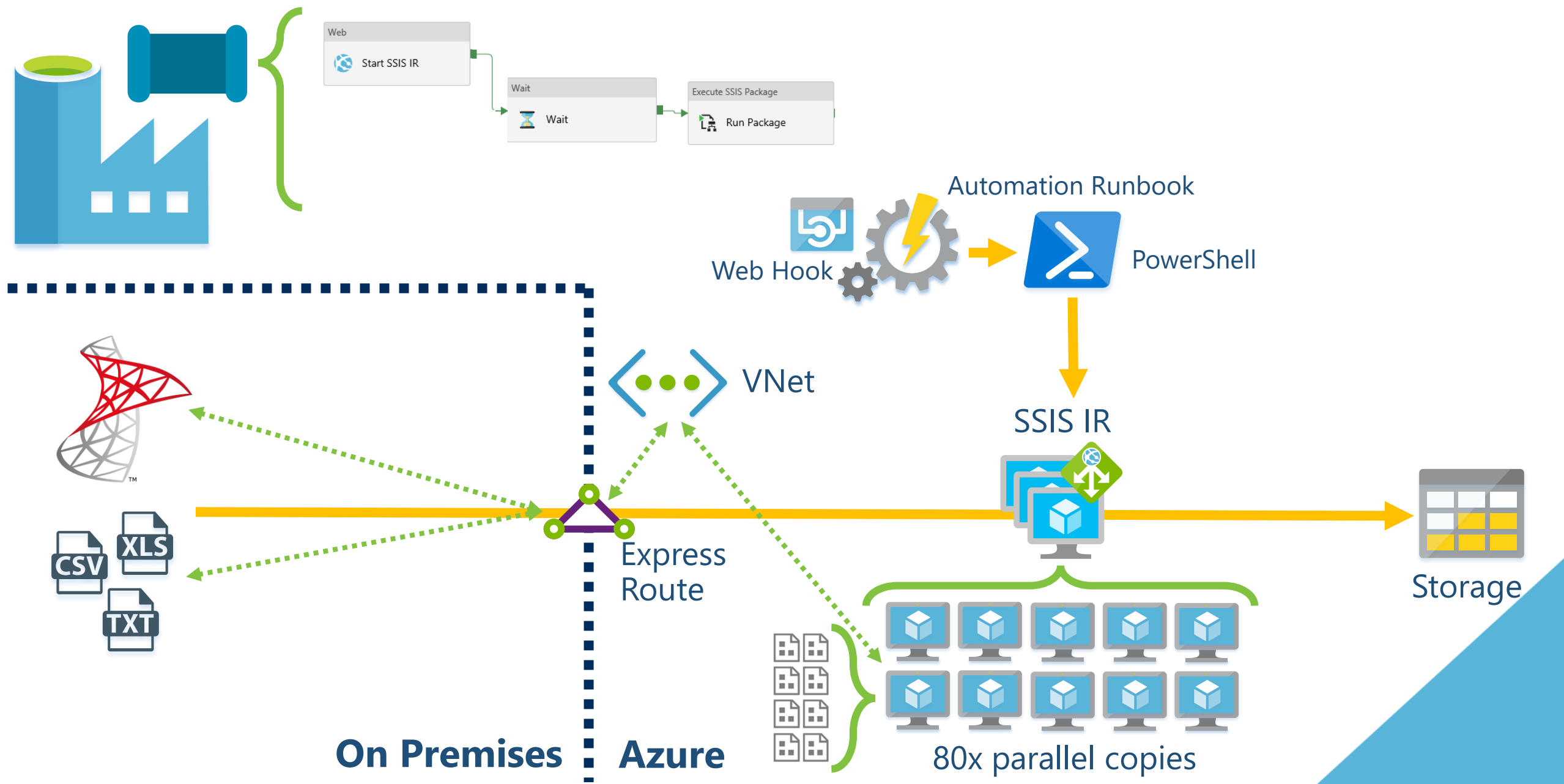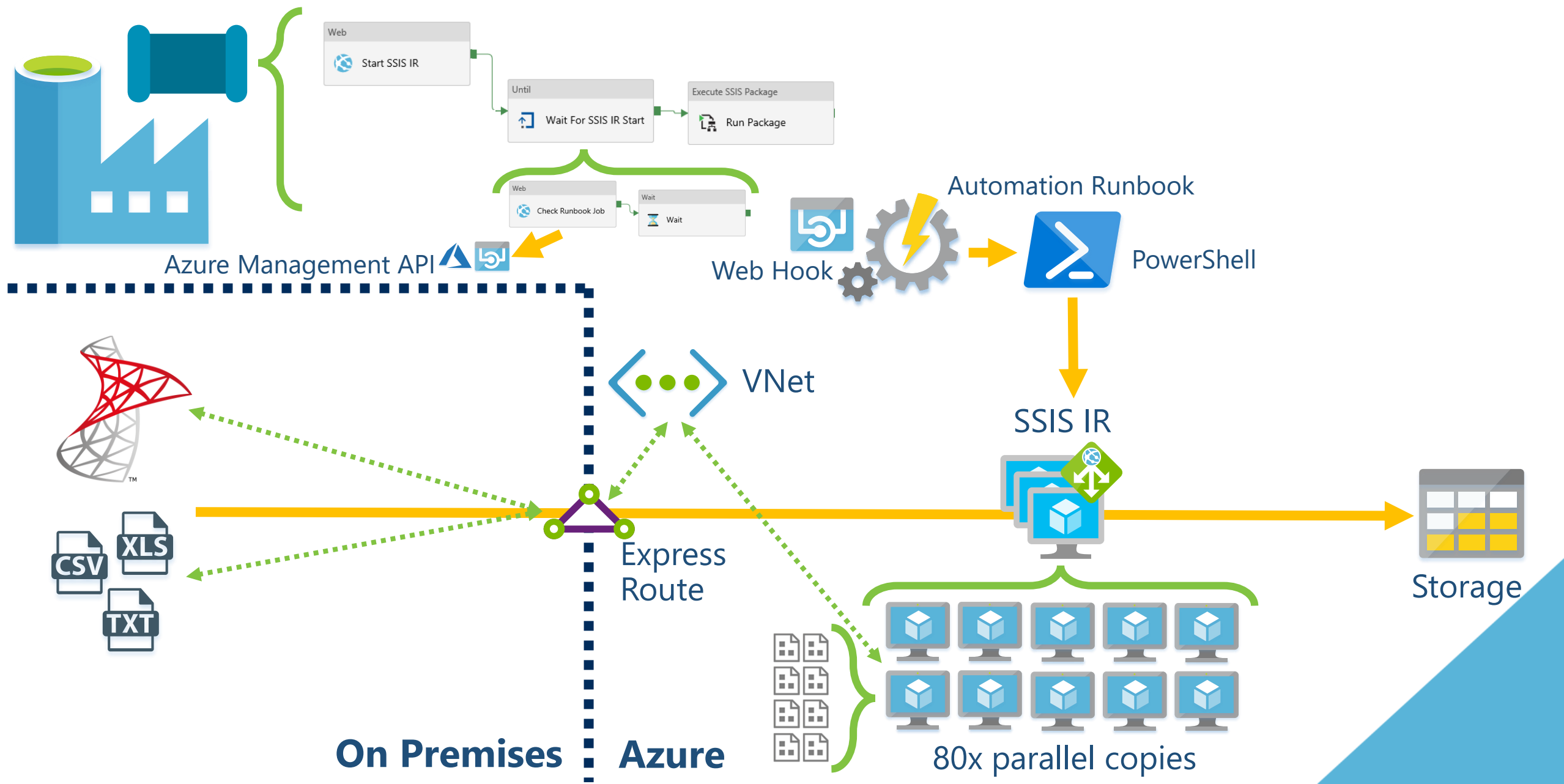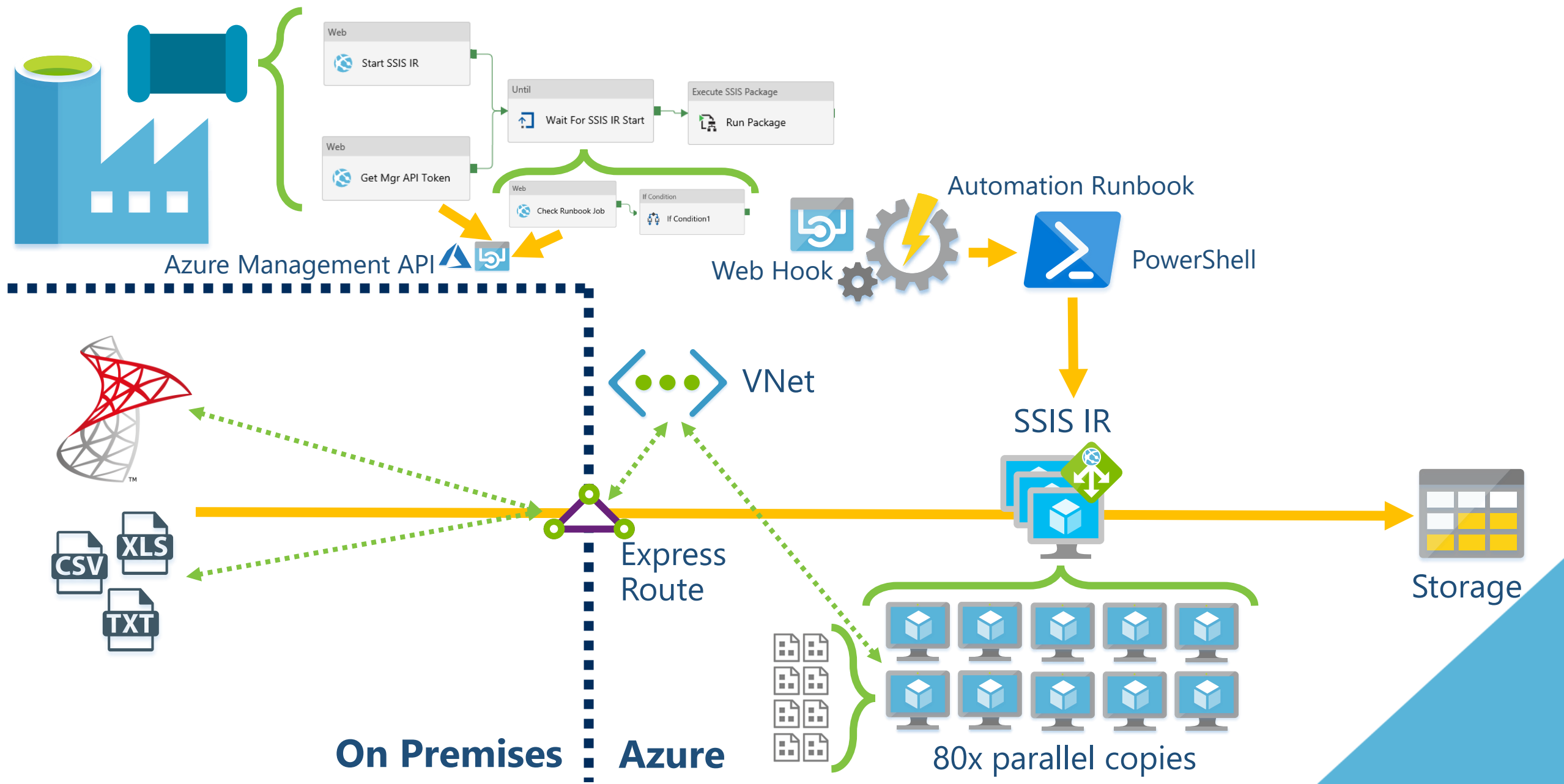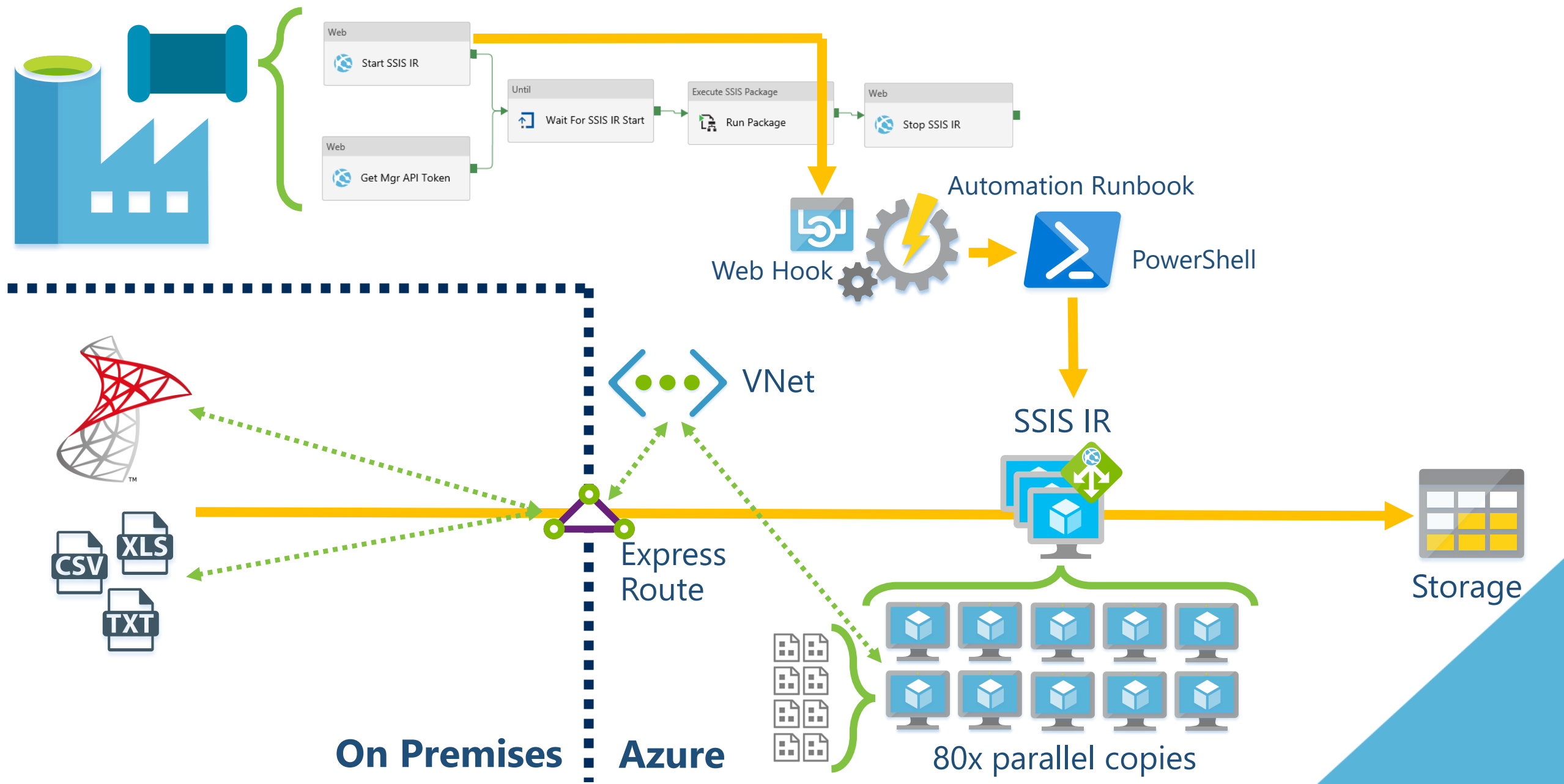On Premises | Azure

80x parallel copies

# The SSIS IR Start/Stop

# The SSIS IR Start/Stop

# The SSIS IR Start/Stop

# The SSIS IR Start/Stop

Web
Start SSIS IR

Until
Wait For SSIS IR Start

Execute SSIS Package
Run Package

Web
Get Mgr API Token

Web
Check Runbook Job

If Condition
If Condition1

Azure Management API

Web Hook

Automation Runbook

PowerShell

VNet

SSIS IR

Express Route

CSV   XLS   TXT

Storage

80x parallel copies

**On Premises**   **Azure**

# The SSIS IR Start/Stop



Web — Start SSIS IR

Web — Get Mgr API Token

Until — Wait For SSIS IR Start

Execute SSIS Package — Run Package

Web — Stop SSIS IR

Web Hook

Automation Runbook

PowerShell

VNet

Express Route

SSIS IR

Storage

CSV  XLS  TXT

80x parallel copies

**On Premises**  :  **Azure**

# The SSIS IR Parallelism

Lookup
Get SSIS Packages

ForEach
Run Packages

* MAXDOP 50

Execute SSIS Package
Run Package

SSISDB

VNet

SSIS IR

Express Route

* MAXDOP 80

Storage

**On Premises** | **Azure**

80x parallel copies

# The SSIS IR Parallelism



Lookup — Get SSIS Packages
ForEach — Run Packages
Execute SSIS Package — Run Package

SSISDB

VNet

SSIS IR

Express Route

Storage

**On Premises** : **Azure**

80x parallel copies

# The SSIS IR Parallelism



Copy Data — Get SSIS Data
Stored Procedure — Assign Buckets
Lookup — Get Bucket Details
ForEach — Run Execution Buckets
Lookup — Get Package IDs for...

Async vs Sync

* MAXDOP 50

Execute Pipeline — Execute Bucket

* New MAXDOP 2500

Lookup — Get SSIS Packages
ForEach — Run Packages
Execute SSIS Package — Run Package

* MAXDOP 50

SSISDB

Metadata

VNet

SSIS IR

Express Route

Storage

CSV  XLS  TXT

On Premises : Azure

80x parallel copies

# The SSIS IR Parallelism

Async vs Sync

| Resource | Default Limit | Maximum Limit |
|----------|---------------|---------------|
| Data factories in an Azure subscription | 50 | Contact support |
| Total number of entities (Pipeline, Datasets, Triggers, Linked Services, Integration runtimes) within a data factory | 5000 | Contact support |
| Total CPU cores for Azure-SSIS Integration Runtime(s) under one subscription | 256 | Contact support |
| Concurrent pipeline runs per data factory (shared among all pipelines in the factory) | 10,000 | Contact support |
| Max activities per pipeline (includes inner activities for containers) | 40 | 40 |
| Max number of Linked Integration Runtime that can be created against a single Self-hosted Integration Runtime | 20 | Contact support |
| Max parameters per pipeline | 50 | 50 |
| ForEach items | 100,000 | 100,000 |
| ForEach parallelism | 20 | 50 |
| Characters per expression | 8,192 | 8,192 |
| Minimum Tumbling Window Trigger interval | 15 min | 15 min |
| Max Timeout for pipeline activity runs | 7 days | 7 days |
| Bytes per object for pipeline objects [1] | 200 KB | 200 KB |
| Bytes per object for dataset and linked service objects [1] | 100 KB | 2000 KB |
| Data integration units per copy activity run [3] | 256 | Contact support |

Execute Pipeline — Execute Bucket

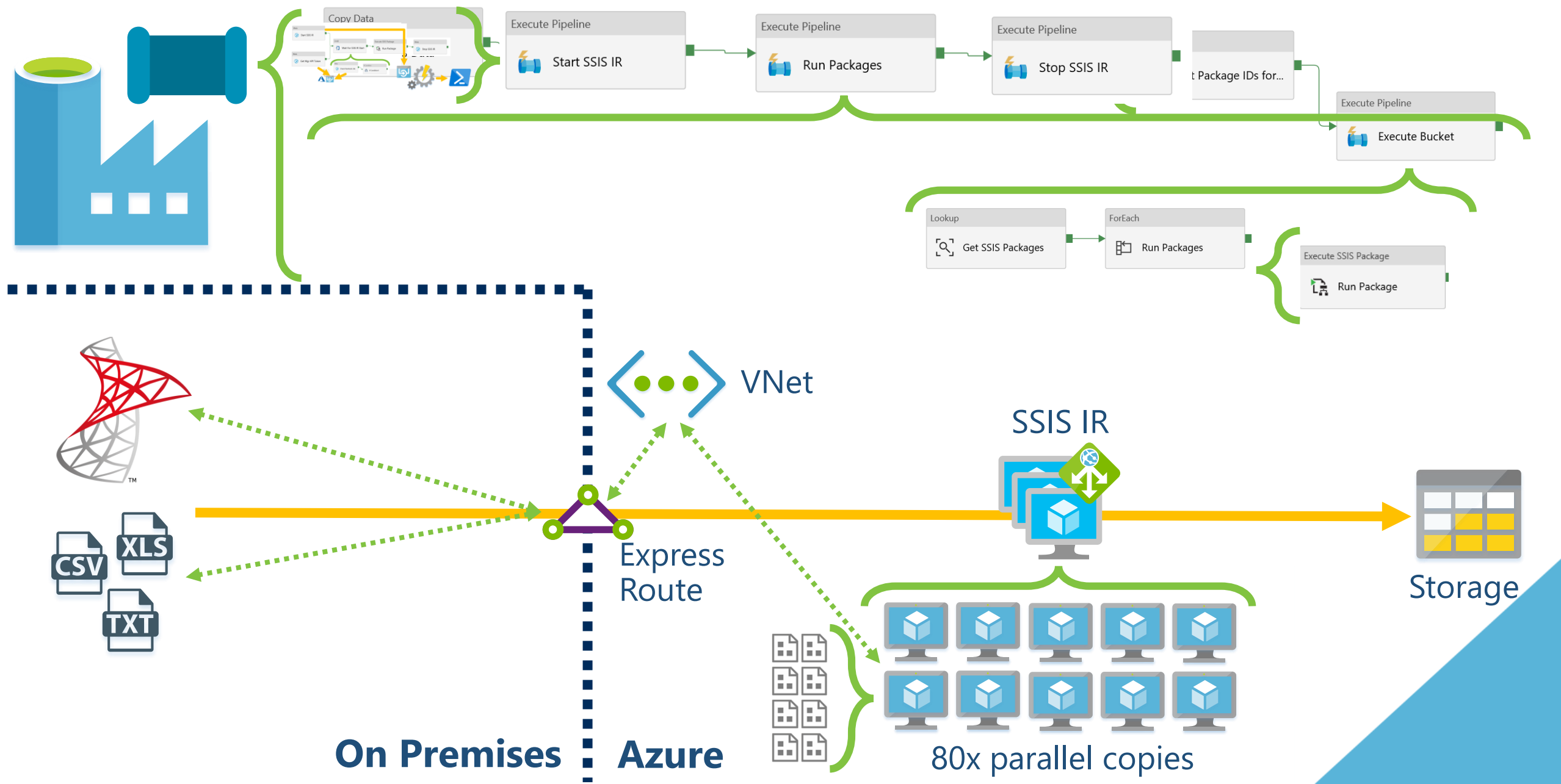Execute SSIS Package — Run Package

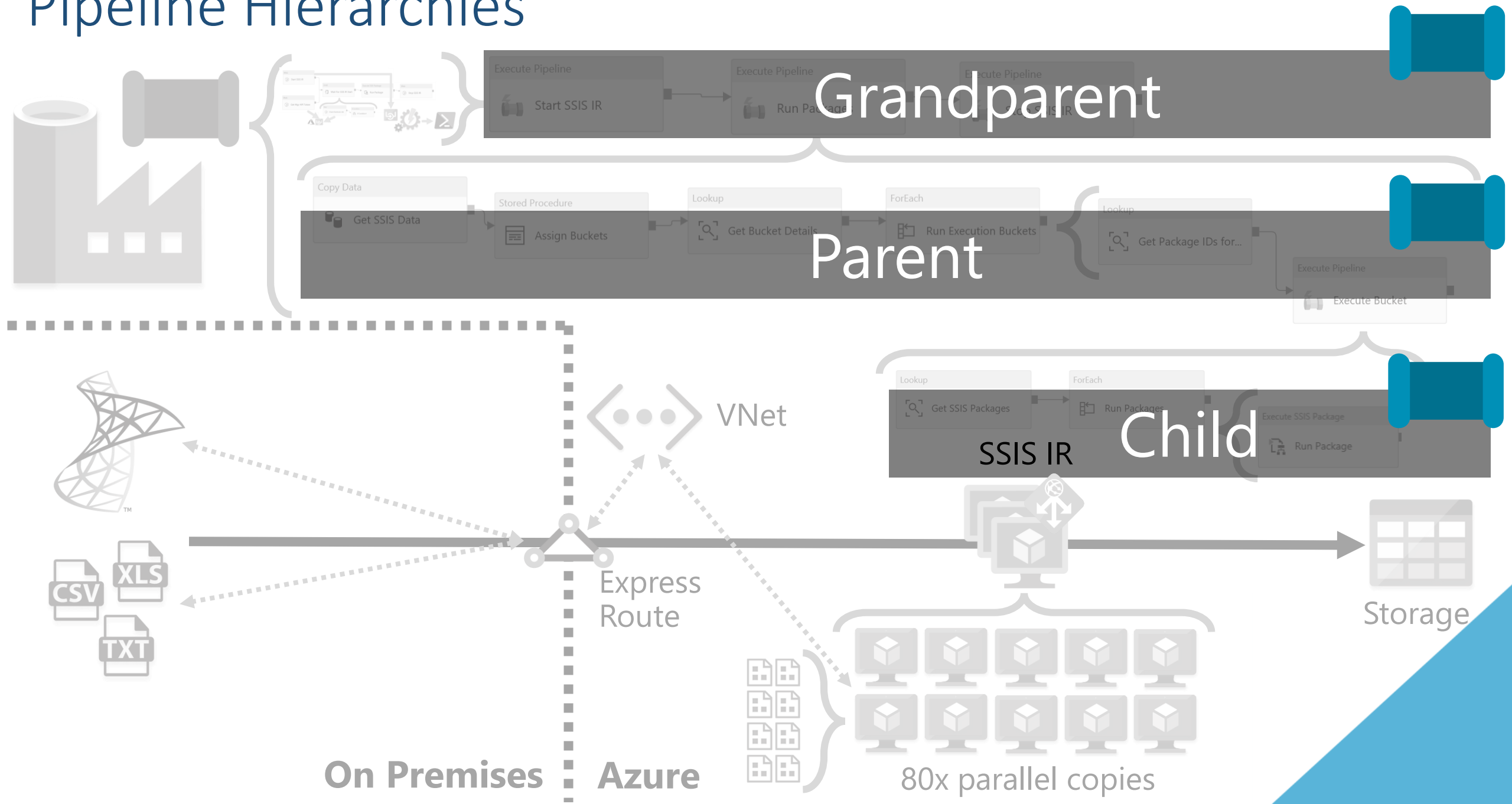Package IDs for...

UP 50

Storage

https://github.com/MicrosoftDocs/azure-docs/blob/master/includes/azure-data-factory-limits.md
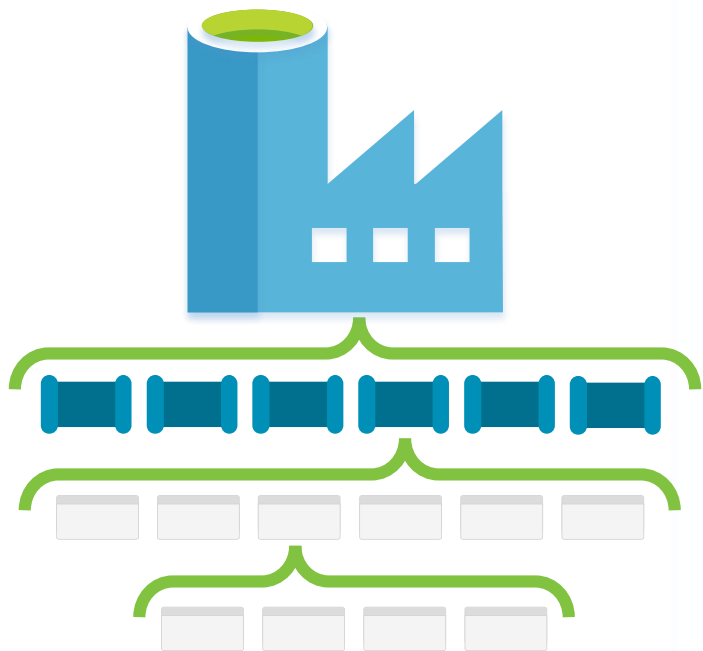
**On Premises : Azure**

80x parallel copies

# SSIS IR & Package Complete Orchestration Solution



**Copy Data**

**Execute Pipeline** — Start SSIS IR

**Execute Pipeline** — Run Packages

**Execute Pipeline** — Stop SSIS IR

...t Package IDs for...

**Execute Pipeline** — Execute Bucket

**Lookup** — Get SSIS Packages

**ForEach** — Run Packages

**Execute SSIS Package** — Run Package

VNet

SSIS IR

Express Route

Storage

On Premises | Azure

80x parallel copies

# Pipeline Hierarchies



**Grandparent**

Execute Pipeline — Start SSIS IR
Execute Pipeline — Run Package
Execute Pipeline — Stop SSIS IR

**Parent**

Copy Data — Get SSIS Data
Stored Procedure — Assign Buckets
Lookup — Get Bucket Details
ForEach — Run Execution Buckets
Lookup — Get Package IDs for...
Execute Pipeline — Execute Bucket

**Child**

Lookup — Get SSIS Packages
ForEach — Run Packages
Execute SSIS Package — Run Package

SSIS IR

VNet

Express Route

Storage

On Premises : Azure

80x parallel copies

CSV  XLS  TXT

# Pattern Summary

**Execute Pipeline**
Grandparent

## High Level Control Flow and Pipeline Triggers

**Execute Pipeline**
Parent

## Platform Component Control
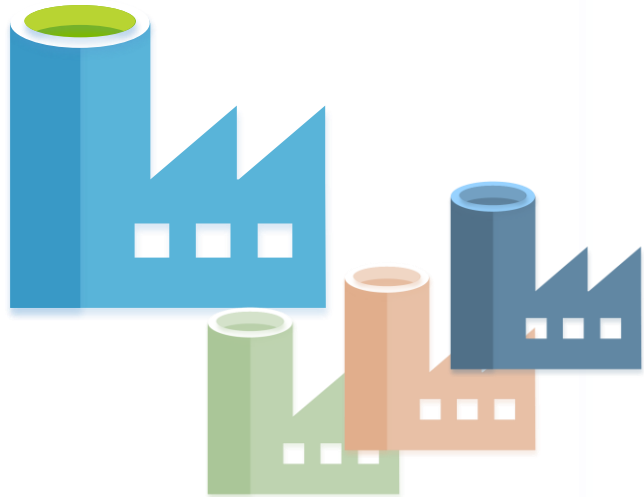
## Manage Parallel Activities
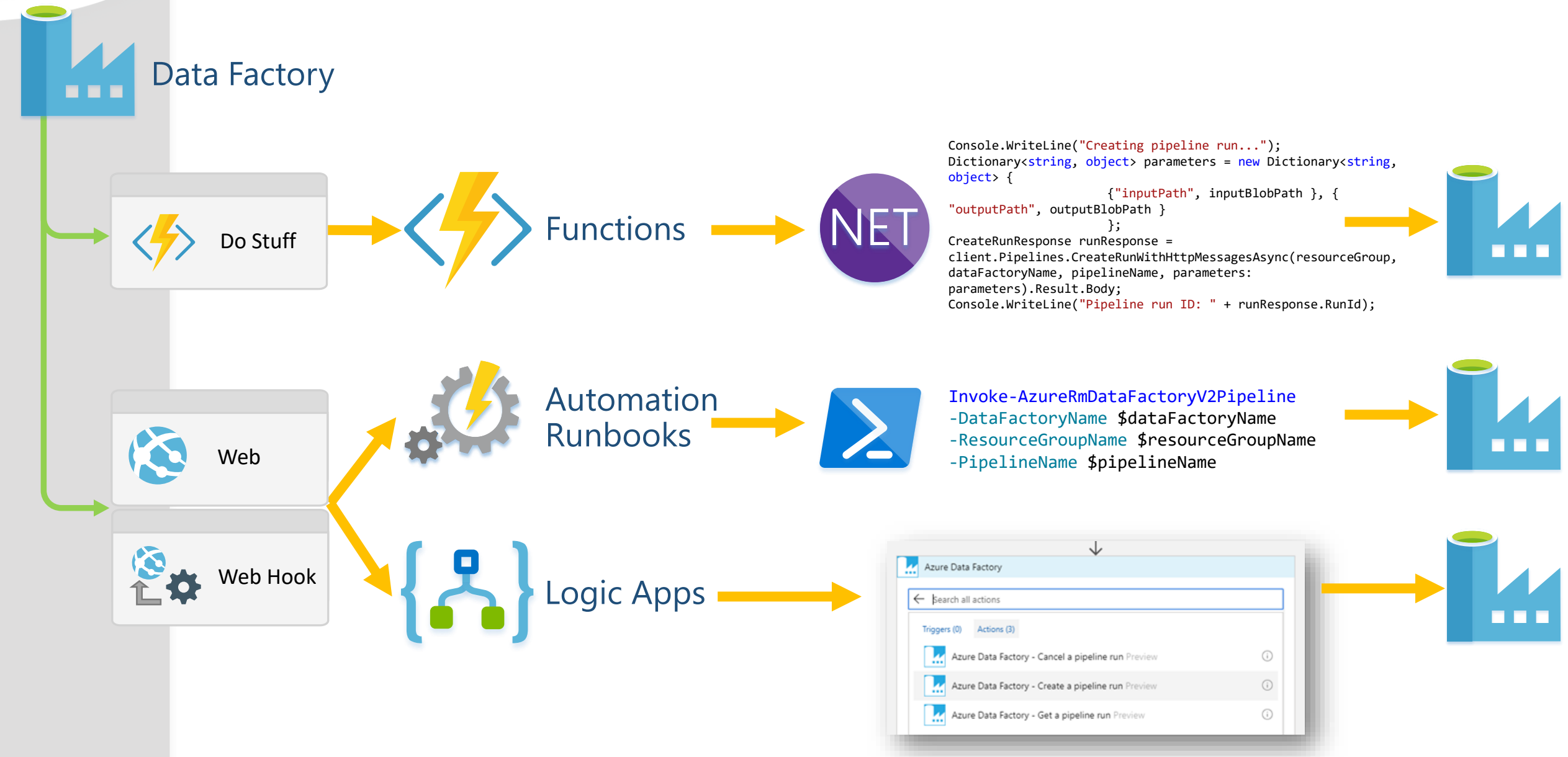
**Execute Pipeline**
Child

## Service Level Executions

# Solution Bootstrapping

# Bootstrapping

Data Factory

## Do Stuff → Functions
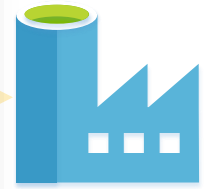
```
Console.WriteLine("Creating pipeline run...");
Dictionary<string, object> parameters = new Dictionary<string,
object> {
                        {"inputPath", inputBlobPath }, {
"outputPath", outputBlobPath }
                        };
CreateRunResponse runResponse =
client.Pipelines.CreateRunWithHttpMessagesAsync(resourceGroup,
dataFactoryName, pipelineName, parameters:
parameters).Result.Body;
Console.WriteLine("Pipeline run ID: " + runResponse.RunId);
```

## Web → Automation Runbooks

```
Invoke-AzureRmDataFactoryV2Pipeline
-DataFactoryName $dataFactoryName
-ResourceGroupName $resourceGroupName
-PipelineName $pipelineName
```

## Web Hook → Logic Apps

Azure Data Factory

← Search all actions

Triggers (0)    Actions (3)

Azure Data Factory - Cancel a pipeline run Preview

Azure Data Factory - Create a pipeline run Preview

Azure Data Factory - Get a pipeline run Preview

Bootstrapping

# Bootstrapping

Data Factory

Do Stuff → Functions → NET

```
Console.WriteLine("Creating pipeline run...");
Dictionary<string, object> parameters = new Dictionary<string,
object> {
                    {"inputPath", inputBlobPath }, {
"outputPath", outputBlobPath }
                    };
CreateRunResponse runResponse =
client.Pipelines.CreateRunWithHttpMessagesAsync(resourceGroup,
dataFactoryName, pipelineName, parameters:
parameters).Result.Body;
Console.WriteLine("Pipeline run ID: " + runResponse.RunId);
```

Web → Automation Runbooks

```
Invoke-AzureRmDataFactoryV2Pipeline
-DataFactoryName $dataFactoryName
-ResourceGroupName $resourceGroupName
-PipelineName $pipelineName
```

Web Hook → Logic Apps

Azure Data Factory

Search all actions

Triggers (0)   Actions (3)

Azure Data Factory - Cancel a pipeline run Preview

Azure Data Factory - Create a pipeline run Preview

Azure Data Factory - Get a pipeline run Preview

# Bootstrapping – Wider Analytics Solution

Data Factory

# Bootstrapping – Wider Analytics Solution



Data Factory

Analytics Solution

# Bootstrapping – Wider Analytics Solution



Logic Apps

Data Factory

Source Systems

Analytics Solution

**On Premises** **Azure**

# Data Factory DevOps – CI/CD

# Data Factory Continuous Integration

Developer

Developer

Developer

Git

{;} JSON

branch

save

merge

master

publish

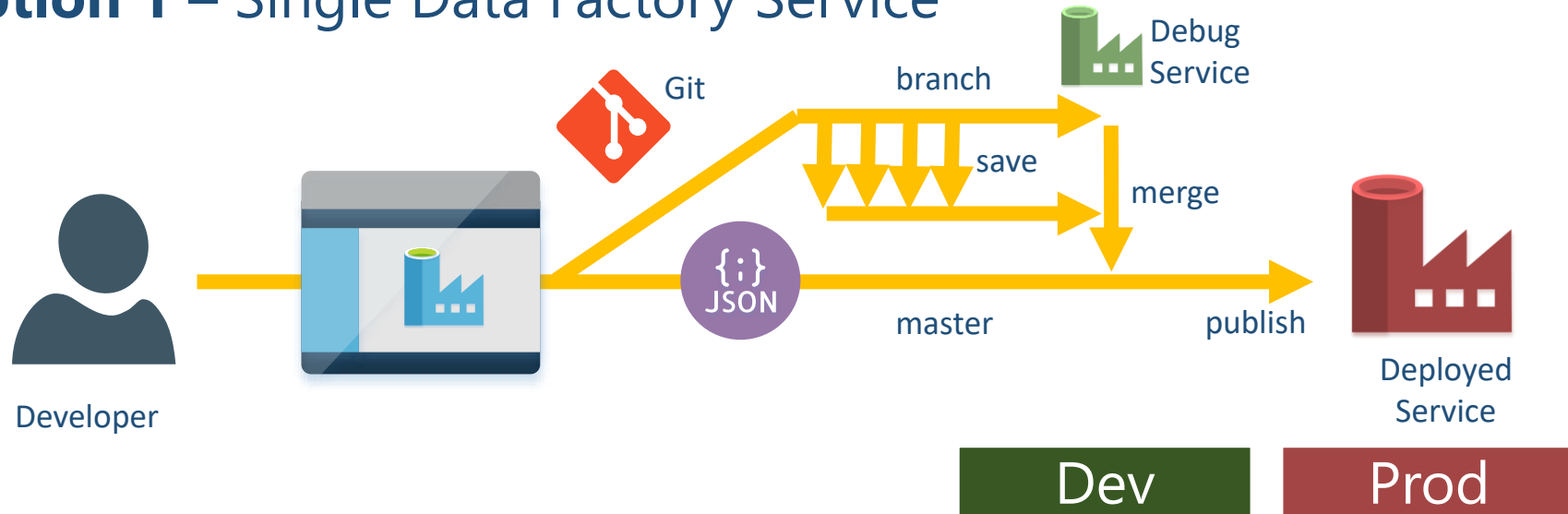Debug Service

Deployed Service

Dev

Prod

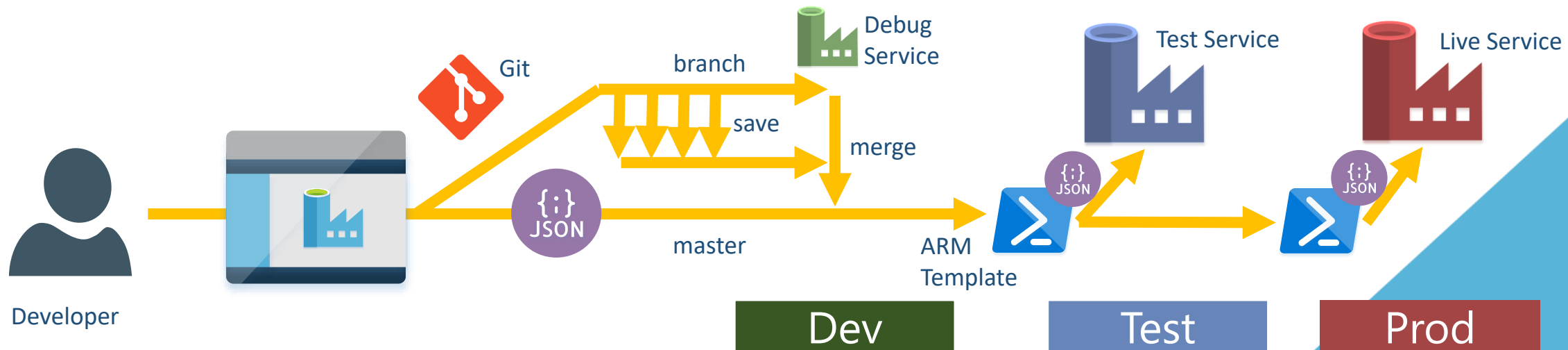# Data Factory Continuous Delivery
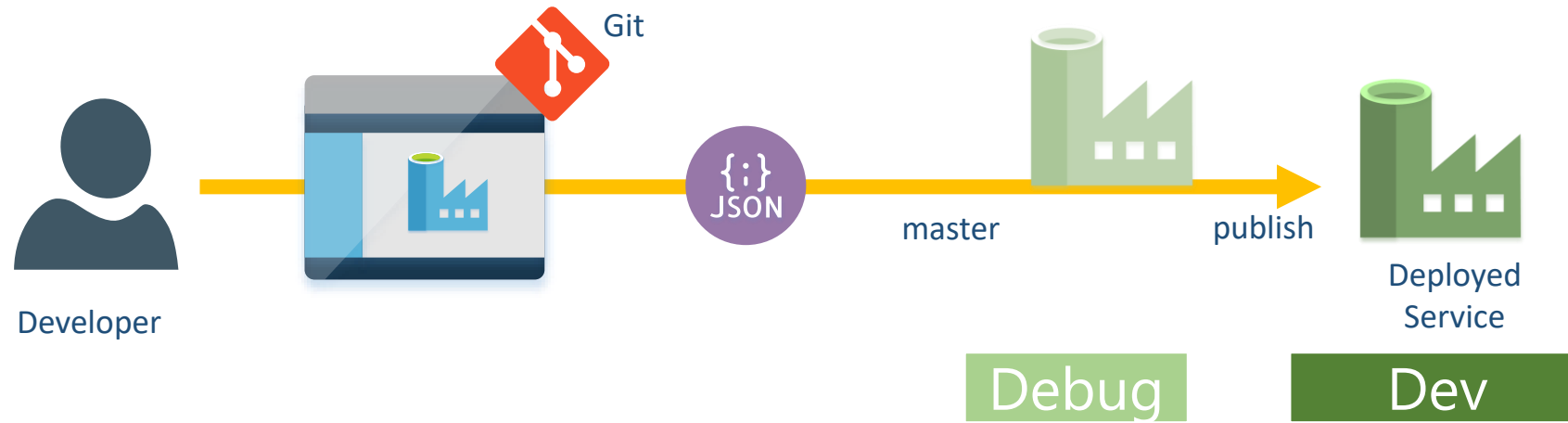
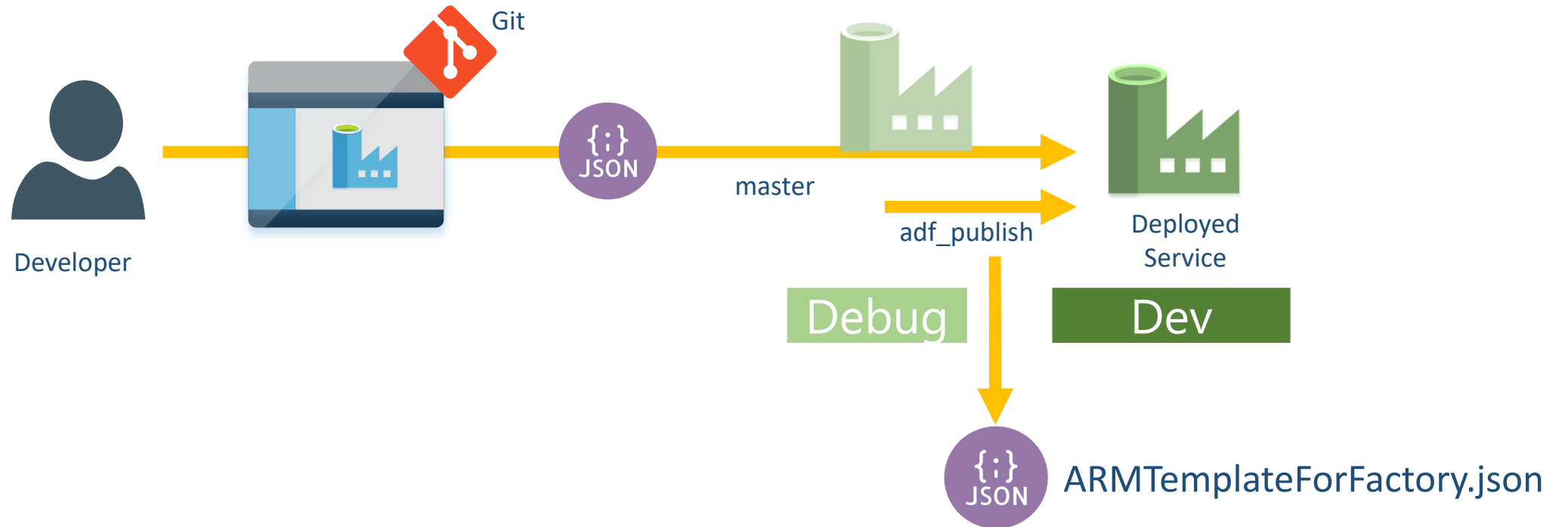# Data Factory Continuous Delivery

## Option 1 – Single Data Factory Service



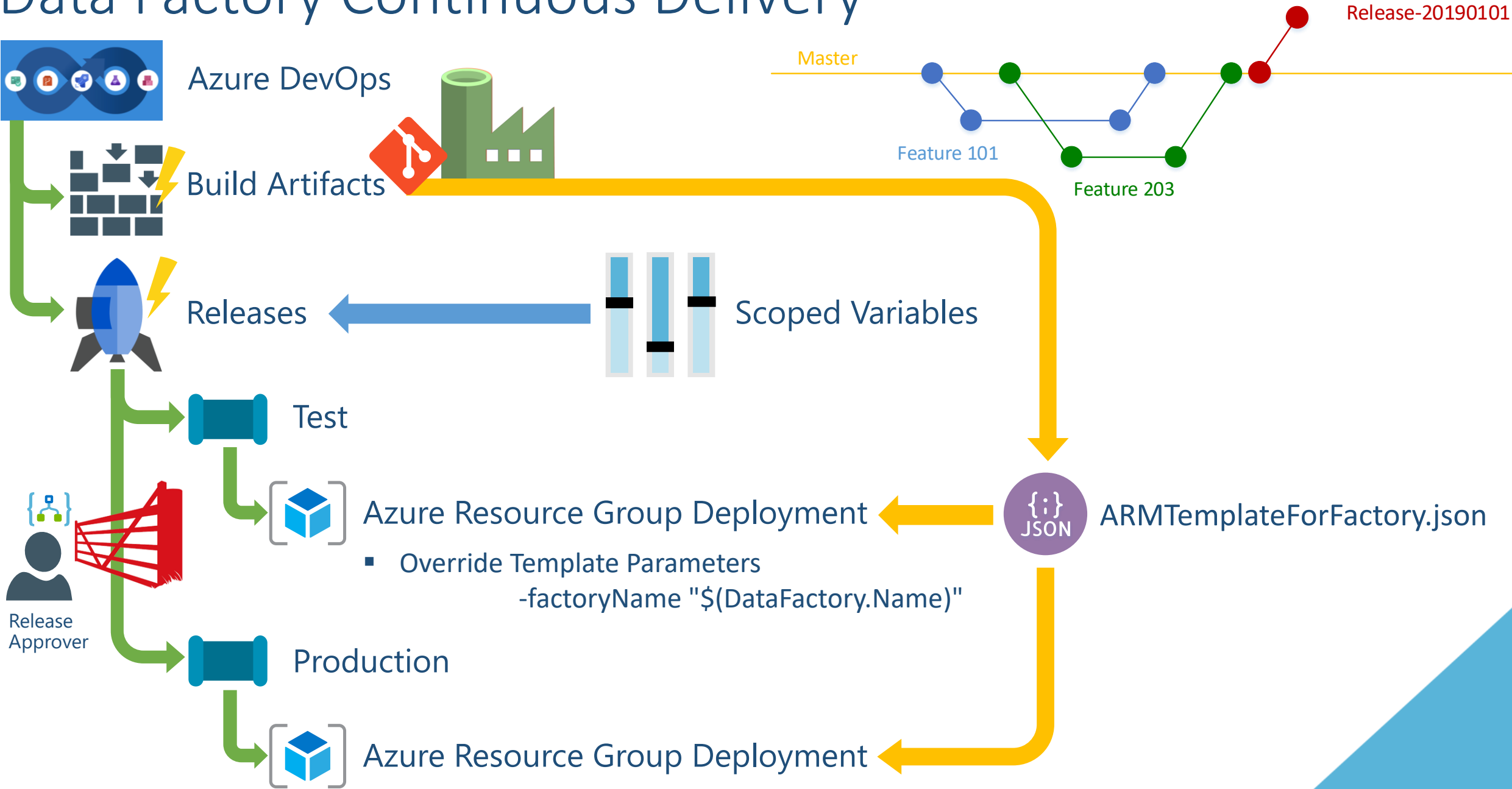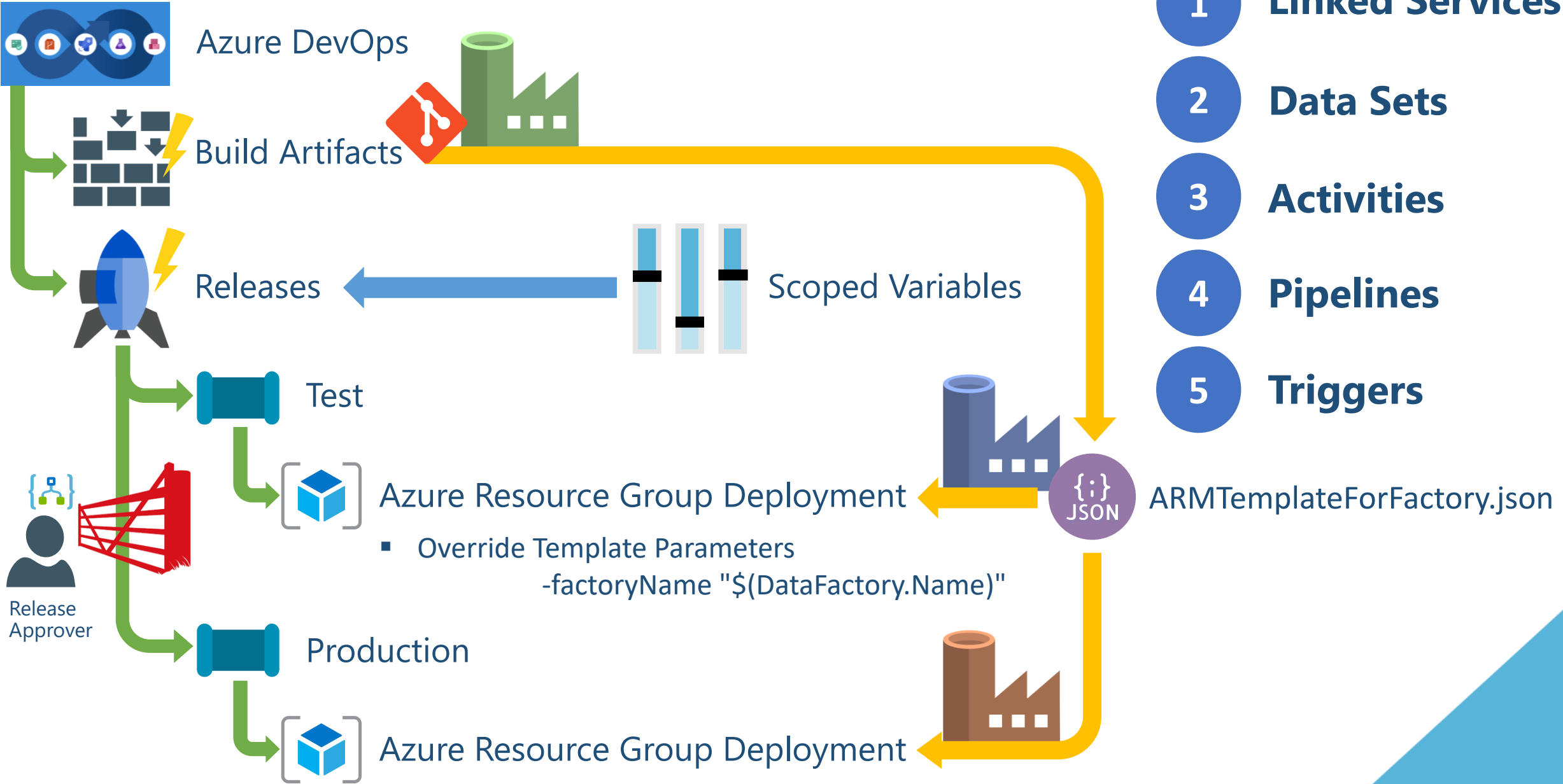## Option 2 – ARM Templates for Multiple Data Factory Services

# Data Factory Publish

Developer

Git

JSON

master    publish

Deployed
Service

Debug    Dev

# Data Factory Publish

Data Factory Continuous Delivery

# Data Factory Continuous Delivery



Azure DevOps

Build Artifacts

Releases

Scoped Variables

Test

Release Approver

Azure Resource Group Deployment

- Override Template Parameters
  -factoryName "$(DataFactory.Name)"

Production

Azure Resource Group Deployment

ARMTemplateForFactory.json

1 **Linked Services**

2 **Data Sets**

3 **Activities**

4 **Pipelines**

5 **Triggers**

# Data Factory Continuous Delivery

Azure DevOps

Build Artifacts

Releases

Scoped Variables

Test

Release Approver

Azure Resource Group Deployment

- Override Template Parameters
  -factoryName "$(DataFactory.Name)"

Production

Azure Resource Group Deployment

**1** **Linked Services**

ARMTemplateForFactory.json

# Session Agenda

- Data Factory – A Quick Overview ✓

- Dynamic Pipelines ✓

- Extending Data Factory ✓
  - Web Activities
  - Custom Activities

- True Scale Out Execution ✓
  - SSIS Integration Runtime

- Data Factory – In Production ✓
  - Bootstrapping
  - DevOps

# Thank you for listening...

Paul Andrew

Microsoft
**MVP** Most Valuable
Professional

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

**Twitter:** @mrpaulandrew
**LinkedIn:** In/mrpaulandrew

**GitHub:** github.com/mrpaulandrew