A Introduction to Azure Data Factory

Integration Pipelines



Paul Andrew | Technical Architect in Azure CoE

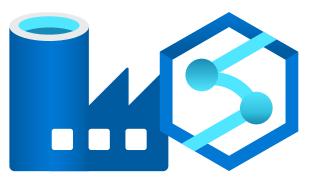












A Introduction to Azure Data Factory

Integration Pipelines



Paul Andrew | Technical Architect in Azure CoE















https://github.com/mrpaulandrew

CommunityEvents

Demo code, content and slides from various community events.

{Event/Location}-{Month}-{Year}

Agenda

What is it and why use it?

DDData Factory Data Flows

DData Factory Components

Source Control

DDCommon Activities

Deployments

DEEXECUTION Dependencies

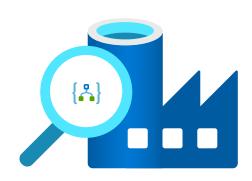
Monitoring & Logging

MIntegration Runtimes

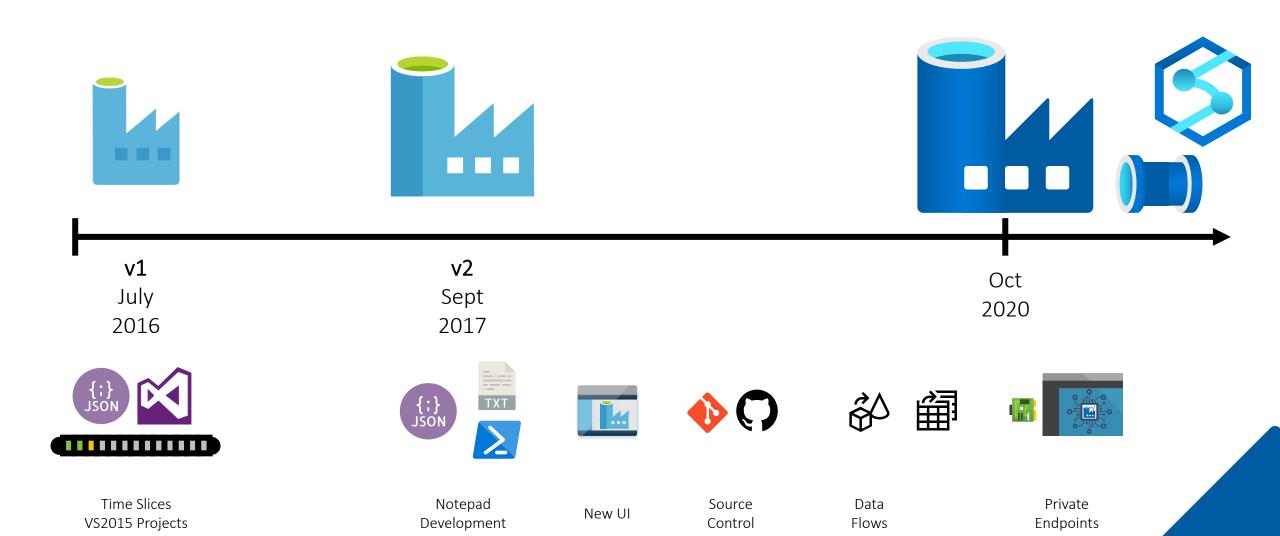
Conclusions

DD Azure/Hosted/SSIS

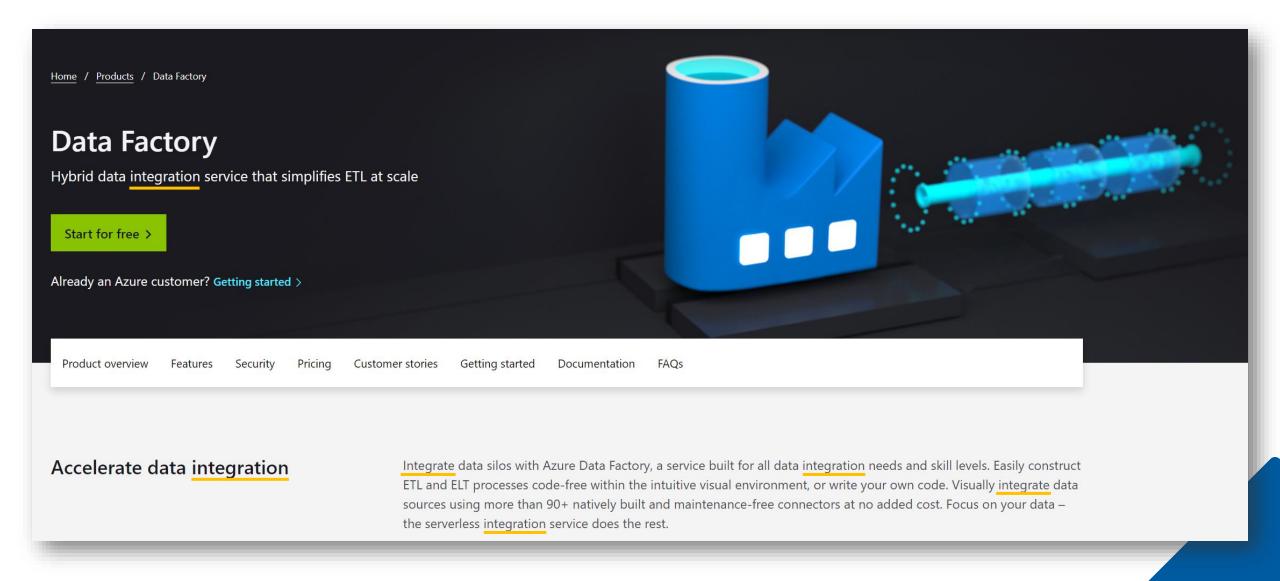
Azure Data Factory — What is it? Why use it?



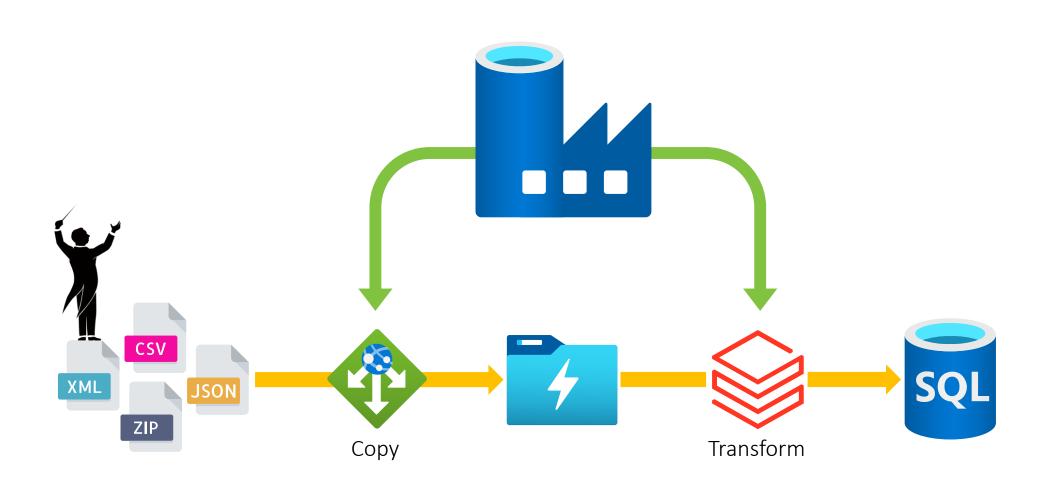
A Quick History Lesson

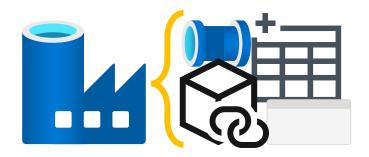


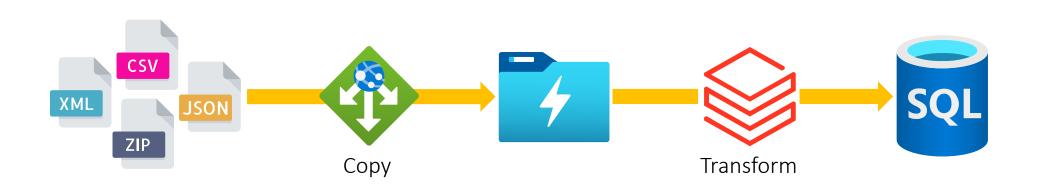
What is Azure Data Factory (ADF)?

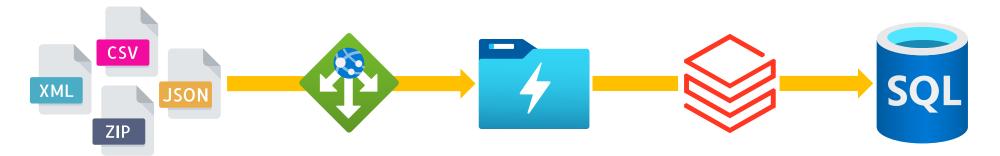


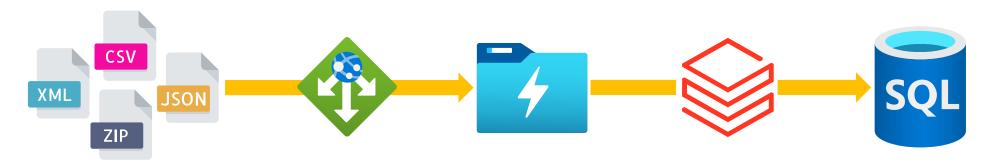
What is Azure Data Factory (ADF)?





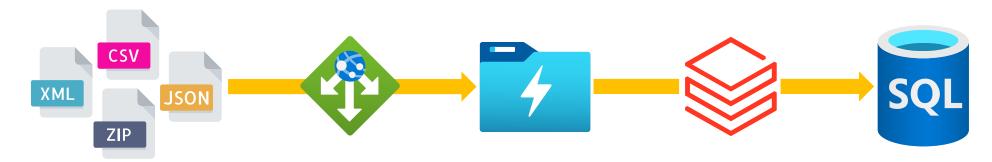




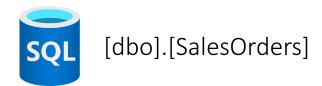


1 Linked Services — What to interact with and how?



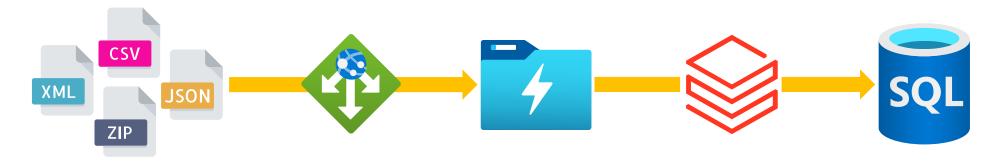


- 1 Linked Services
- Datasets Where is my data? What format? What file path/table do I need?

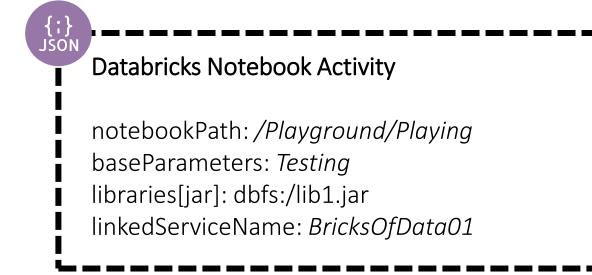


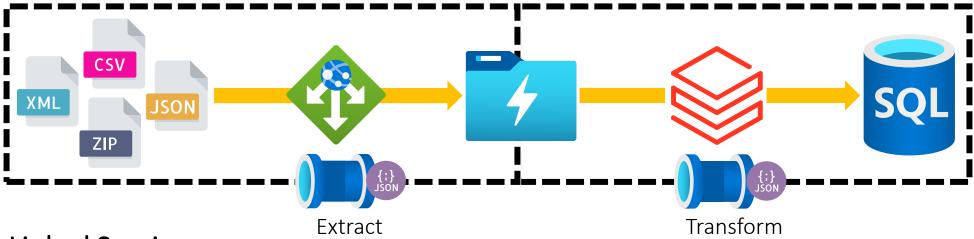


/RAW/Orders/2018/01/01/SalesOrders.csv

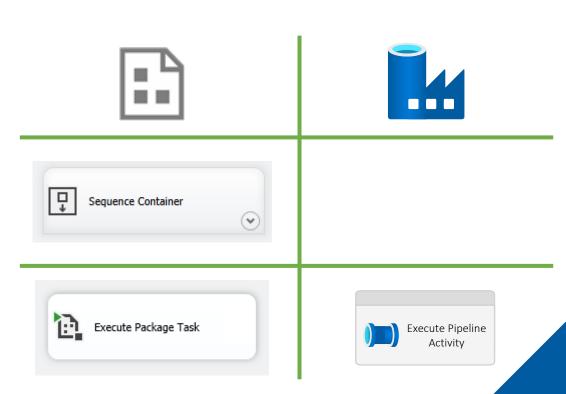


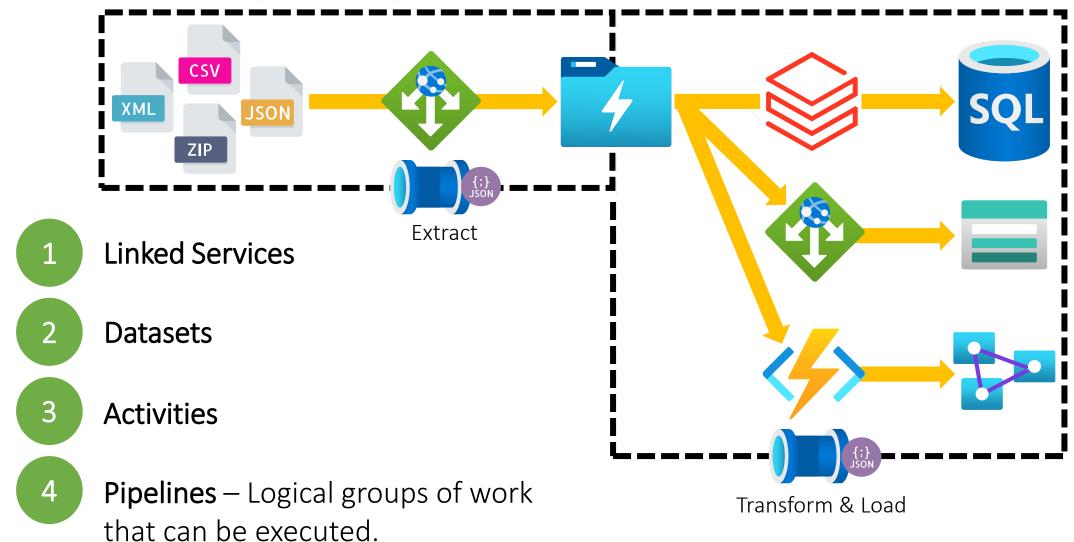
- 1 Linked Services
- 2 Datasets
- Activities What do we want to happen when we invoke a Linked Service?
 With what conditions?

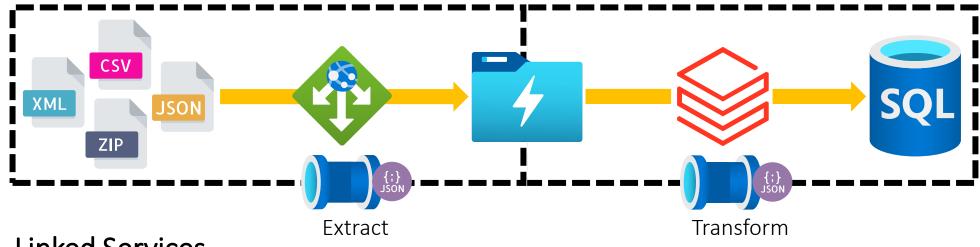




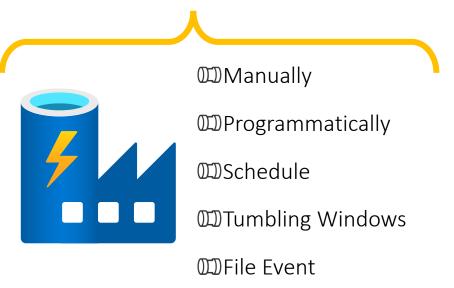
- 1 Linked Services
- 2 Datasets
- 3 Activities
- 4 Pipelines Logical groups of work that can be executed.

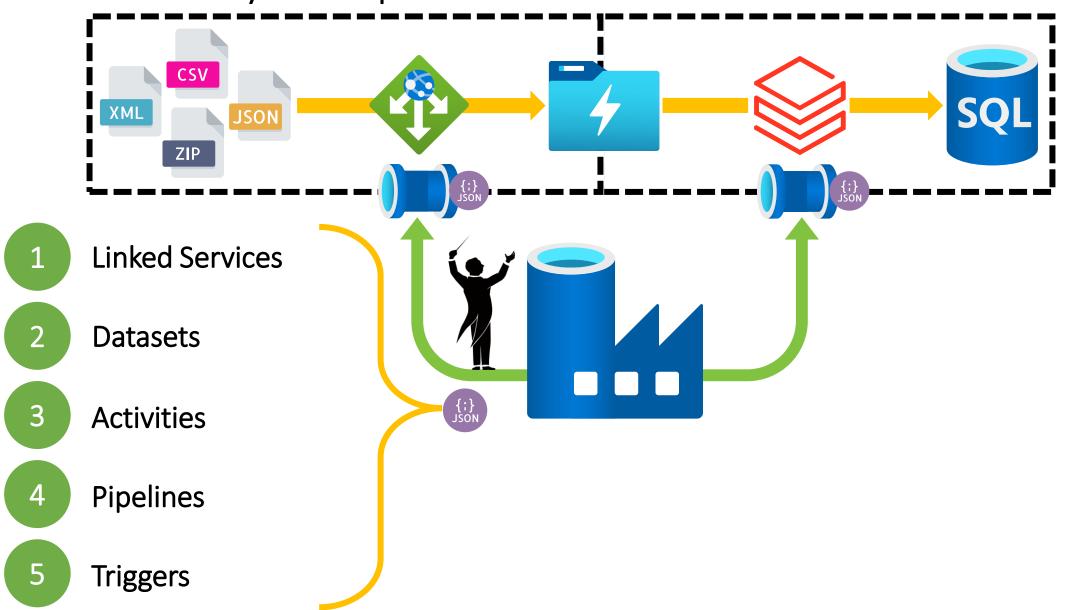






- 1 Linked Services
- 2 Datasets
- 3 Activities
- 4 Pipelines
- 5 Triggers Telling our when pipelines to run.

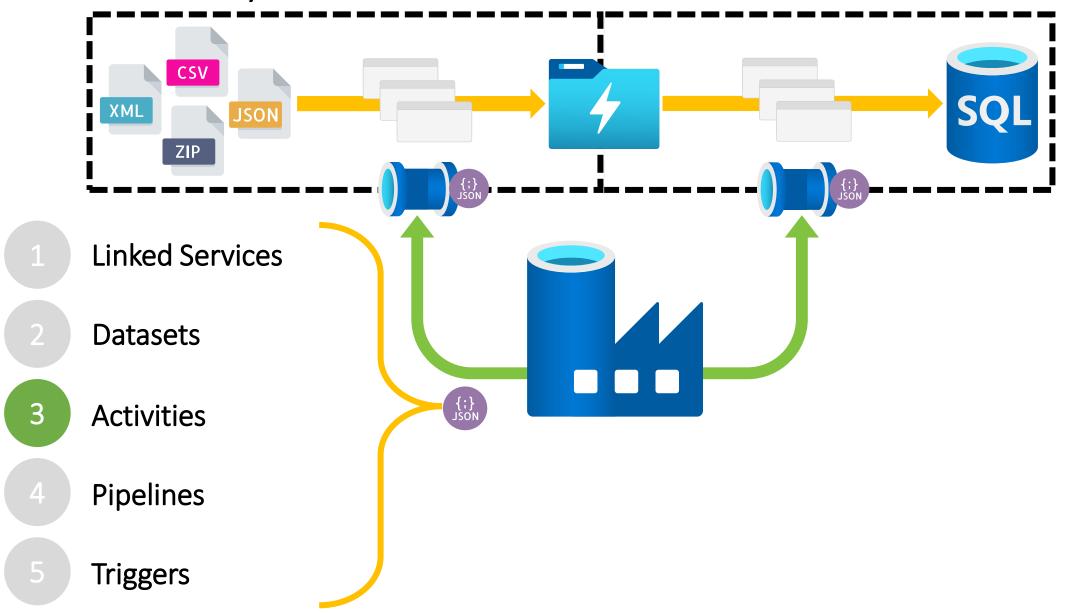




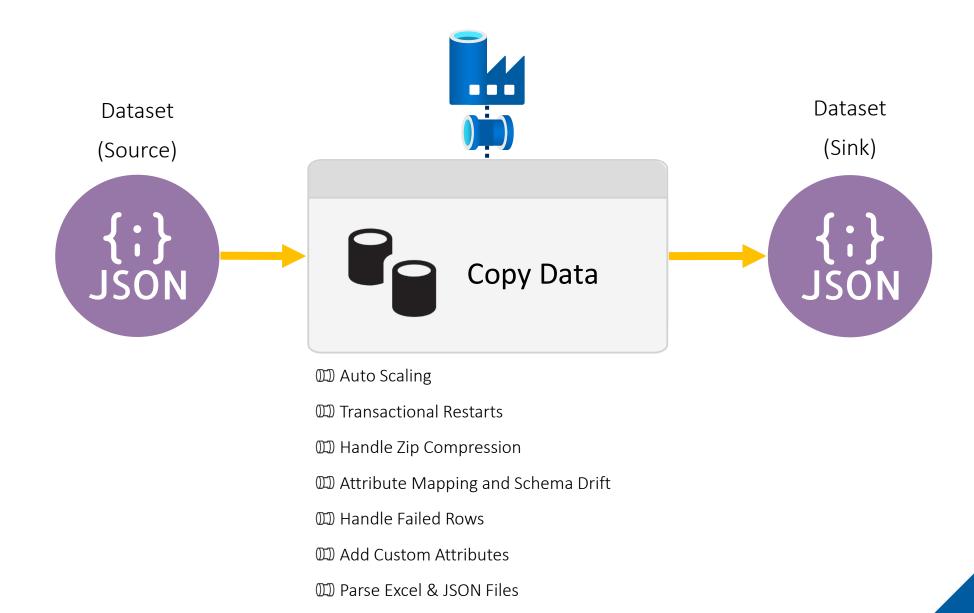
Common Activities

```
SELECT TOP 5
    [ActivityName],
    [Inputs],
    [Outputs],
    [Details]
FROM
    [metadata].[AdfActivities]
WHERE
    [Notes] = 'Pauls Favourites';
```

Data Factory Common Activities

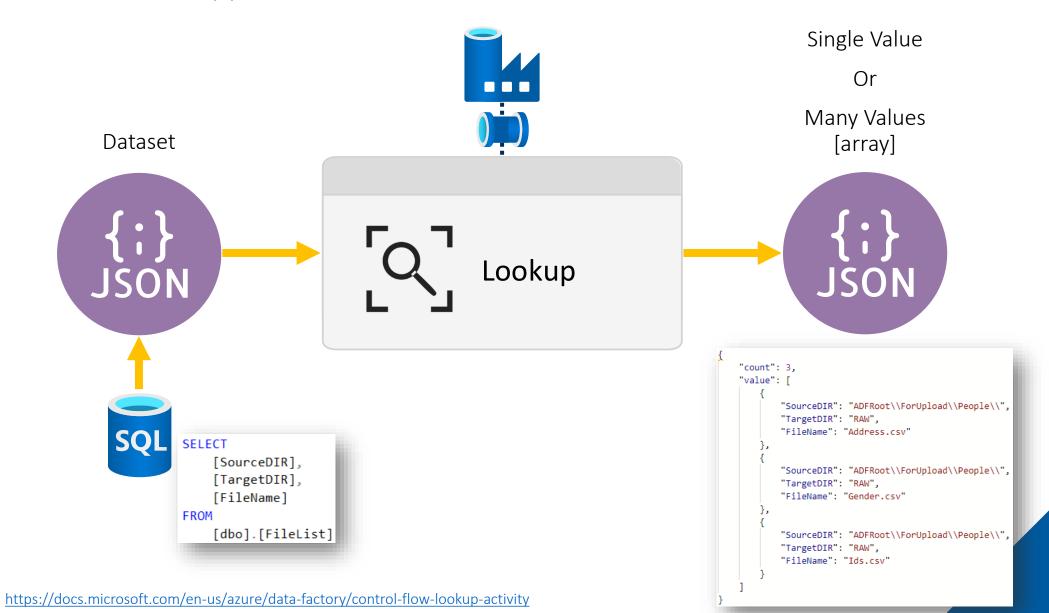


Copy



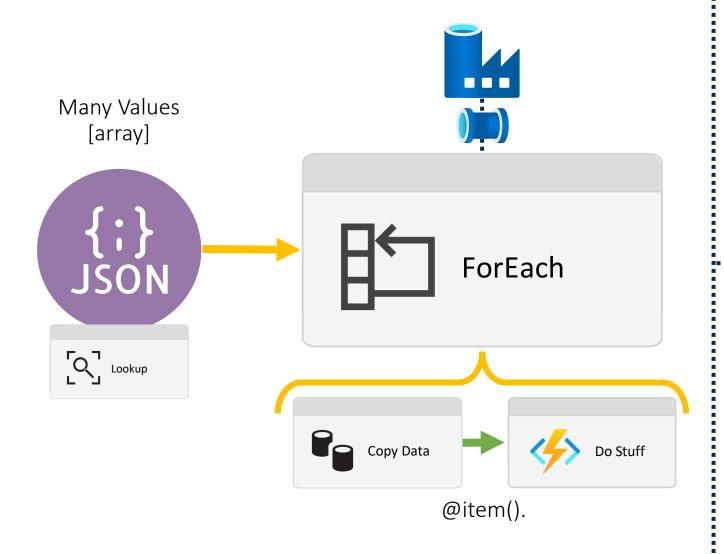
Lookup

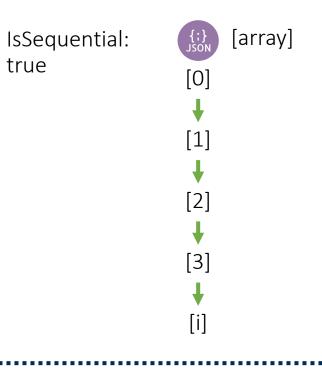
Get value to support other control flow activities

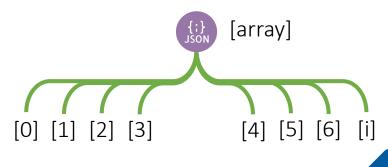


ForEach

Scaling Out Control Flow Activities



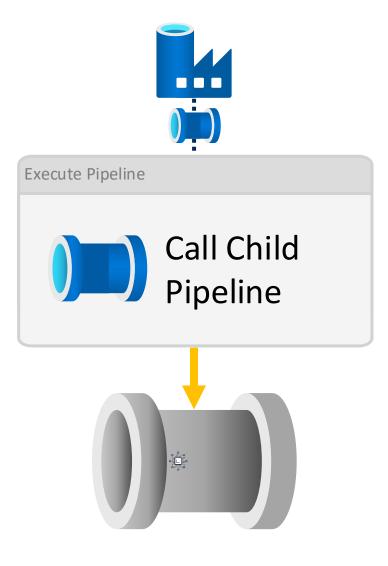




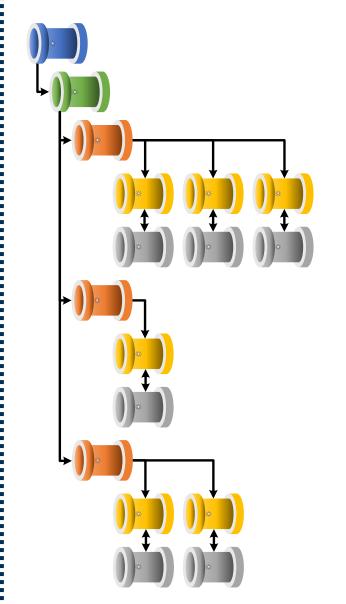
Batch Count Default: 20

Batch Count Max: 50

Execute Pipeline

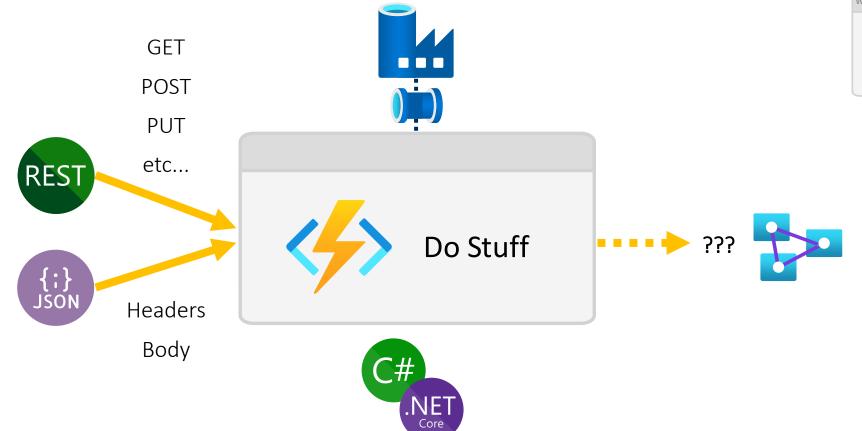


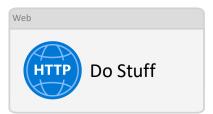


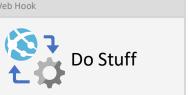


Azure Function

Extend Data Factory with Rest Calls







Custom

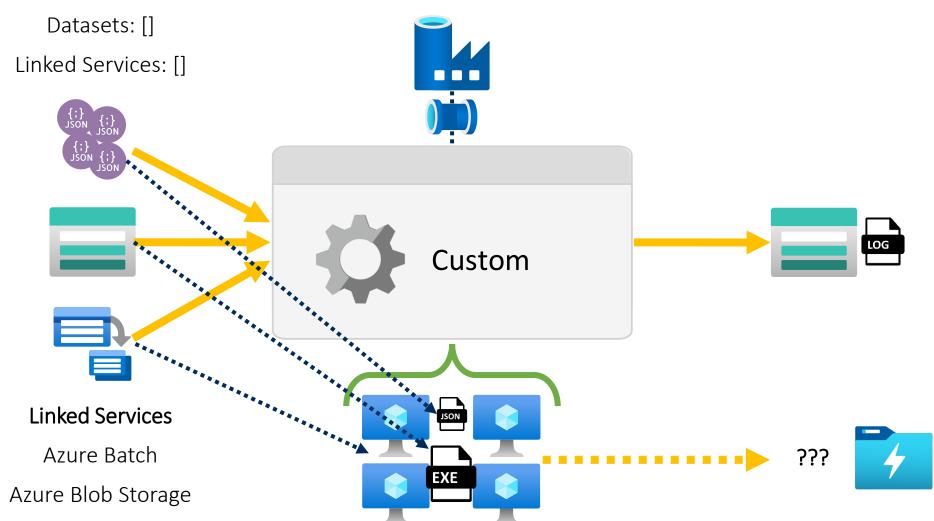
Extend Data Factory with Custom Code



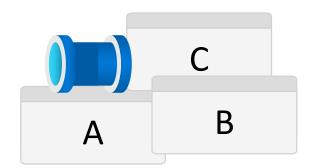
Creating a Custom Activity

https://mrpaulandrew.com/2018/11/12/c reating-an-azure-data-factory-v2-customactivity/

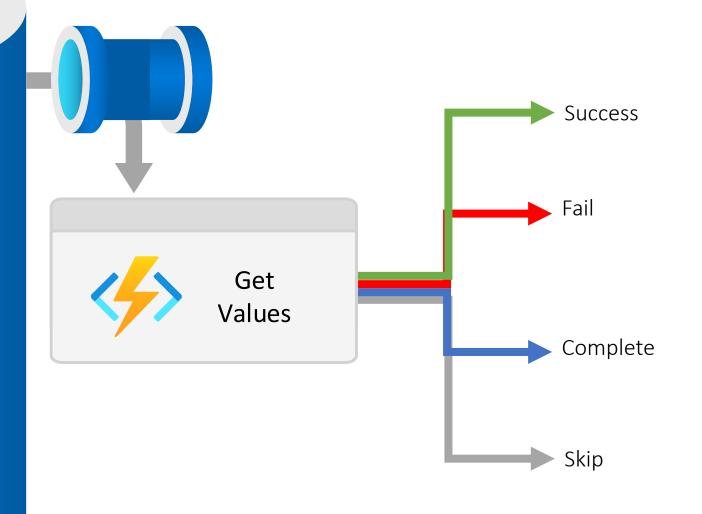
References Objects



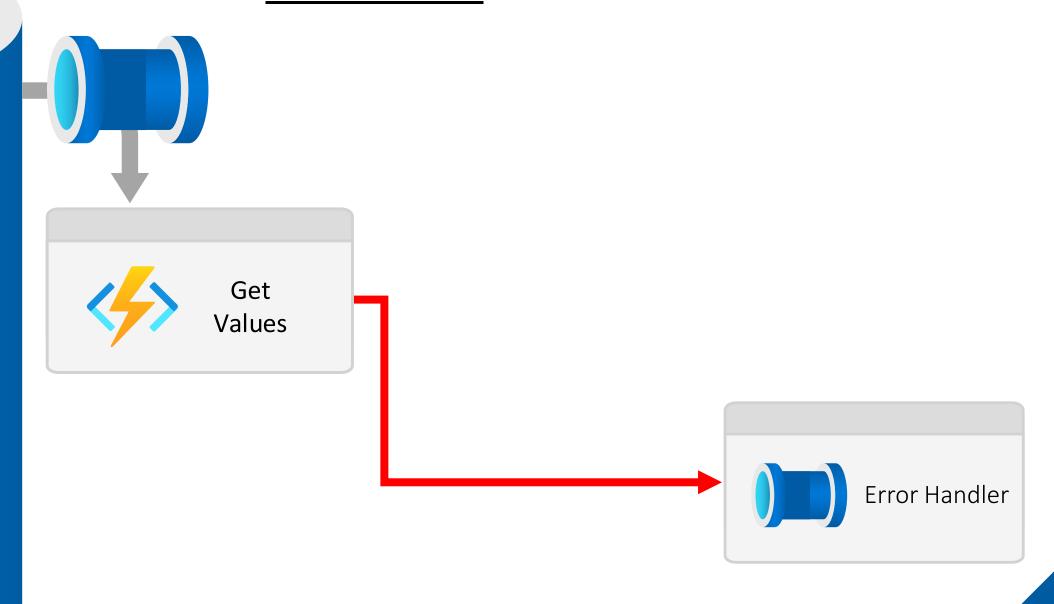
Execution Dependencies



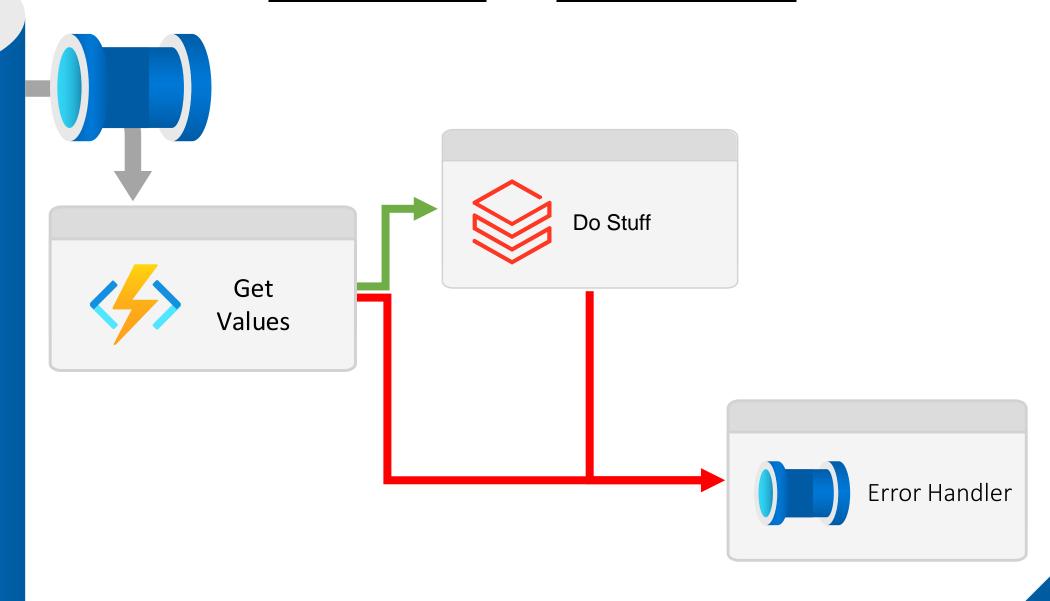
Execution Dependency Options



Execution On Failure

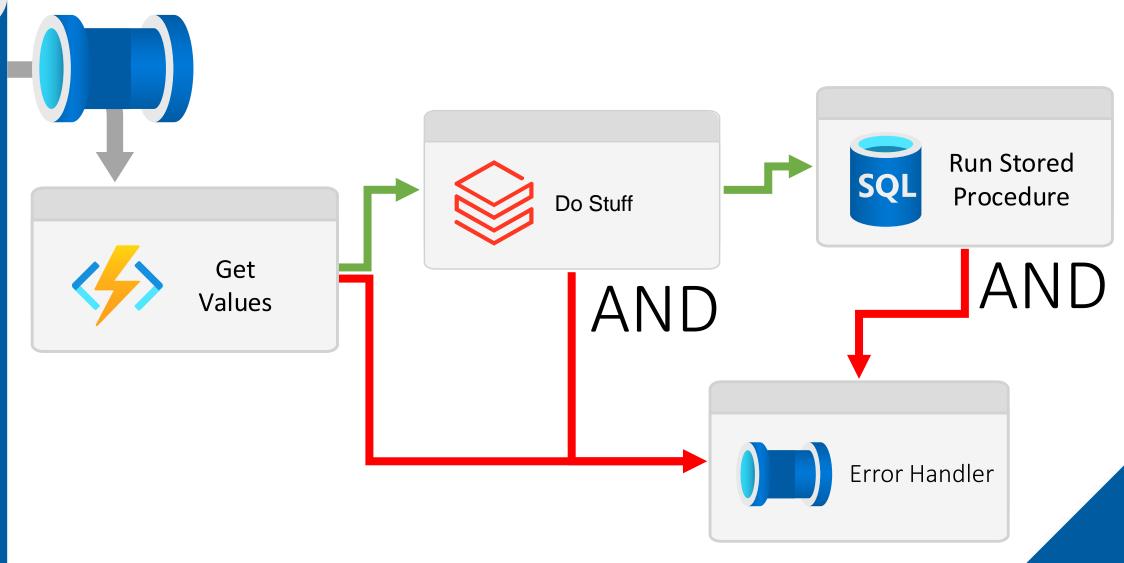


Execution On Failure or On Success

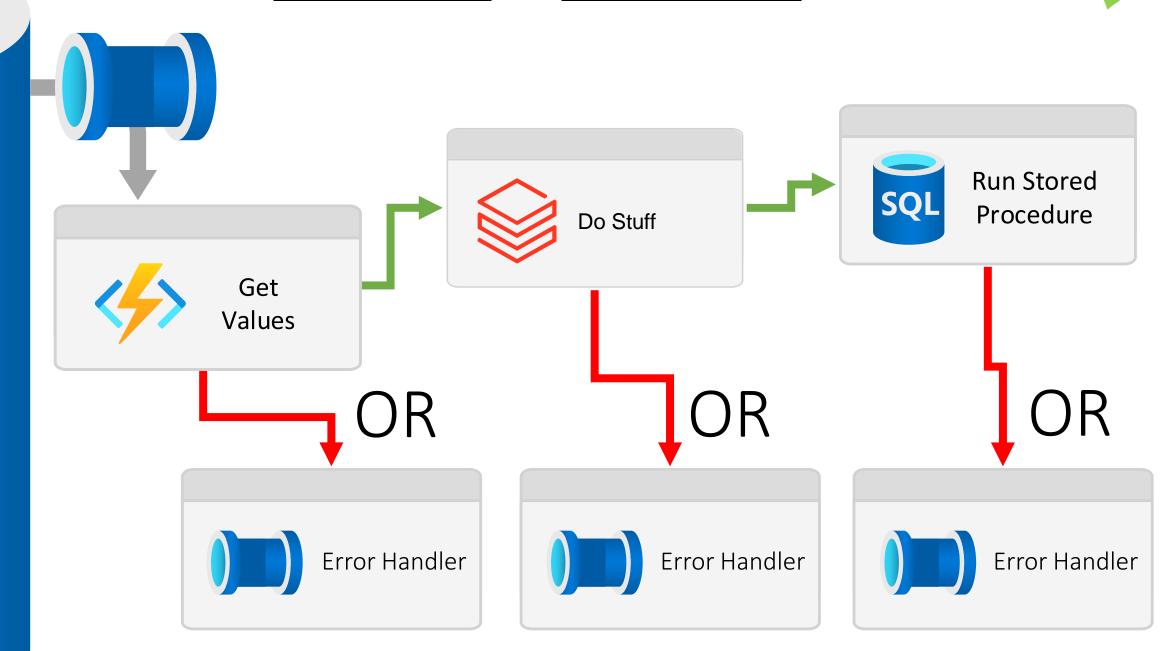


Execution On ???

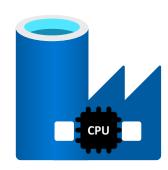




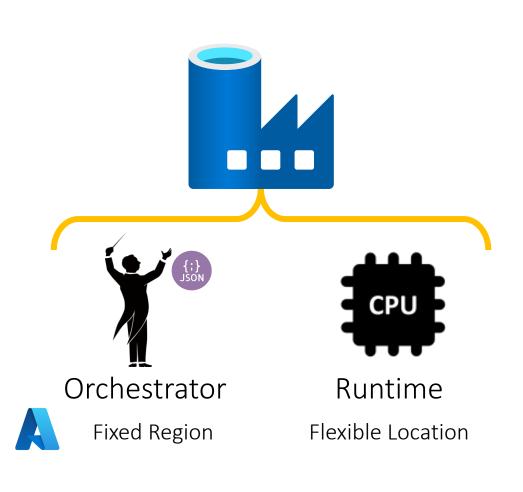
Execution On Failure or On Success

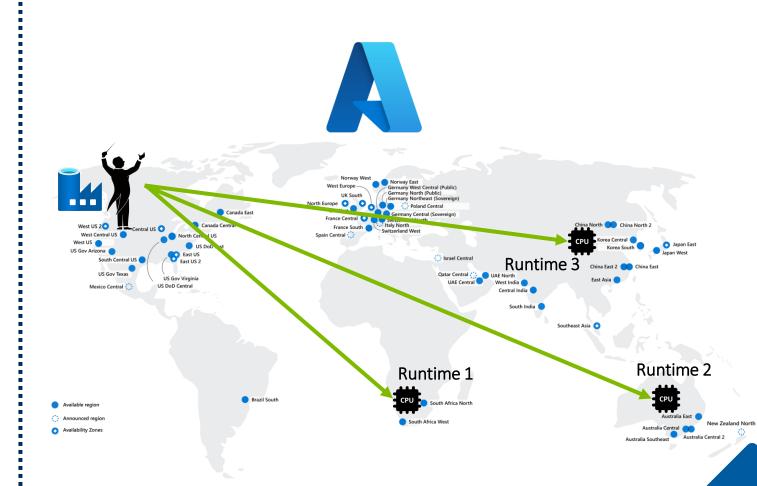


Integration Runtimes

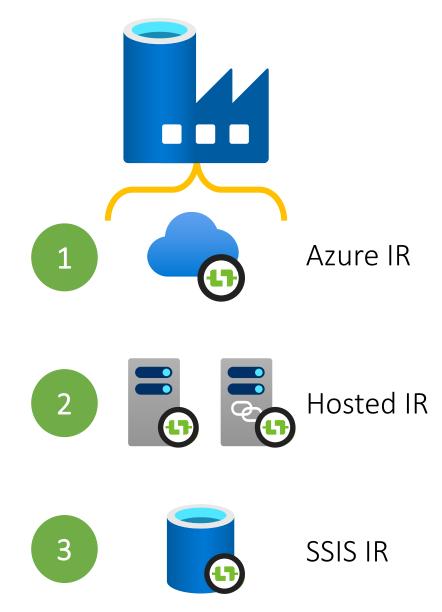


What is an Integration Runtime?

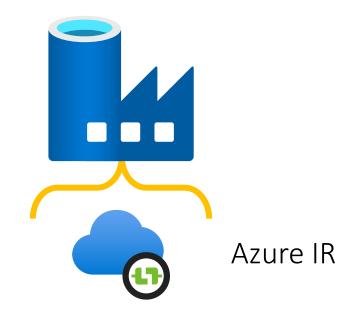




What can an Integration Runtime do?



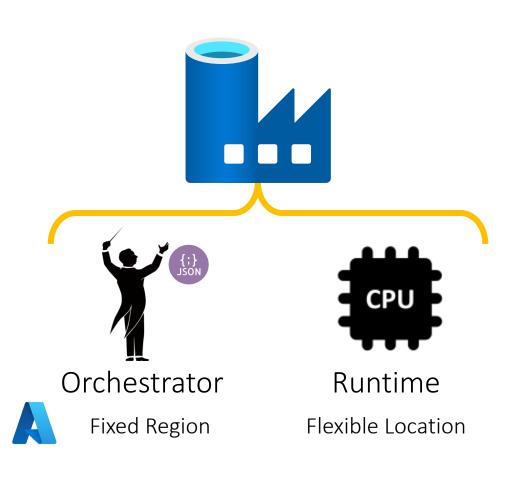
Azure Integration Runtime

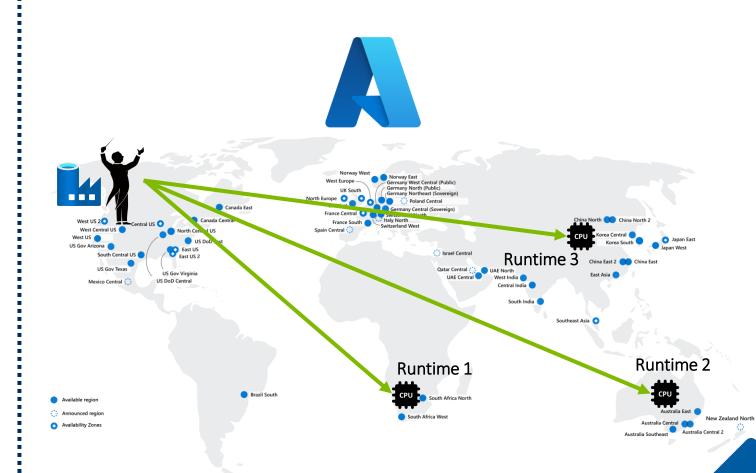




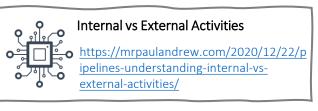


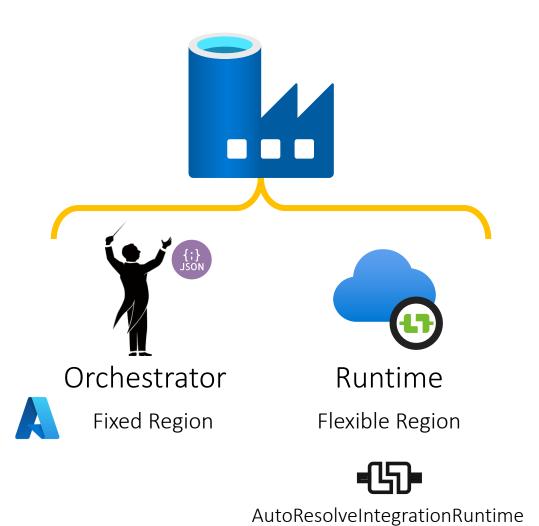
Azure Integration Runtime





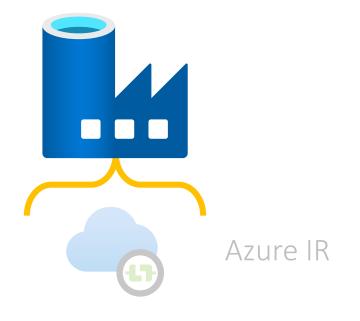
Azure Integration Runtime







Hosted Integration Runtime





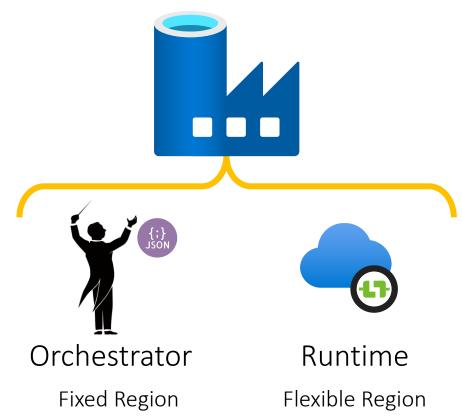


Hosted Integration Runtime

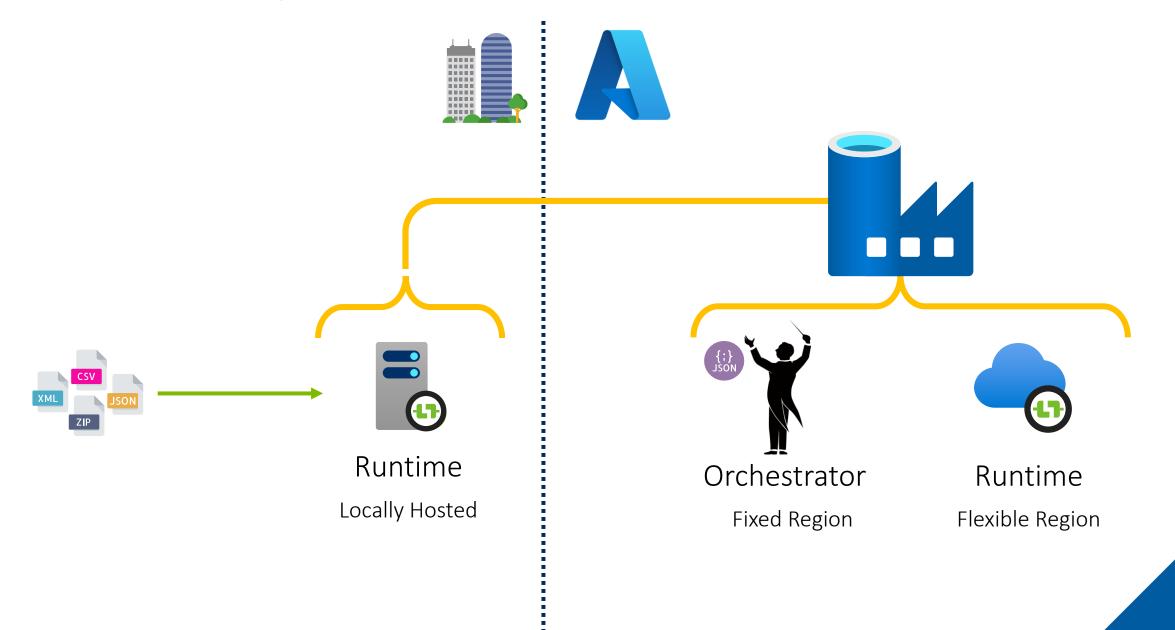




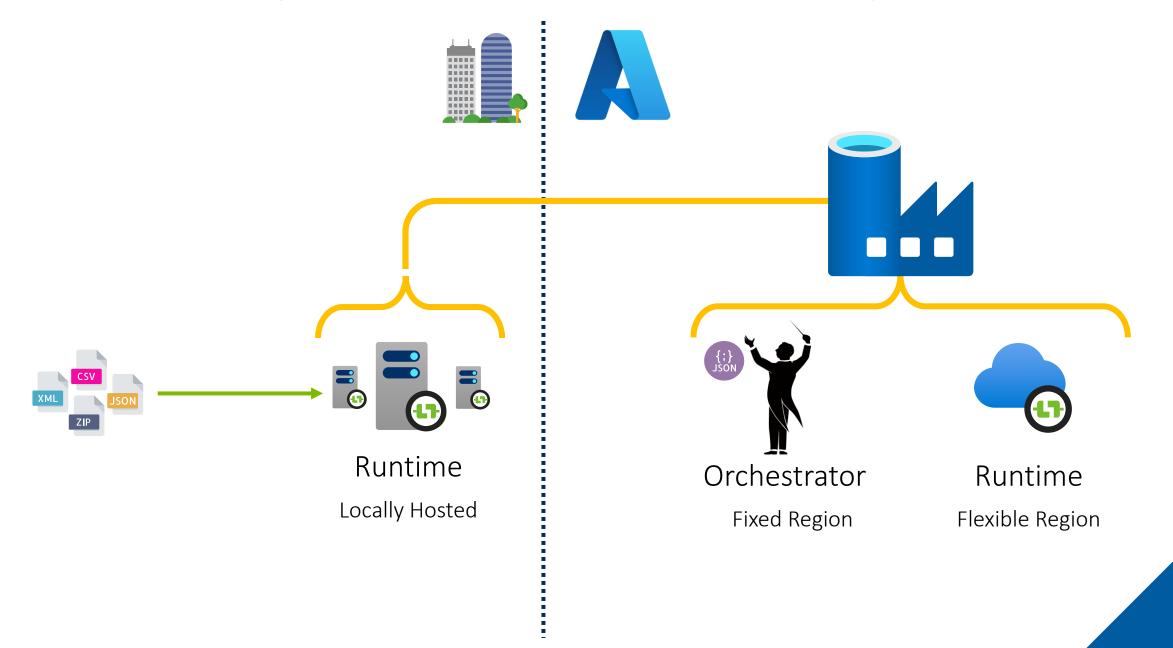




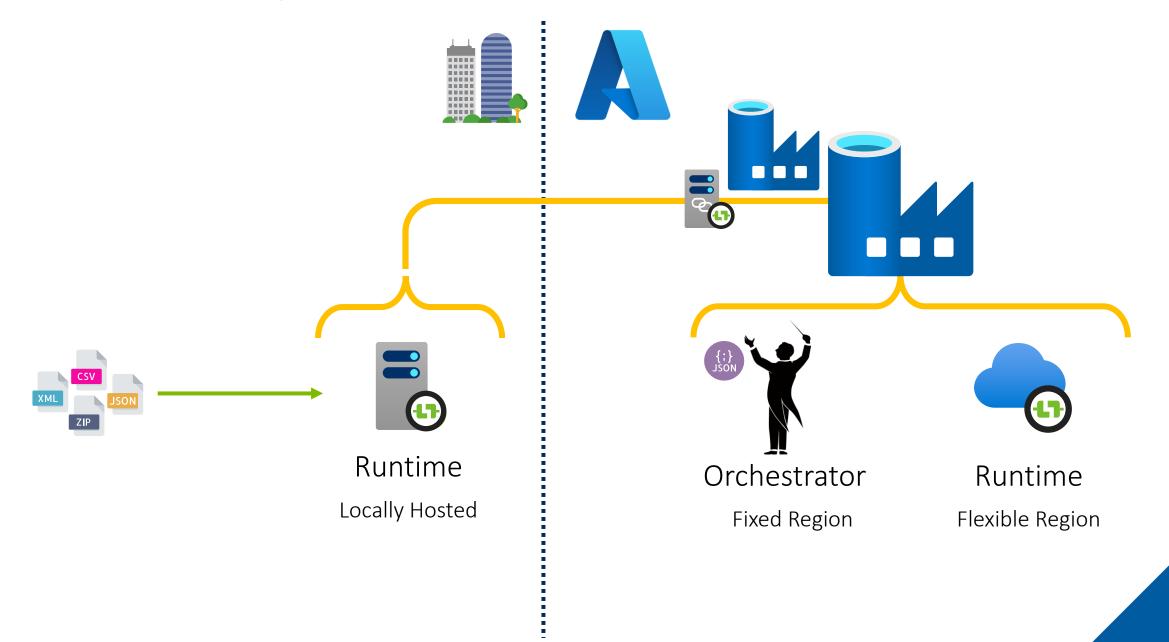
Hosted Integration Runtime



Hosted Integration Runtime – Secondary Nodes



Hosted Integration Runtime – Linked

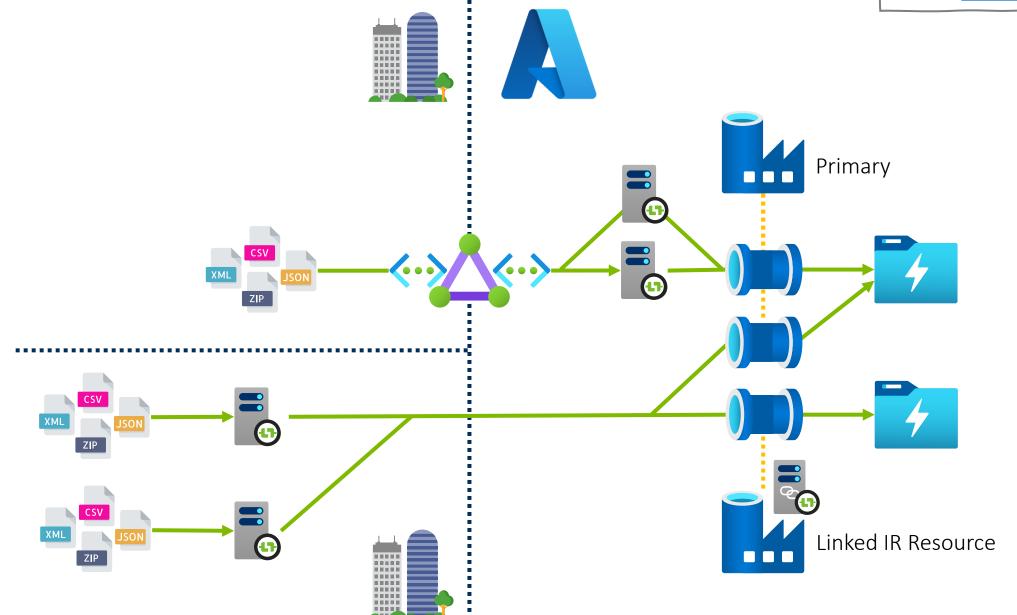


Hosted IR Advanced Patterns

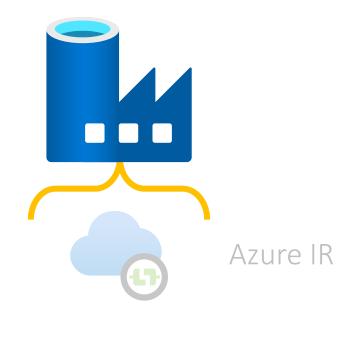


Scaling Azure Data Integration Pipelines

caling-azure-data-integration-pipelines-decoupling-data-extract-and-transform/



SSIS Integration Runtime





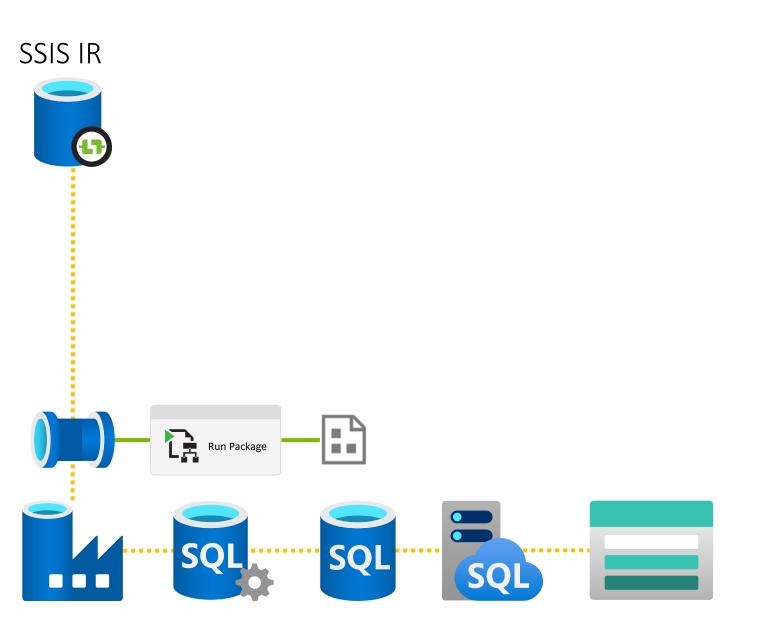


Running an SSIS Package in Azure

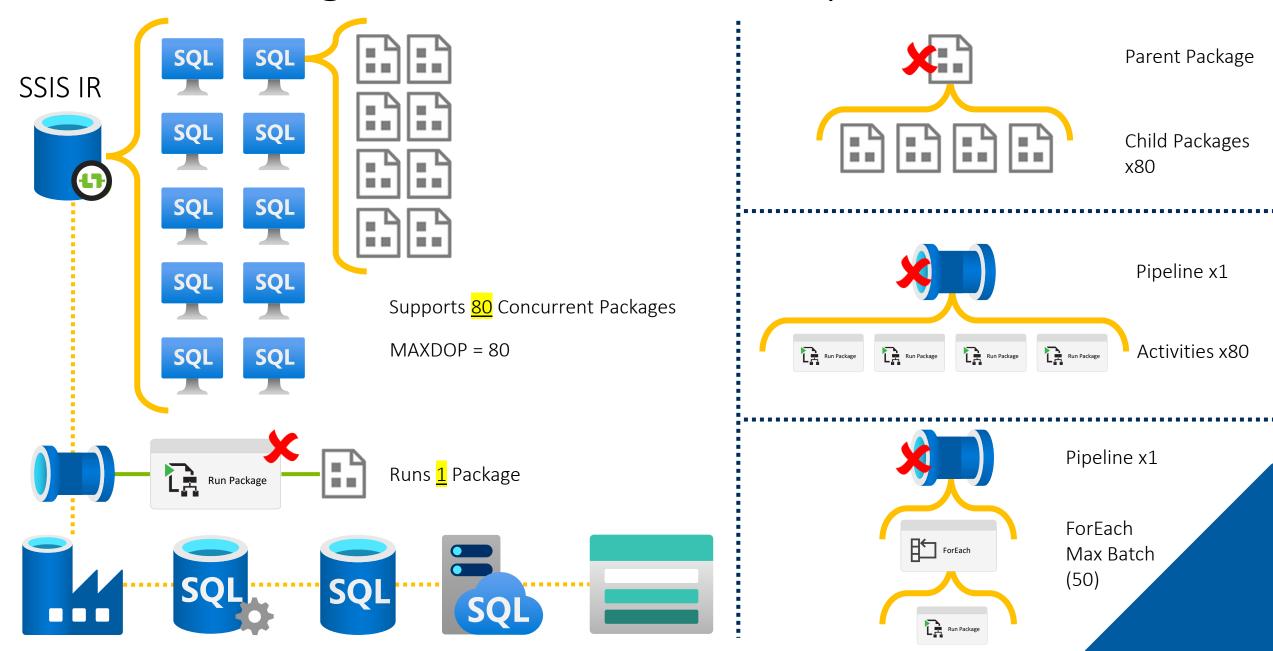


SSIS IR

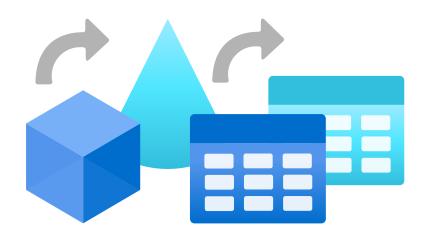
Running an SSIS Package in Azure



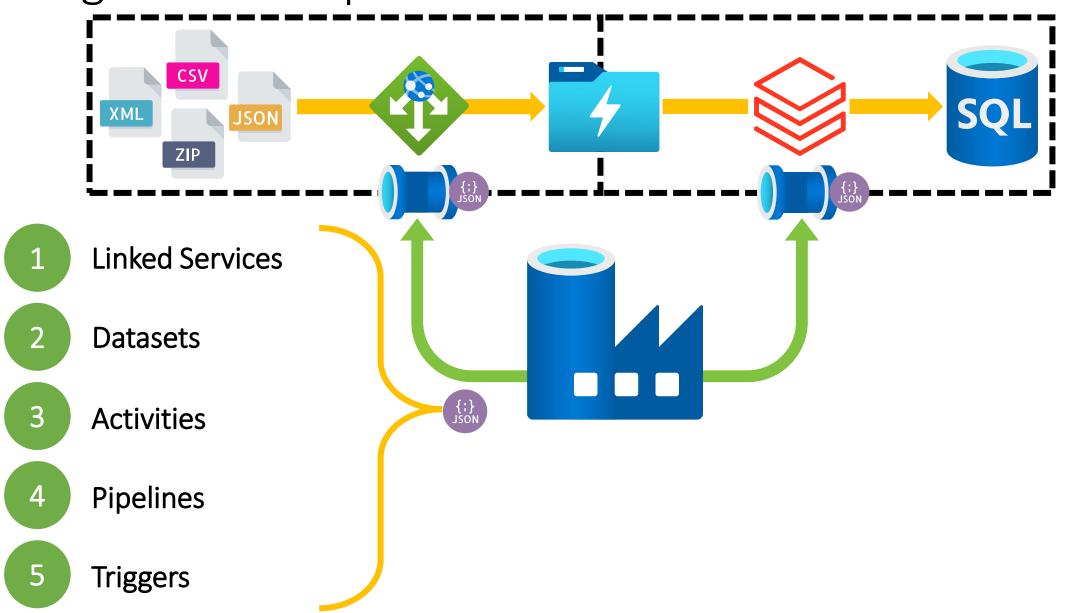
Problem: Using All Of The SSIS IR Compute



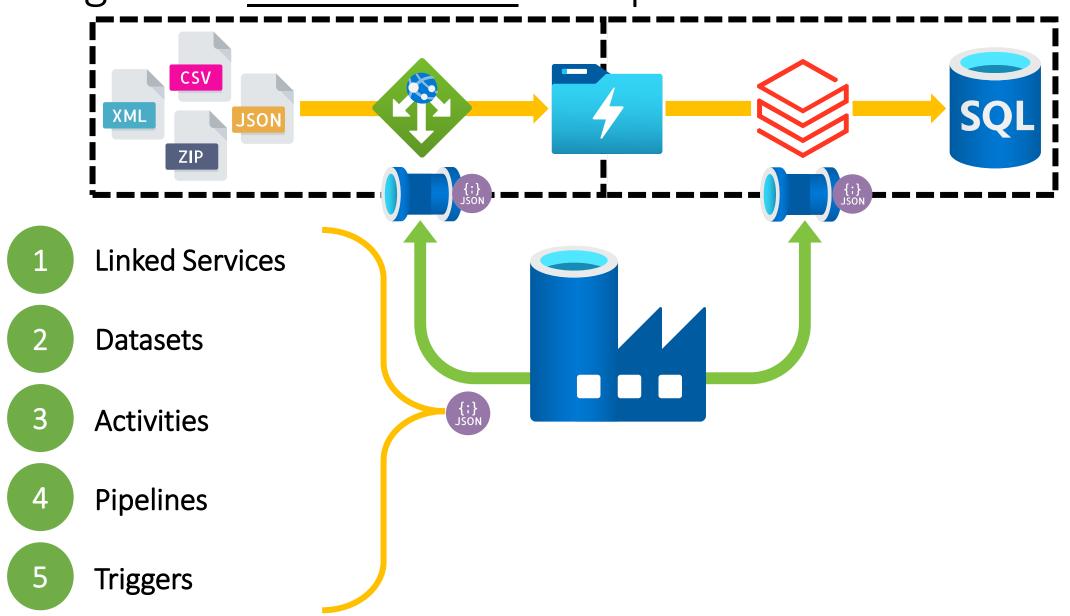
Data Flows

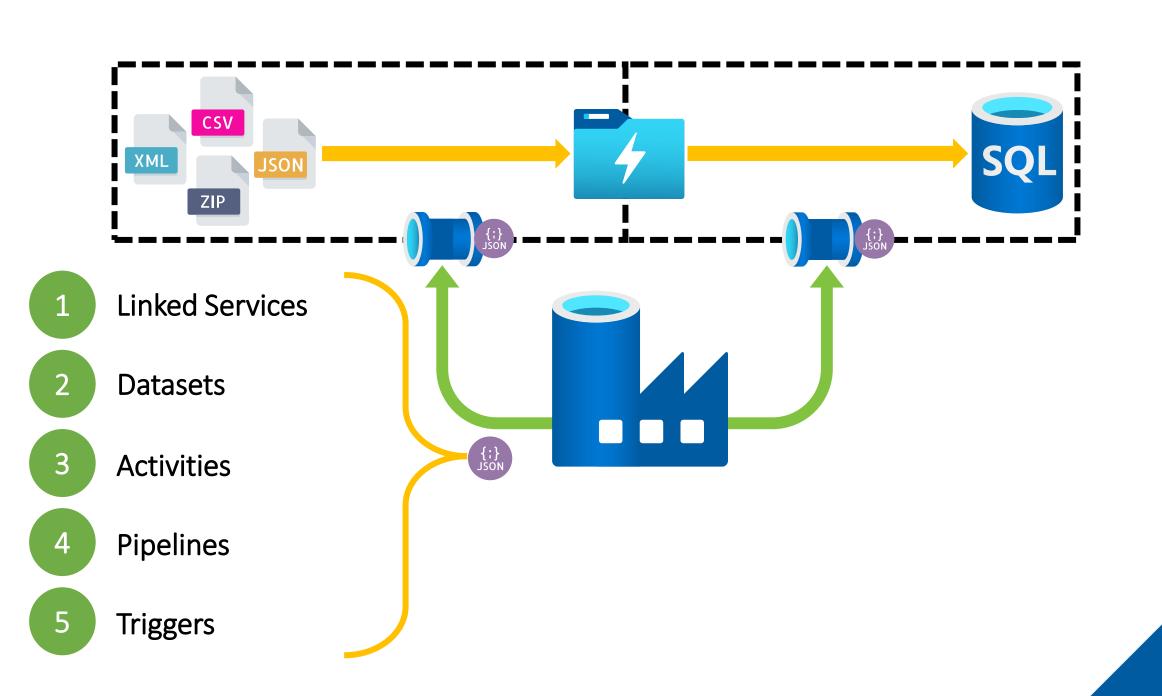


Integration Components

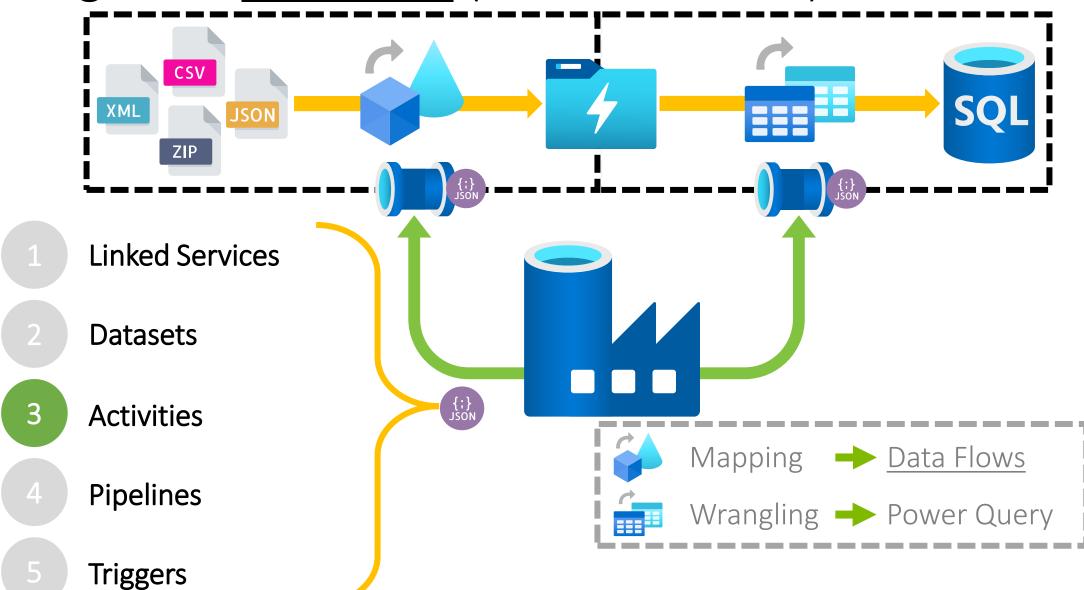


Integration Control Flow Components

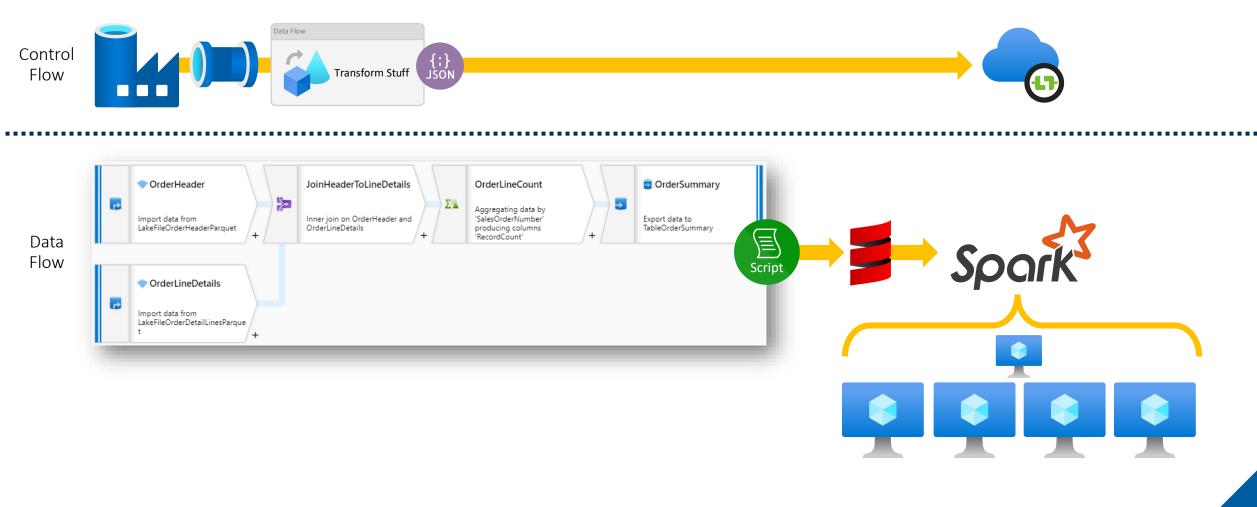




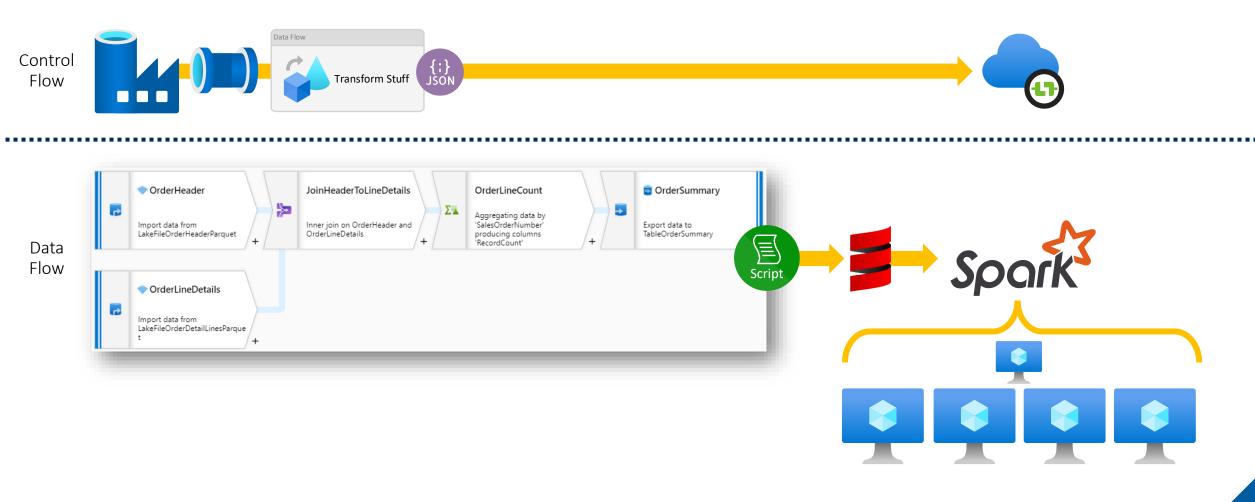
Integration Data Flow (Transformation) Activities



What is a Mapping Data Flow?



Q: What is a Mapping Data Flow?



A: Graphic no low/low code data transformation tool that sits on top of Apache Spark.

Data Flows – Inputs & Outputs

Source & Sink



Linked Services



















Dataset 📰















Source Types

Inline







Data Flows — Transformations



New Branch

Conditional Split



Derived Column



Flatten



Filter



Parse



Sort



Alter Row



Join

Exists

Union



Select



Aggregate



Surrogate Key



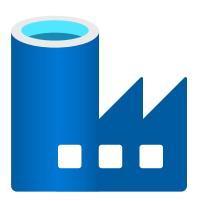
Pivot



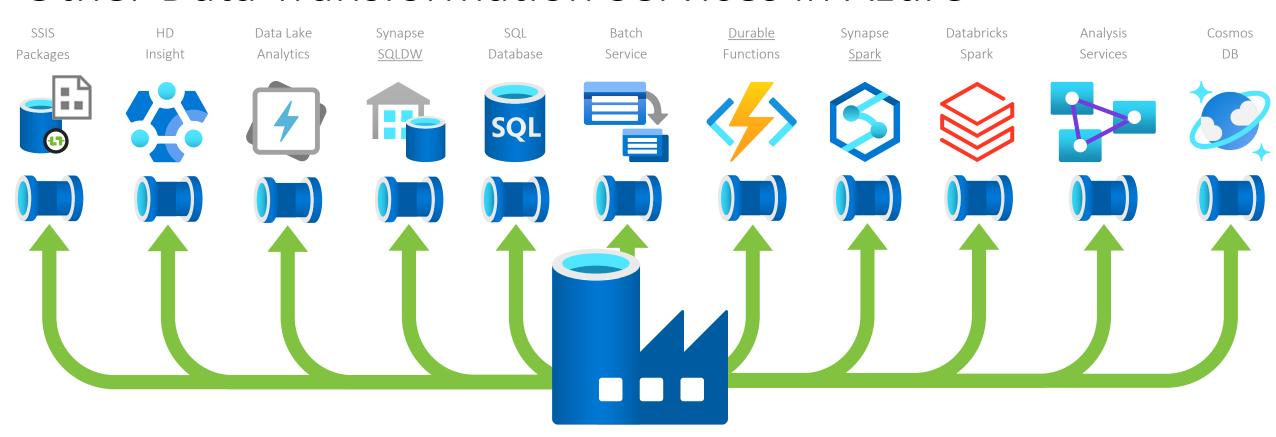




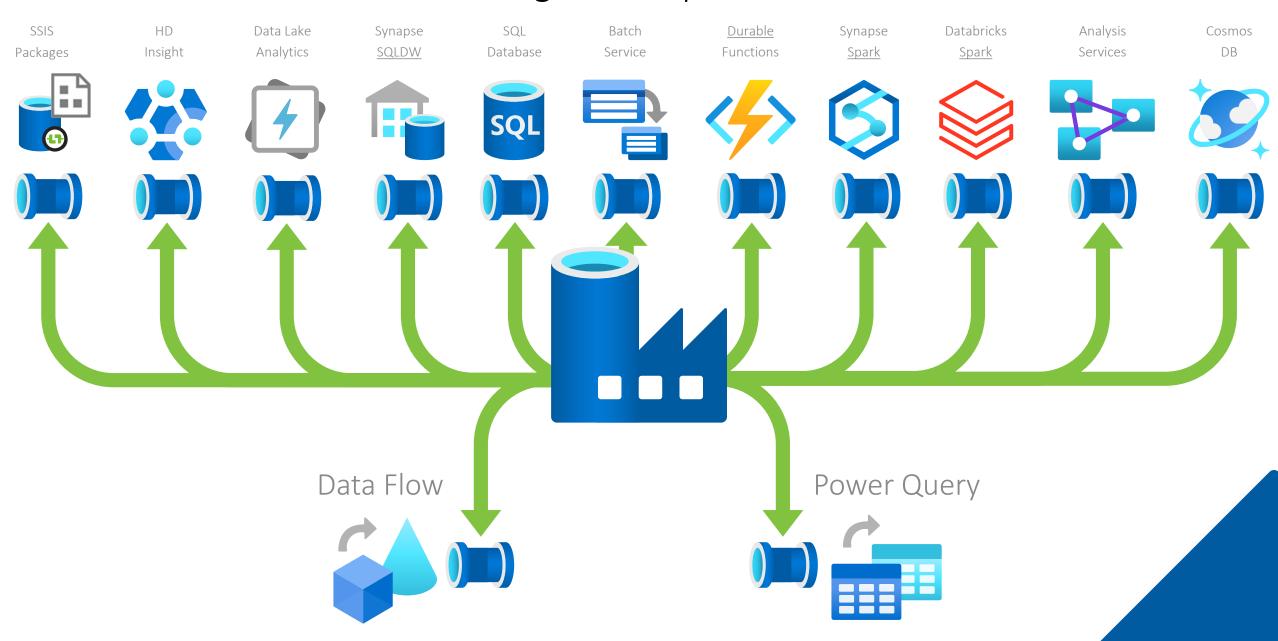
Key
Input & Output Modifiers
Schema Modifiers
Formatters
Row Modifiers



Other Data Transformation Services in Azure



When Should We Use These Integration Pipeline Transformation Activities?



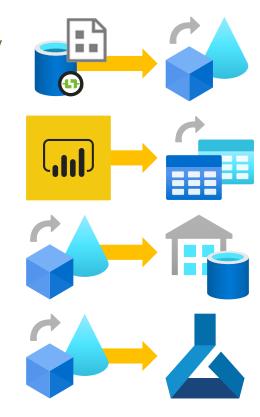
Use Cases

SSIS developers who are transferring existing skills to cloud native technologies have a very low barrier to entry and don't need to worry about distributed compute to get started.

Data engineering made easy for the <u>power users who has grown out of Power BI</u> following a series of Data Lake exploration sessions.

Data insight teams needing to do <u>rapid prototyping and data warehouse loading</u> within a single Azure Resource making deployments simple and release cycles short.

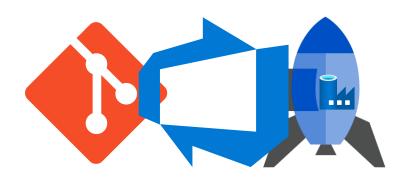
Simpler and quicker data wrangling for <u>data scientists</u> that want to <u>quickly prepare multiple raw</u> <u>datasets</u> ready for model training and testing, also with the ability to use large amounts of compute.

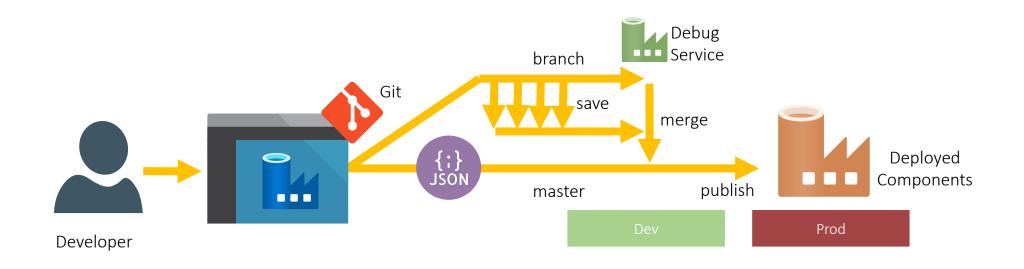


Data Flows used to deliver all data transformation workloads as part of a end to end cloud based data analytics/warehouse solution.

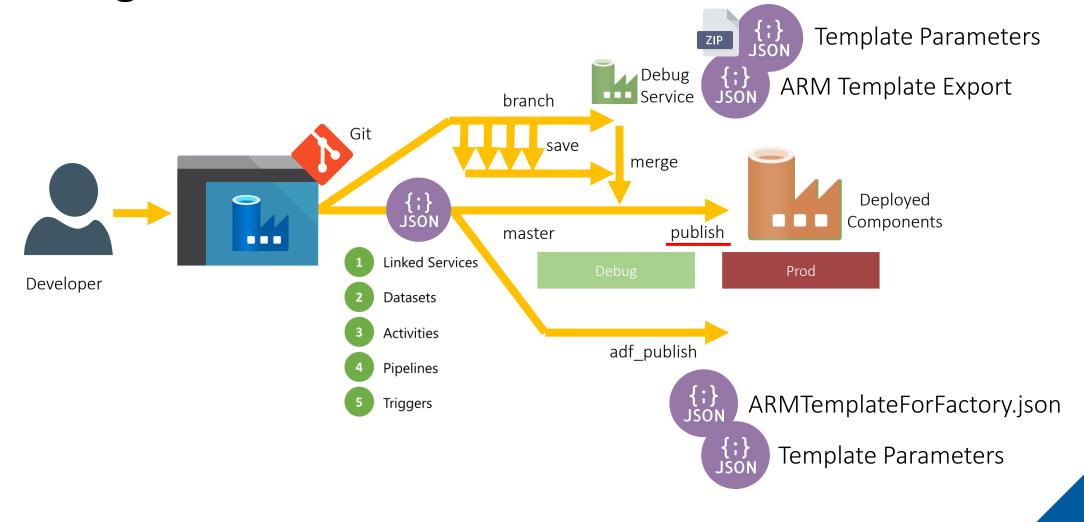
Data Flows script dynamically generated from external metadata and injected into like we once did with BIML for SSIS packages.

Source Control & Deployments

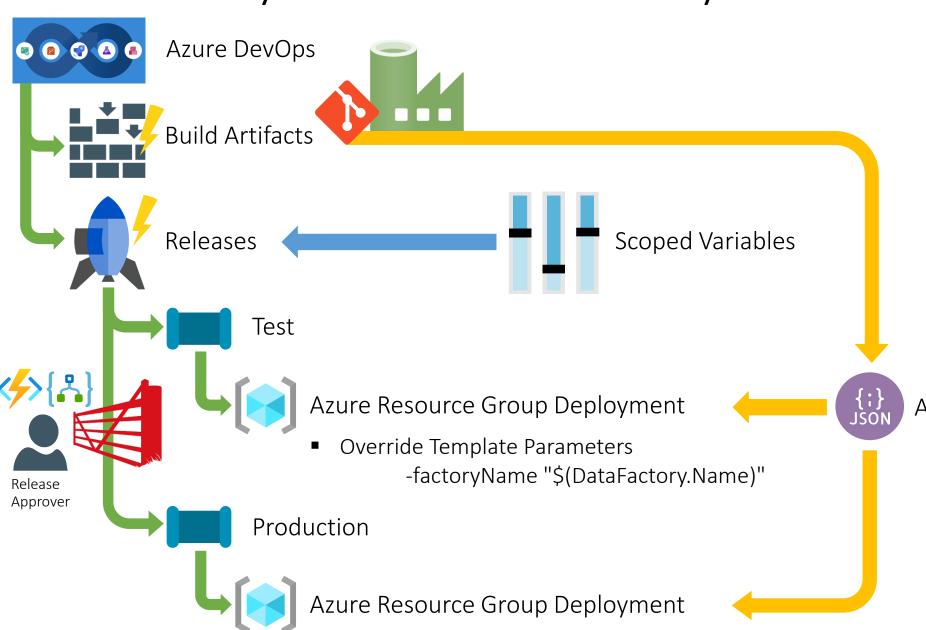


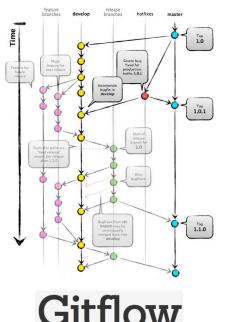


Getting Our ADF Source Code



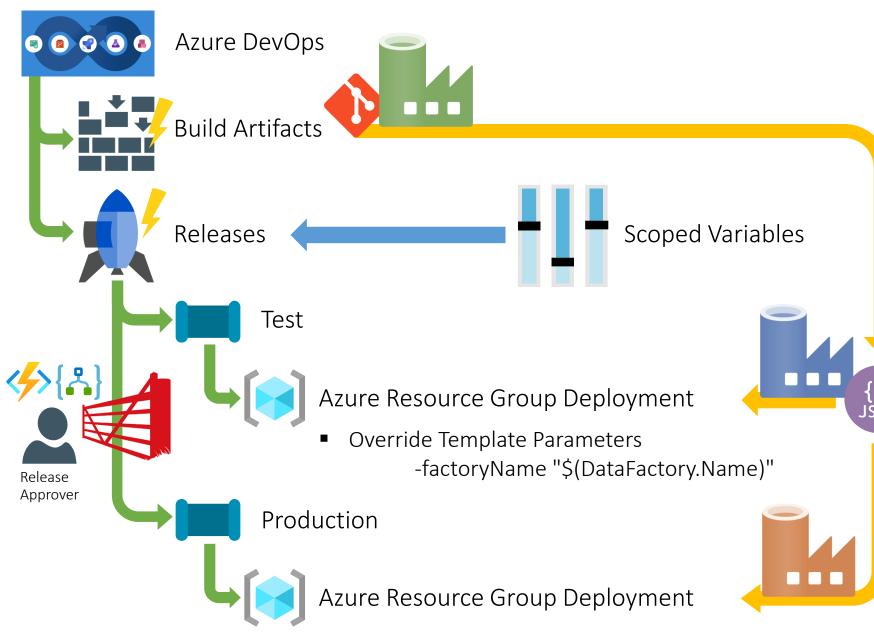
Data Factory Continuous Delivery





ARMTemplateForFactory.json

Data Factory Continuous Delivery - Simple



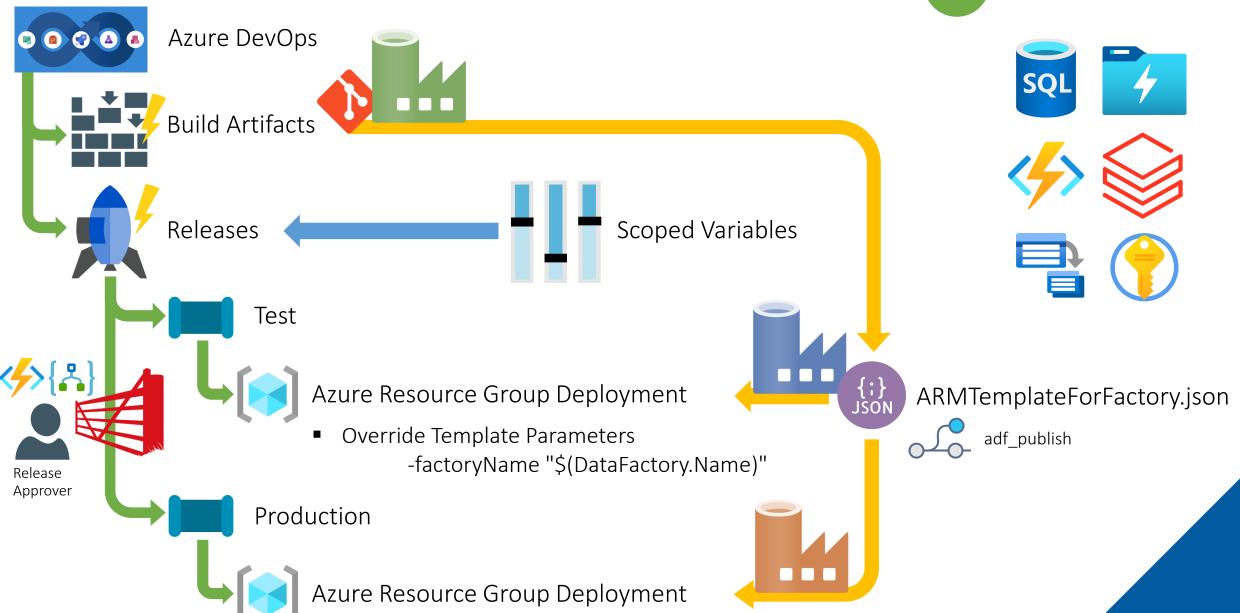
- 1 Linked Services
- 2 Datasets
- 3 Activities
- 4 Pipelines
- 5 Triggers

ARMTemplateForFactory.json

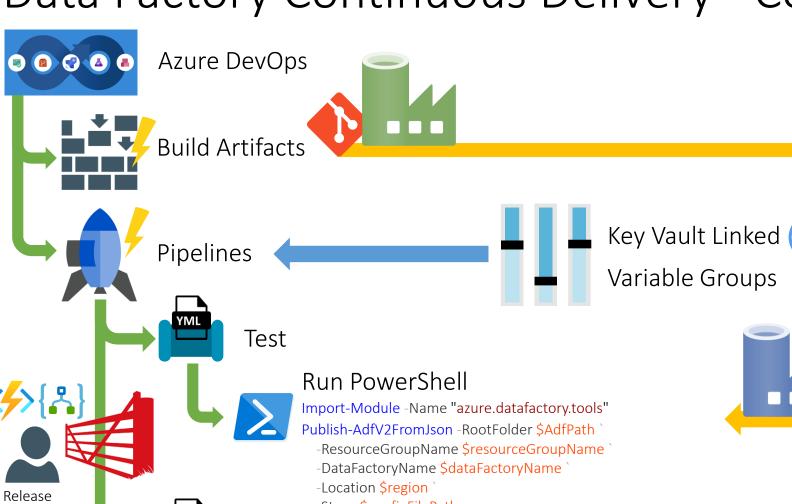
Data Factory Continuous Delivery - Simple



Linked Services



Data Factory Continuous Delivery - Complex





Linked Services



Datasets



Activities



Pipelines



Triggers

linkedservices.json pipelines & activites.json datasets.json triggers.json



{release} / {feature} / {tag}

-Stage \$configFilePath

Production



Approver

Publish Azure Data Factory

Data Factory DevOps Story Summary

What is your code branching strategy?

Which source control tool to use?

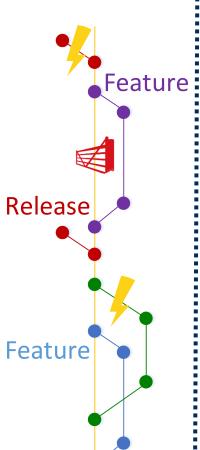
How many environments do we want?

What deployment method do we want to use?

What artifacts are we going to use?...

OR

How much control do you want?



Master

















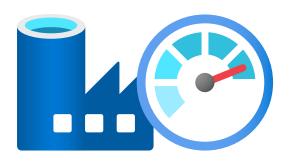




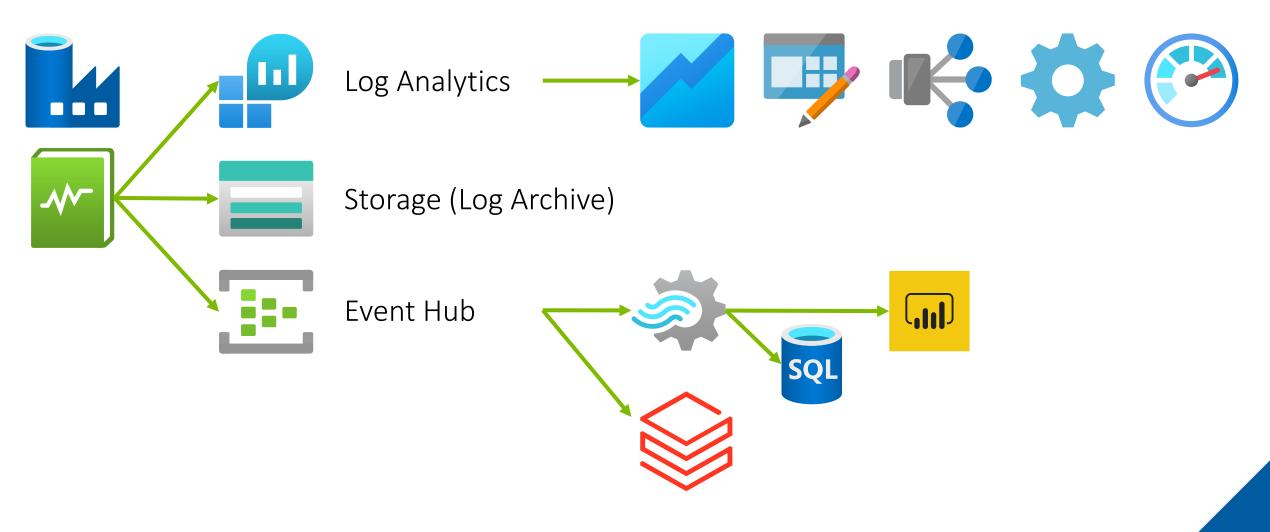


linkedservices.json pipelines & activites.json datasets.json triggers.json

Monitoring & Logging



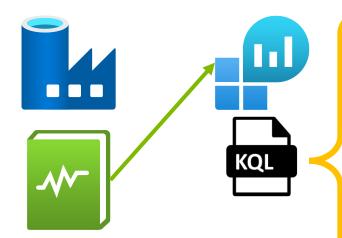
Diagnostic Settings

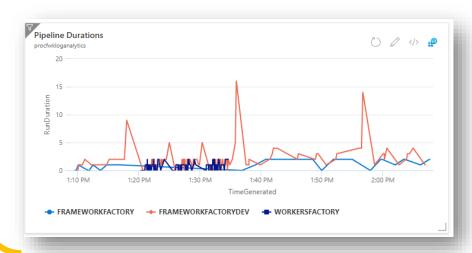


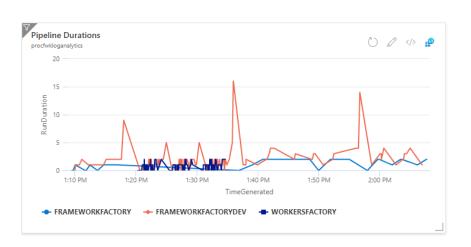
Diagnostic Settings



Using Log Analytics







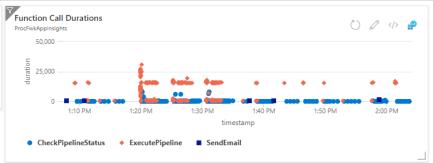


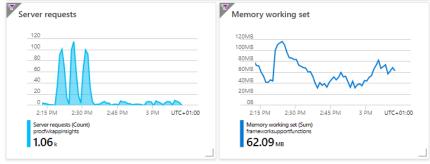


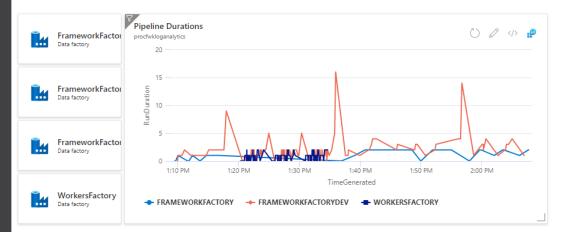












2:30 PM

2:45 PM

Compute utilization

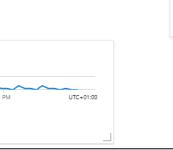
DTU percentage (Max)

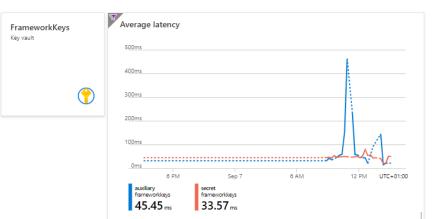
13 %

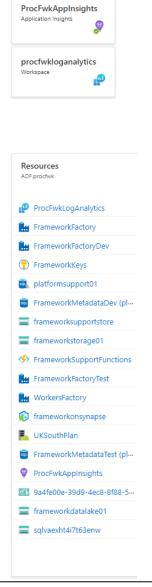
platformsupport01/frameworkmetadatadev

FrameworkMetadat...

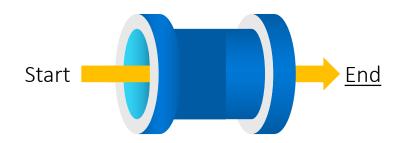
SQL database Online



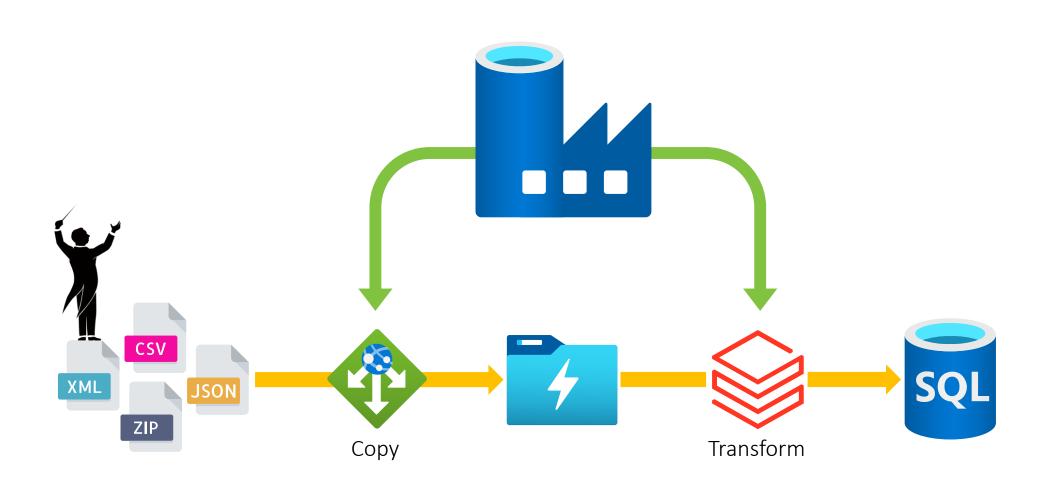




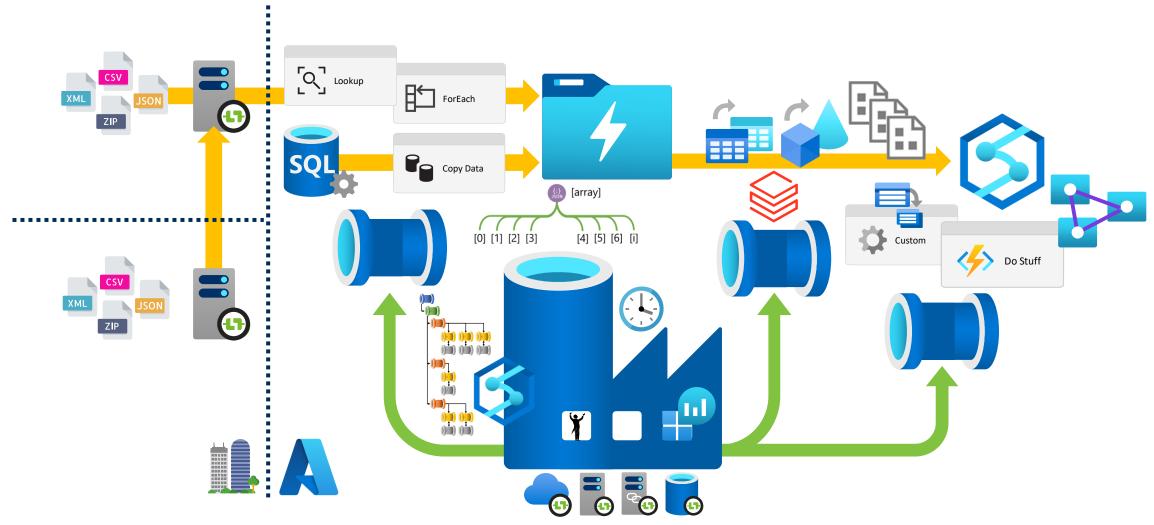
Conclusions



What is Azure Data Factory (ADF)?



What are Azure Data Factory Integration Pipelines?



- 1. A complete Microsoft Azure integration tool.
- 2. Orchestrator of our <u>Control Flow</u> operations with scale out Activities.
- 3. Orchestrator of our <u>Data Flow</u> transformations using cloud native services.
- 4. The scheduler of solutions using a variety of Pipeline Triggers and dynamic frameworks.

What Next?

Best Practices for Implementing Azure Data Factory



- D Environment Setup
- Multiple Data Factory Instance's
- Deployments
- DD Automated Testing
- Maming Conventions
- D Pipeline Hierarchies
- DD Pipeline & Activity Descriptions
- M Annotations
- DD Linked Service Security via Azure Key Vault
- Security Custom Roles
- Dynamic Linked Services

- Generic Datasets
- Metadata Driven Processing
- D Parallel Execution
- M Hosted Integration Runtimes
- Azure Integration Runtimes
- Wider Platform Orchestration
- Custom Error Handler Paths
- Monitoring via Log Analytics
- DD Timeouts & Retry
- Service Limitations
- **W** Using Templates
- Documentation

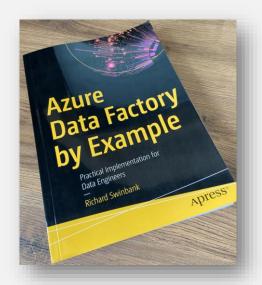


Best Practices for Implementing ADF

https://mrpaulandrew.com/2019/12/18/best-practices-for-implementing-azuredata-factory/

What Next?

Azure Data Factory by Example



Author: Richard Swinbank @RichardSwinbank

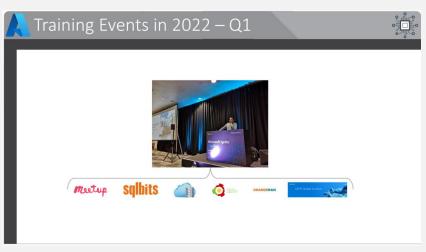
Technical Reviewer: Paul Andrew

ISBN-13978-1484270288

What Next?

More Azure Data Platform Training and Sessions With Me!





Thank you for listening...

Paul Andrew





Blog: mrpaulandrew.com

YouTube: c/mrpaulandrew

Email: paul@mrpaulandrew.com

Twitter: @mrpaulandrew

LinkedIn: In/mrpaulandrew

GitHub: github.com/mrpaulandrew

/CommunityEvents /ContentCollateral