# GROUPBY CODE OF CONDUCT

**The Quick Version**

We are dedicated to a harassment-free experience for everyone, regardless of who you are and what makes you *you.* We recognize the right of any individual to attend and participate. Anyone. This is included but not limited to gender identity and expression, sexual orientation, disability, physical appearance, body size, race, religion, or any other classification, affiliation, or label.

**We do not tolerate harassment in any form.** For the duration of your engagement with GroupBy and its programs, you are expected to act appropriately and to adhere to this Code of Conduct. This includes conduct in-person and online, at the conference itself, as well as any non-conference programs that may include participants: including talks, workshops, parties, on social media, and other online forums. GroupBy participants violating these rules may be sanctioned or expelled without a refund (if that applies) at the discretion of the conference organizers.

You can review the full policy at: **GroupBy.org/Code-of-Conduct**

# A Complete Introduction to Azure Data Factory

Paul Andrew | Principal Consultant & Solution Architect

Microsoft MVP Most Valuable Professional

altius

@MrPaulAndrew     In/MrPaulAndrew     MrPaulAndrew.com
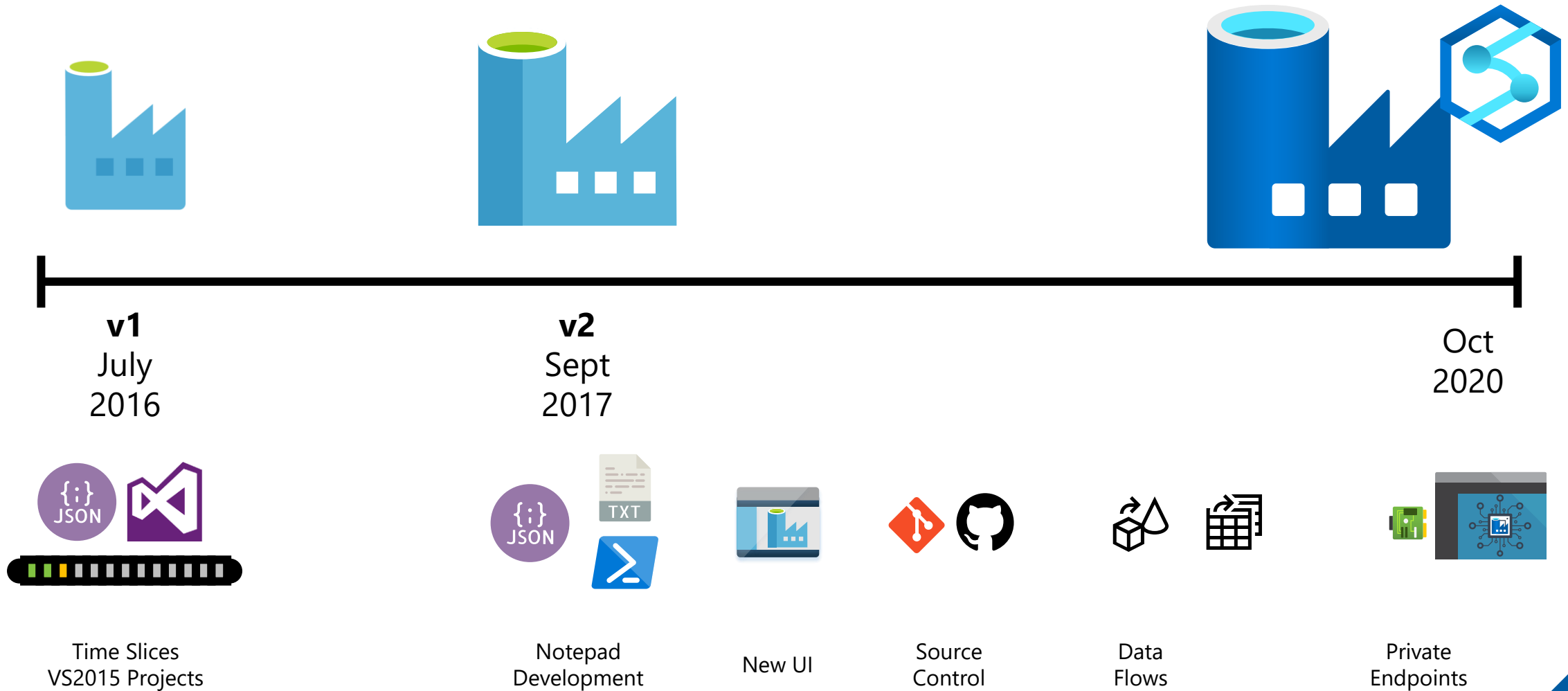
https://github.com/mrpaulandrew

CommunityEvents

# Agenda

What is it and why use it?

Data Factory Components

Common Activities

Execution Dependencies

Integration Runtimes

Running SSIS Packages in Azure

Data Factory Data Flows

Source Code & ARM Deployments

Monitoring & Logging

Conclusions

# Azure Data Factory –
What is it?
Why use it?

# A Quick History Lesson



**v1**
July
2016

**v2**
Sept
2017

Oct
2020

Time Slices
VS2015 Projects

Notepad
Development

New UI

Source
Control

Data
Flows

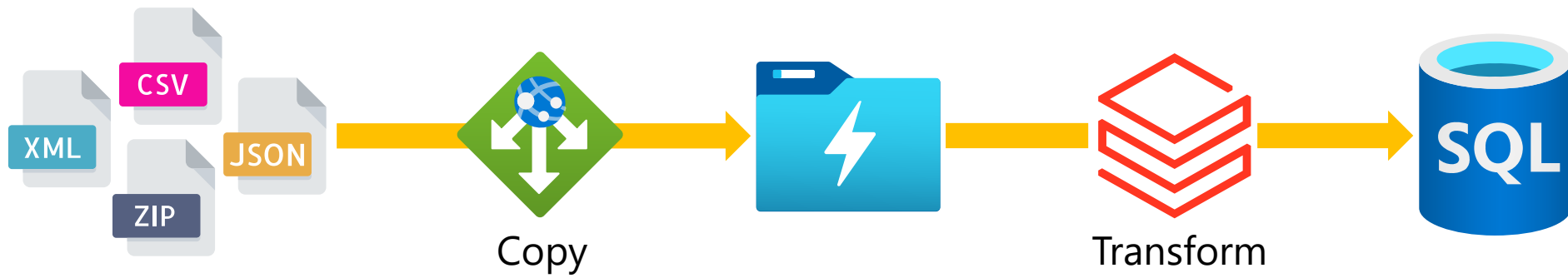Private
Endpoints

# What is Azure Data Factory (ADF)?

# What is Azure Data Factory (ADF)?



Copy

Transform

# Data Factory Components

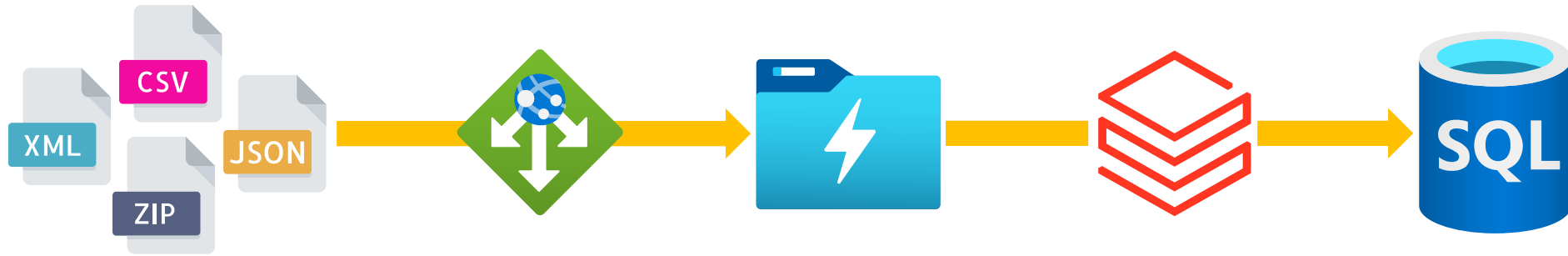# Data Factory Components



CSV
XML
JSON
ZIP
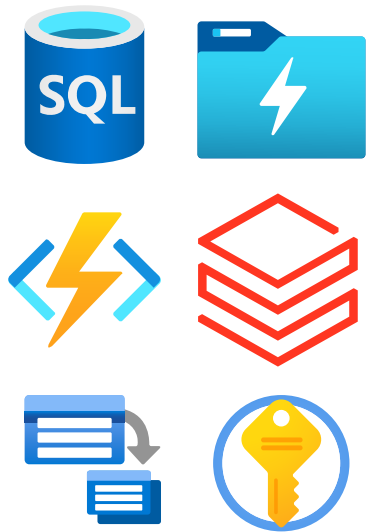→ Copy → Transform → SQL

# Data Factory Components
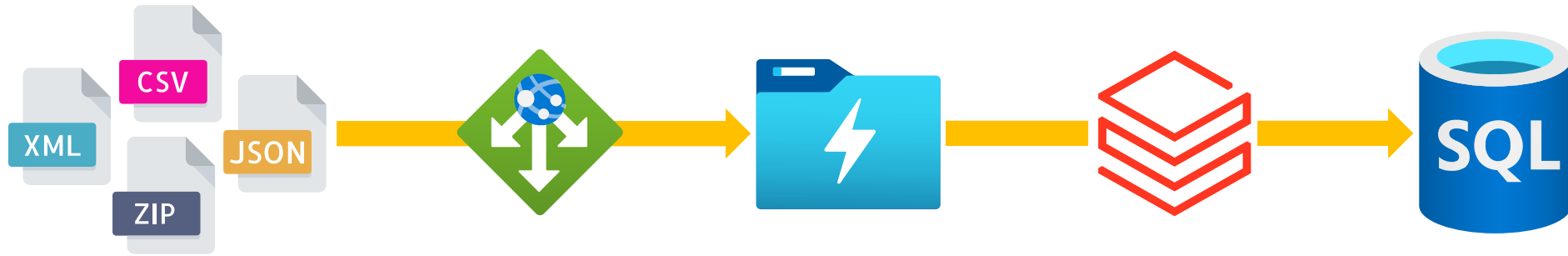
# Data Factory Components



**1** **Linked Services** – What to interact with and how?



SQLDBLinkedService

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

# Data Factory Components



**1** **Linked Services**

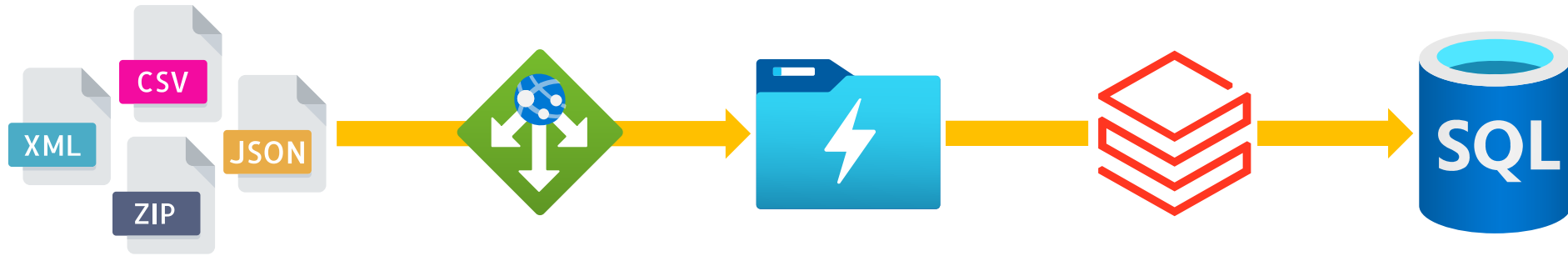**2** **Datasets** – Where is my data? What format? What file path/table do I need?

[dbo].[SalesOrders]

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

# Data Factory Components
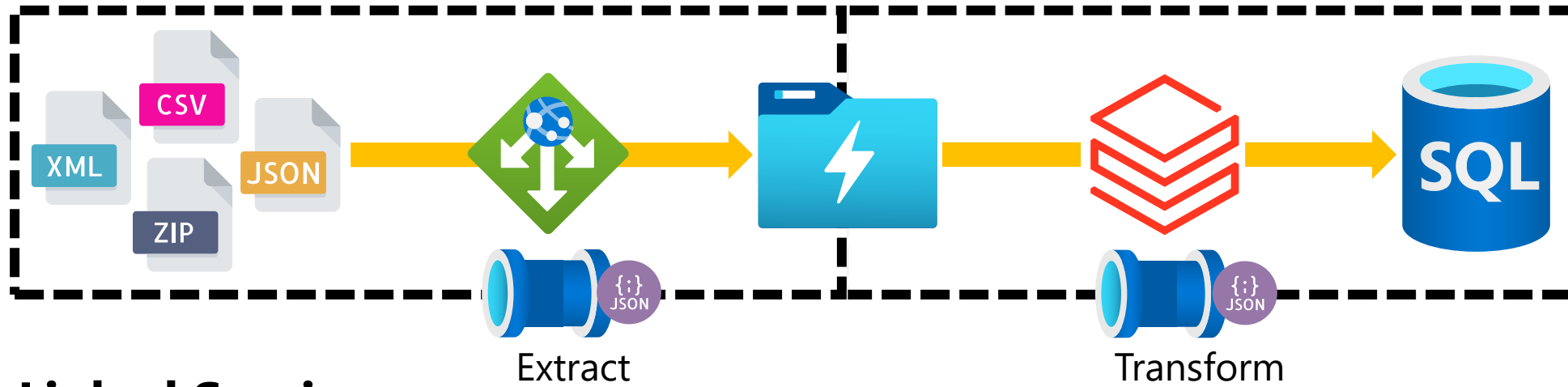


**1** **Linked Services**

**2** **Datasets**

**3** **Activities** – What do we want to happen when we invoke a Linked Service? With what conditions?

**Databricks Notebook Activity**

notebookPath: *ataPlayground/Playing*
baseParameters: *Testing*
libraries[jar]: dbfs:/lib1.jar
linkedServiceName: *BricksOfData01*

# Data Factory Components



Extract        Transform

**1**   **Linked Services**

**2**   **Datasets**

**3**   **Activities**

**4**   **Pipelines** – Logical groups of work that can be executed.

Sequence Container

Execute Package Task

Execute Pipeline Activity

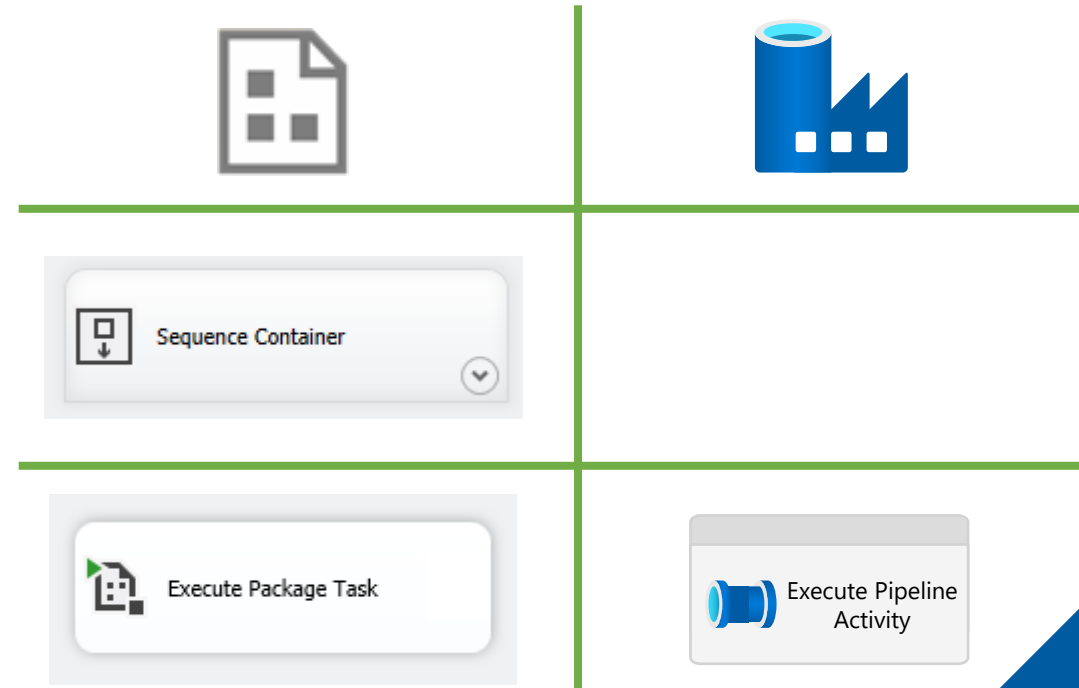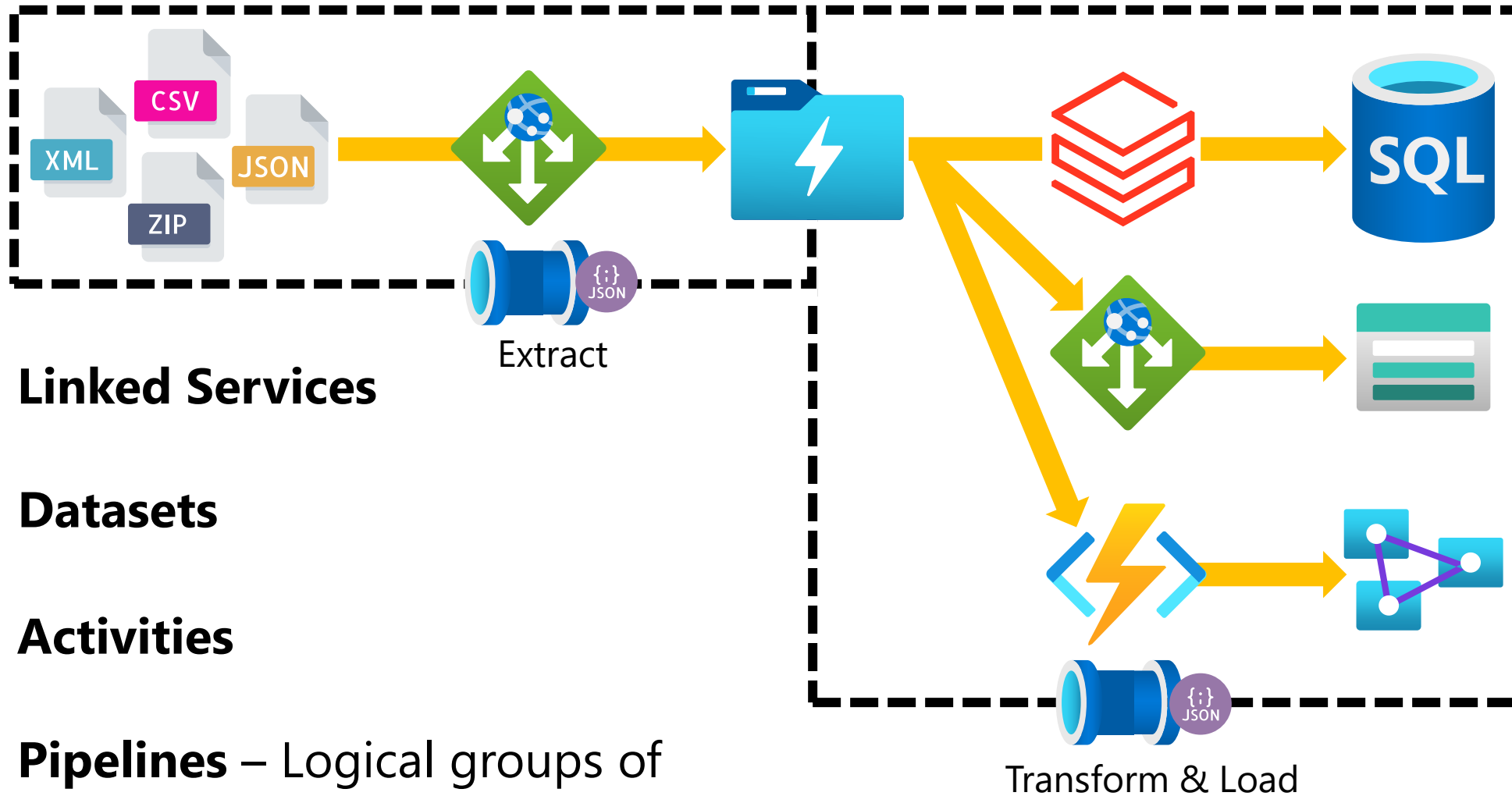# Data Factory Components



Extract

Transform & Load

1. **Linked Services**

2. **Datasets**

3. **Activities**

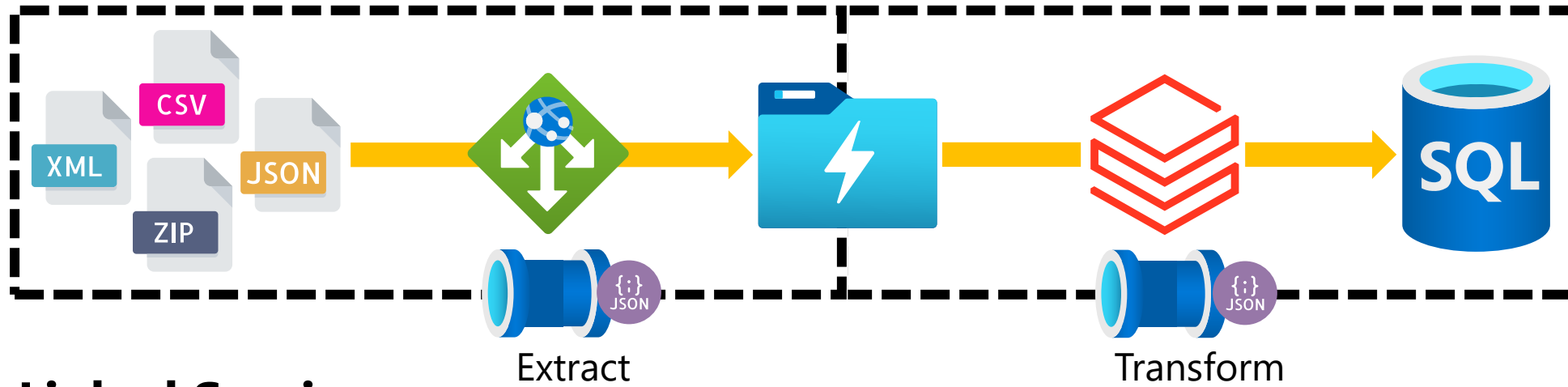4. **Pipelines** – Logical groups of work that can be executed.
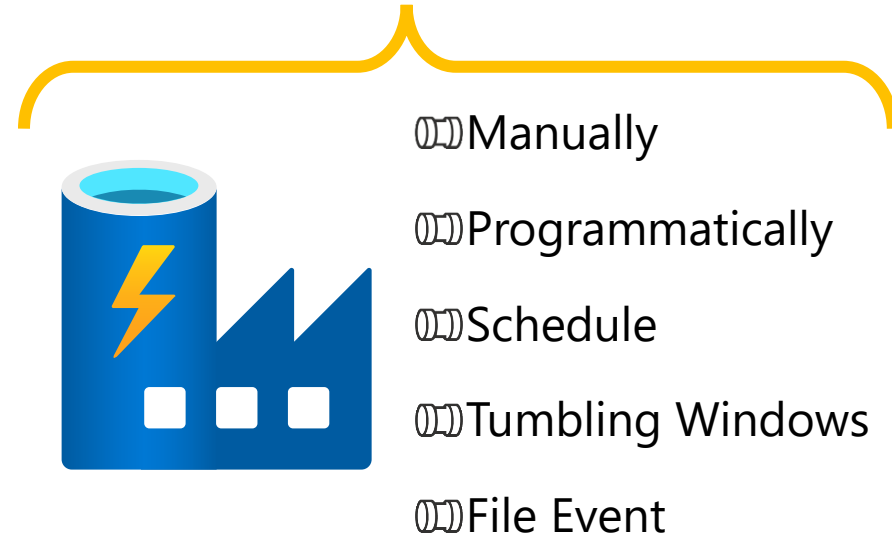
# Data Factory Components



Extract

Transform

1. **Linked Services**

2. **Datasets**

3. **Activities**

4. **Pipelines**

5. **Triggers** – Telling our when pipelines to run.

Manually

Programmatically
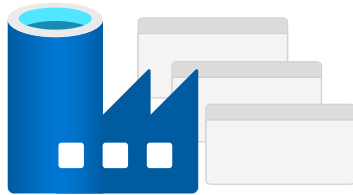
Schedule

Tumbling Windows

File Event

# Data Factory Components



1. **Linked Services**
2. **Datasets**
3. **Activities**
4. **Pipelines**
5. **Triggers**

# Data Factory Control Flow Components



1. **Linked Services**
2. **Datasets**
3. **Activities**
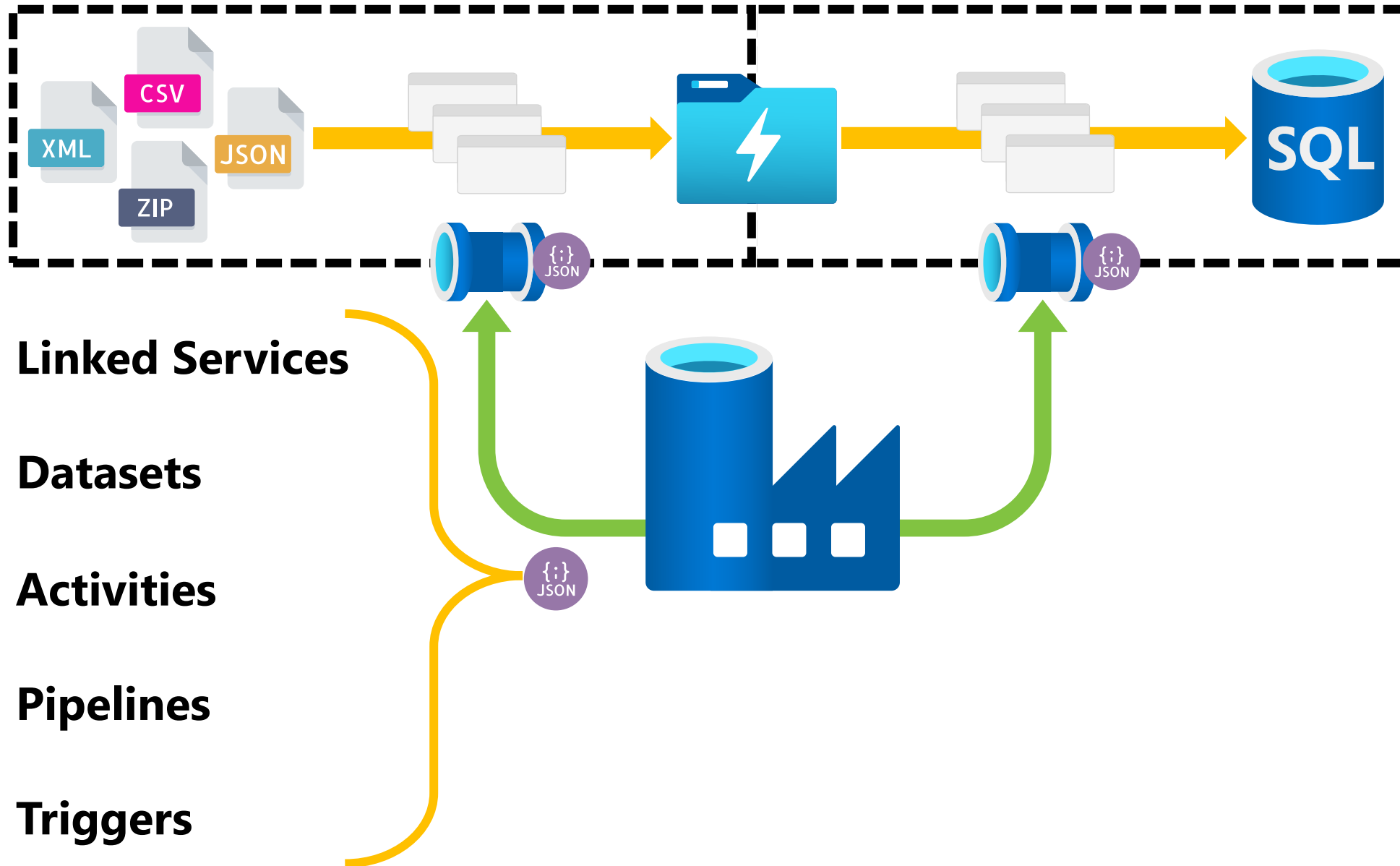4. **Pipelines**
5. **Triggers**
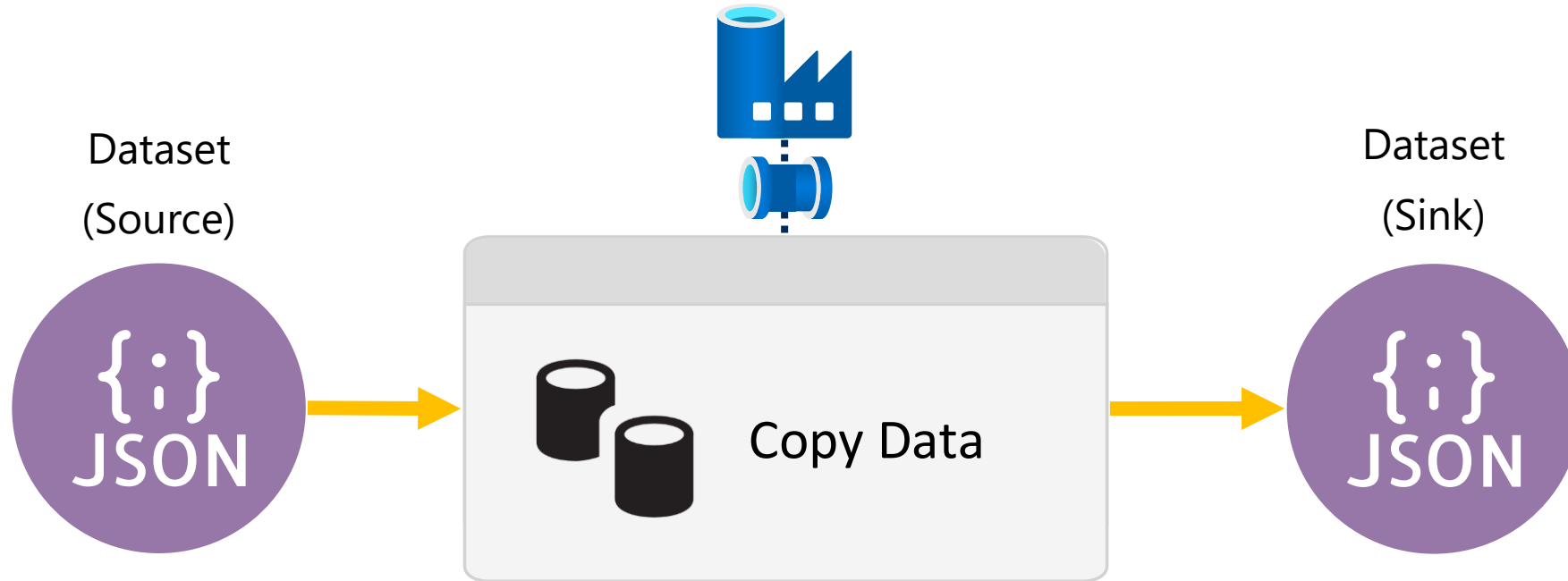
# Common Activities



```sql
SELECT TOP 5
    [ActivityName],
    [Inputs],
    [Outputs],
    [Details]
FROM
    [metadata].[AdfActivites]
WHERE
    [Notes] = 'Pauls Favourites';
```
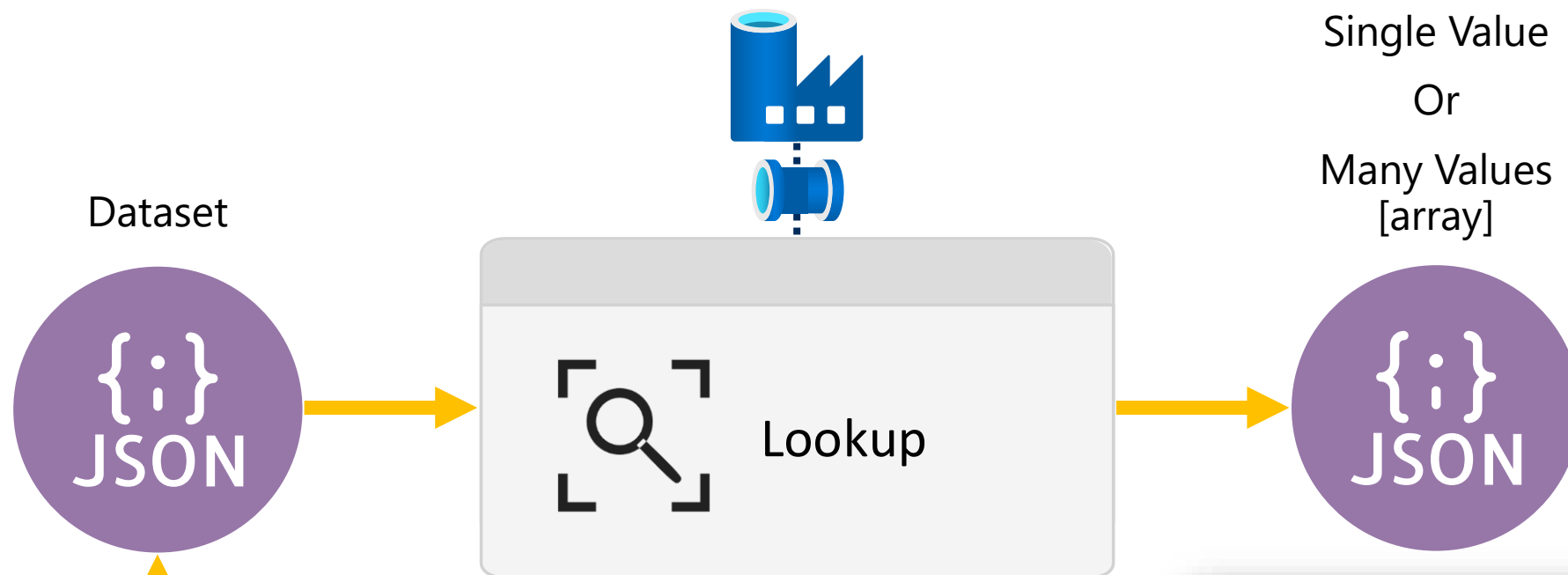
# Data Factory Common Activities



1 **Linked Services**

2 **Datasets**

3 **Activities**

4 **Pipelines**

5 **Triggers**

# Copy

Dataset
(Source)

Dataset
(Sink)

Copy Data

Auto Scaling

Transactional Restarts

Handle Zip Compression

Attribute Mapping and Schema Drift

Handle Failed Rows

Add Custom Attributes

Parse Excel & JSON Files

# Lookup
Get value to support other control flow activities

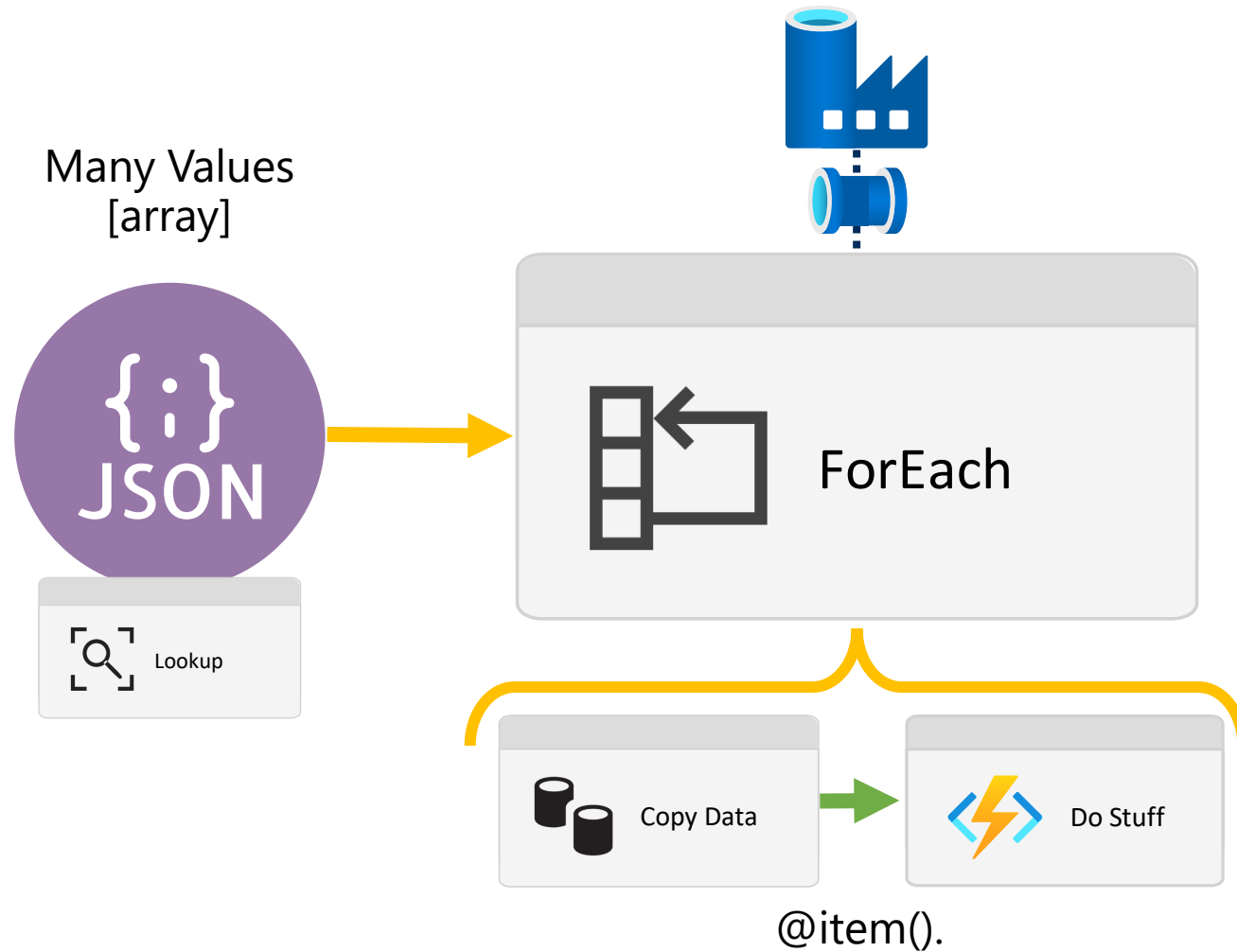Single Value

Or

Many Values
[array]

Dataset

**Lookup**

```
SELECT
    [SourceDIR],
    [TargetDIR],
    [FileName]
FROM
    [dbo].[FileList]
```

```
{
    "count": 3,
    "value": [
        {
            "SourceDIR": "ADFRoot\\ForUpload\\People\\",
            "TargetDIR": "RAW",
            "FileName": "Address.csv"
        },
        {
            "SourceDIR": "ADFRoot\\ForUpload\\People\\",
            "TargetDIR": "RAW",
            "FileName": "Gender.csv"
        },
        {
            "SourceDIR": "ADFRoot\\ForUpload\\People\\",
            "TargetDIR": "RAW",
            "FileName": "Ids.csv"
        }
    ]
}
```
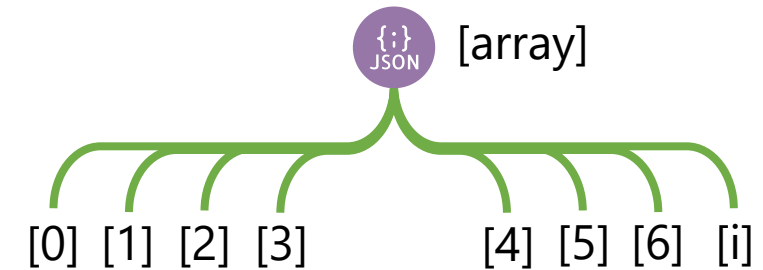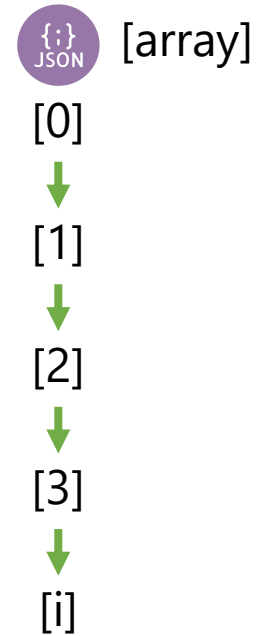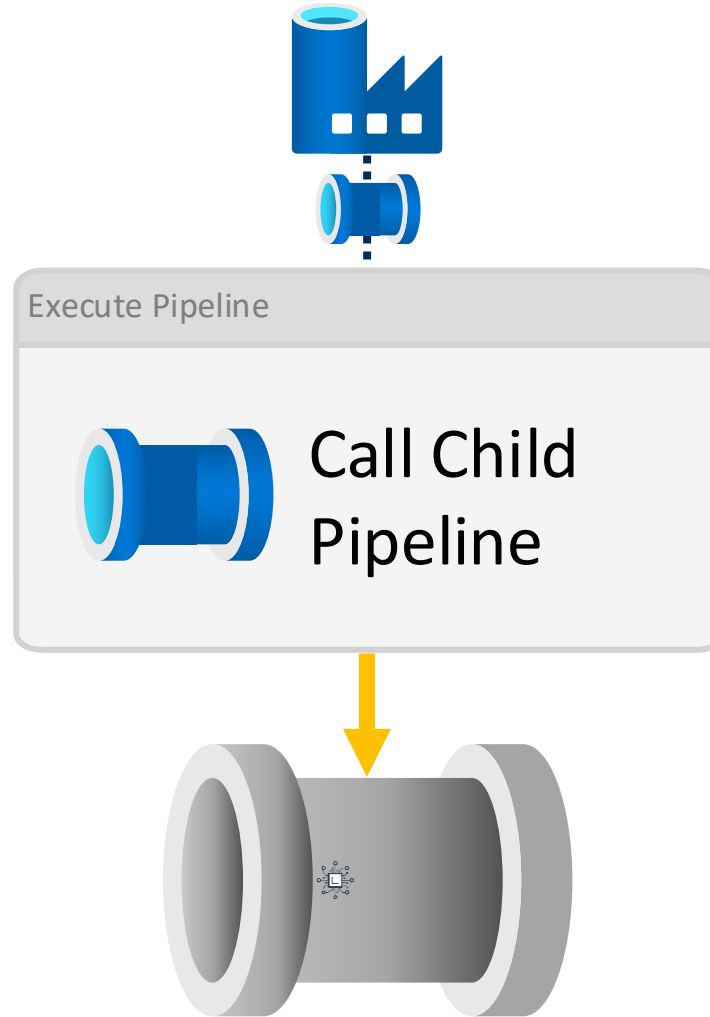
https://docs.microsoft.com/en-us/azure/data-factory/control-flow-lookup-activity

# ForEach
## Scaling Out Control Flow Activities

Many Values
[array]

JSON

Lookup

ForEach

Copy Data → Do Stuff

@item().

IsSequential:
true

JSON [array]

[0]

[1]

[2]

[3]

[i]

JSON [array]

[0] [1] [2] [3]     [4] [5] [6] [i]

Batch Count Default: 20
Batch Count Max: 50

https://docs.microsoft.com/en-us/azure/data-factory/control-flow-for-each-activity

# Execute Pipeline



Execute Pipeline

Call Child Pipeline

# Azure Function
## Extend Data Factory with Rest Calls

GET

POST

PUT

etc...

REST

{ ; }
JSON

Headers

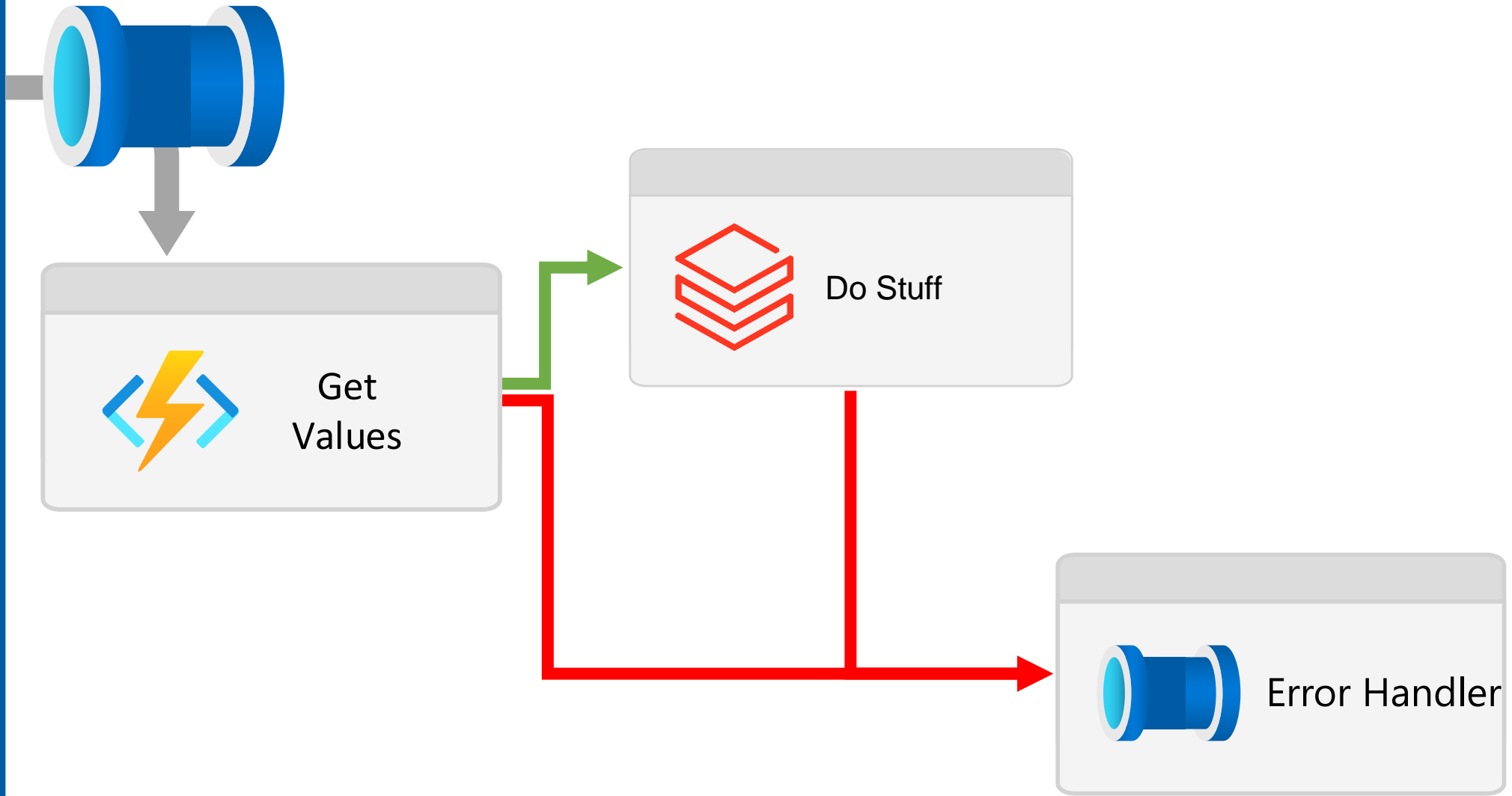Body

Do Stuff

C#

.NET
Core

???

# Execution Dependencies

# Execution Dependency Options

# Execution On Failure

# Execution On Failure or On Success

Get Values

Do Stuff

Error Handler

Execution On ???
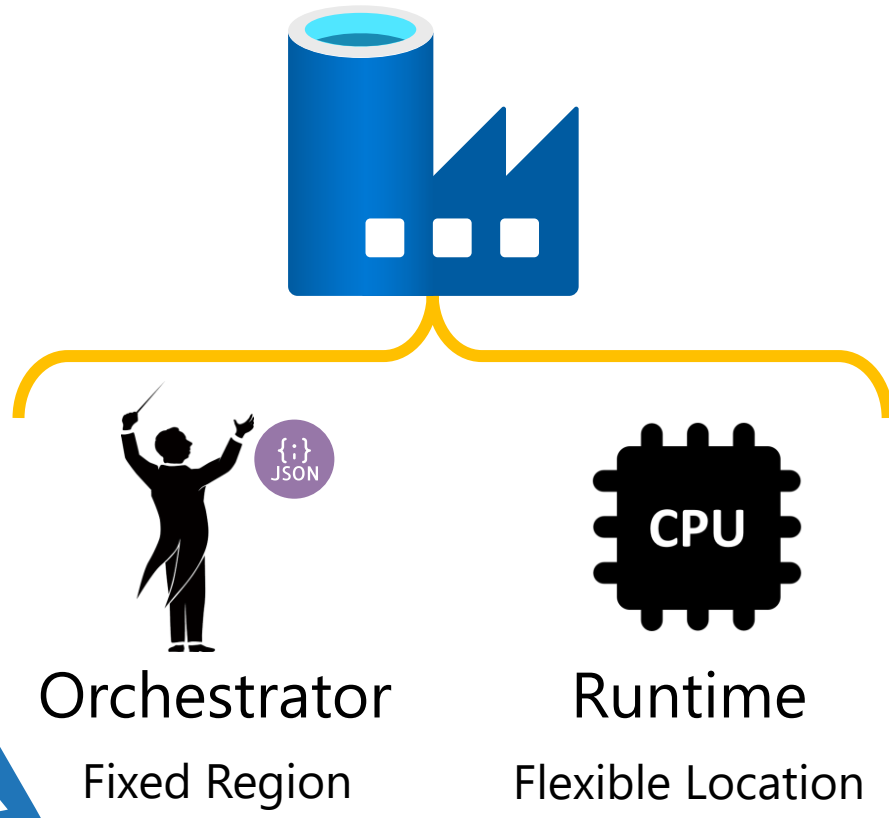
Get Values

Do Stuff

Run Stored Procedure

AND

AND

Error Handler

# Execution On Failure or On Success ✓

# Integration Runtimes

# What is an Integration Runtime?



Orchestrator
Fixed Region

Runtime
Flexible Location

Runtime 1

Runtime 2

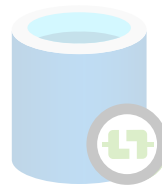Runtime 3

# What can an Integration Runtime do?

1. Azure IR

2. Hosted IR

3. SSIS IR

# Azure Integration Runtime



Azure IR

Hosted IR
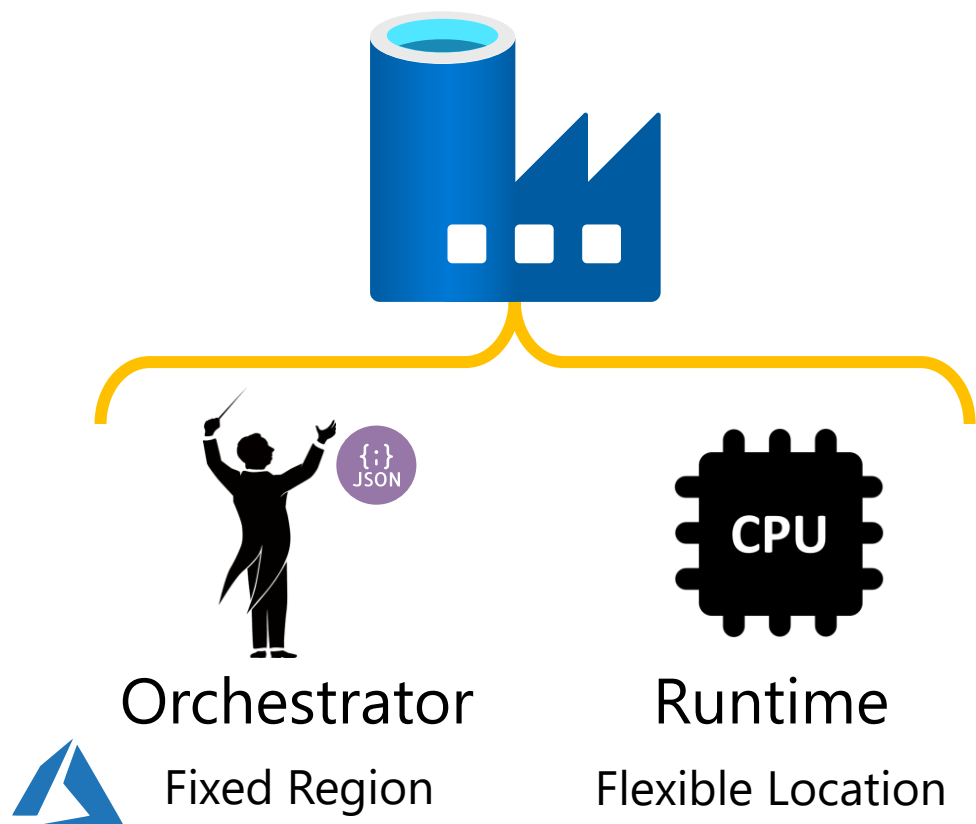
SSIS IR

# Azure Integration Runtime



Orchestrator — Fixed Region

Runtime — Flexible Location

# Azure Integration Runtime

Orchestrator — Fixed Region

Runtime — Flexible Region

AutoResolveIntegrationRuntime

Runtime 1

Runtime 2

Runtime 3

Internal ✓ | External ✗

# Hosted Integration Runtime

Azure IR

Hosted IR

SSIS IR

# Hosted Integration Runtime

Runtime

Orchestrator

Fixed Region

Runtime

Flexible Region

# Hosted Integration Runtime



Runtime
Locally Hosted

Orchestrator
Fixed Region

Runtime
Flexible Region

# Hosted Integration Runtime – Secondary Nodes



Runtime
Locally Hosted

Orchestrator
Fixed Region

Runtime
Flexible Region

# Hosted Integration Runtime – Linked



Runtime
Locally Hosted

Orchestrator
Fixed Region

Runtime
Flexible Region

# SSIS Integration Runtime



Azure IR
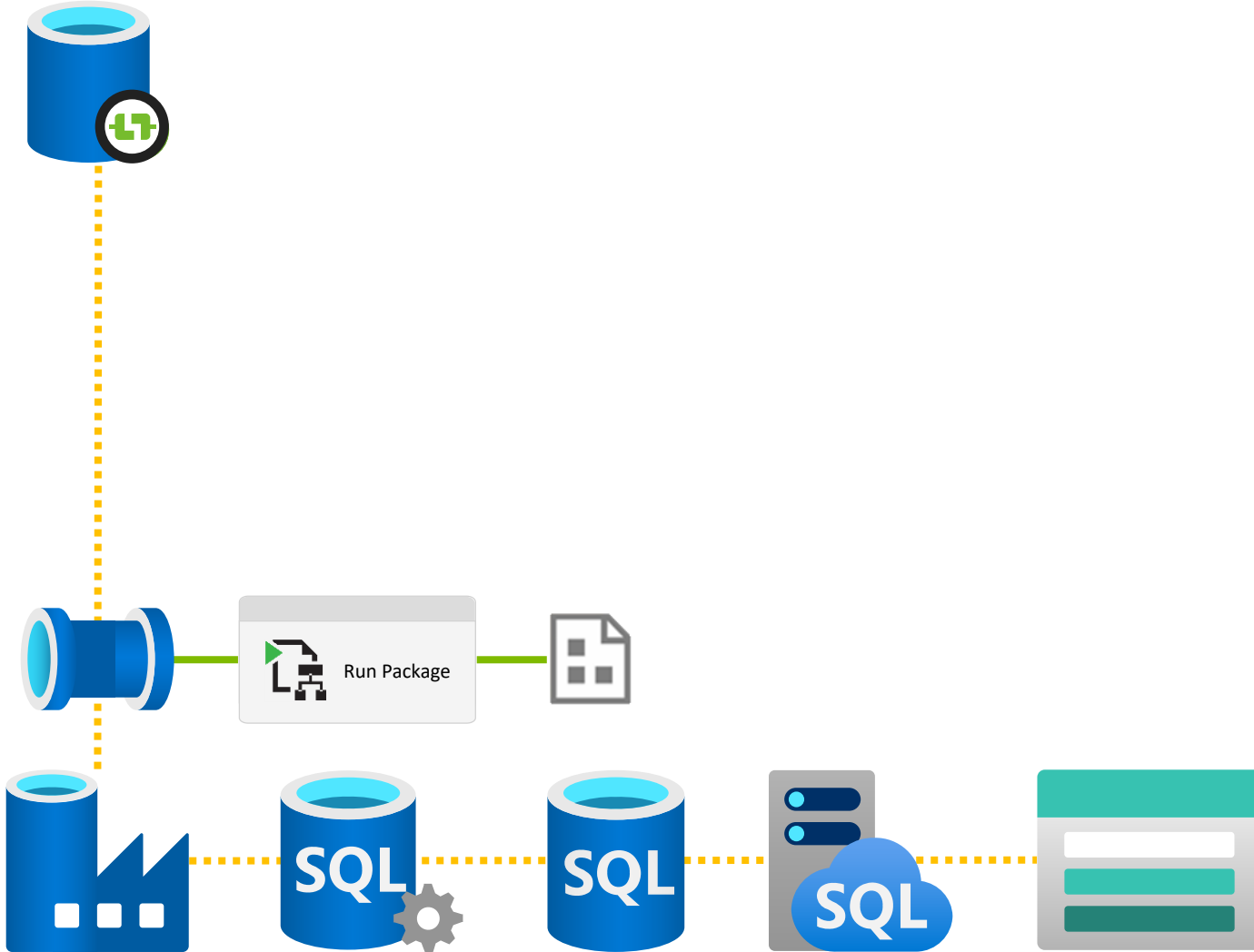
Hosted IR

SSIS IR

# Running an SSIS Package in Azure

SSIS IR

# Running an SSIS Package in Azure

SSIS IR



Run Package

# Problem: Using All Of The SSIS IR Compute

SSIS IR

Supports 80 Concurrent Packages

MAXDOP = 80

Runs 1 Package

Parent Package

Child Packages x80

Pipeline x1

Activities x80

Pipeline x1

ForEach Max Batch (50)

Run Package

ForEach

Run Package

# Data Factory
# Data Flows

# Data Factory Control Flow Components



1. **Linked Services**
2. **Datasets**
3. **Activities**
4. **Pipelines**
5. **Triggers**

# Data Factory Data Flow Components
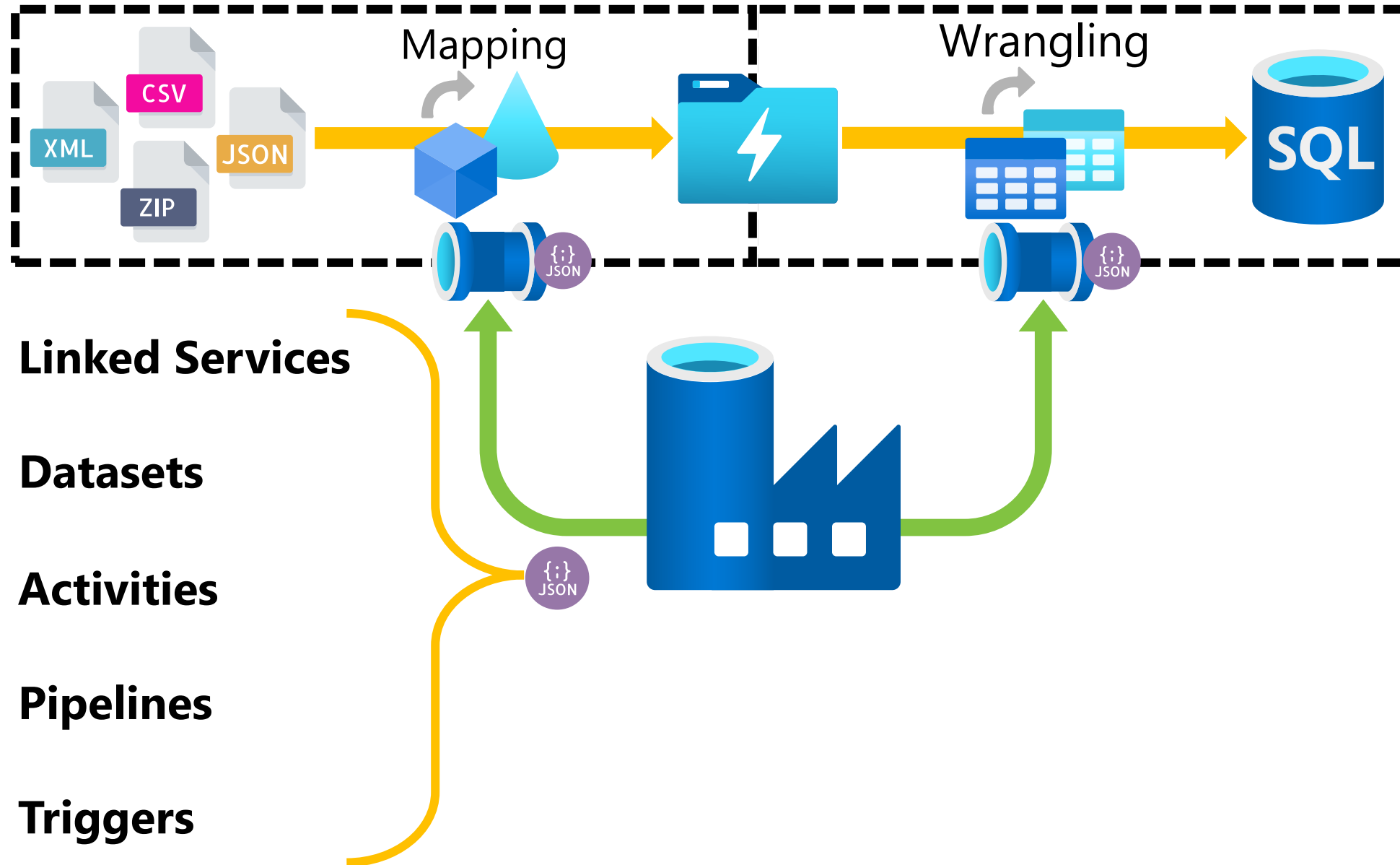


Mapping
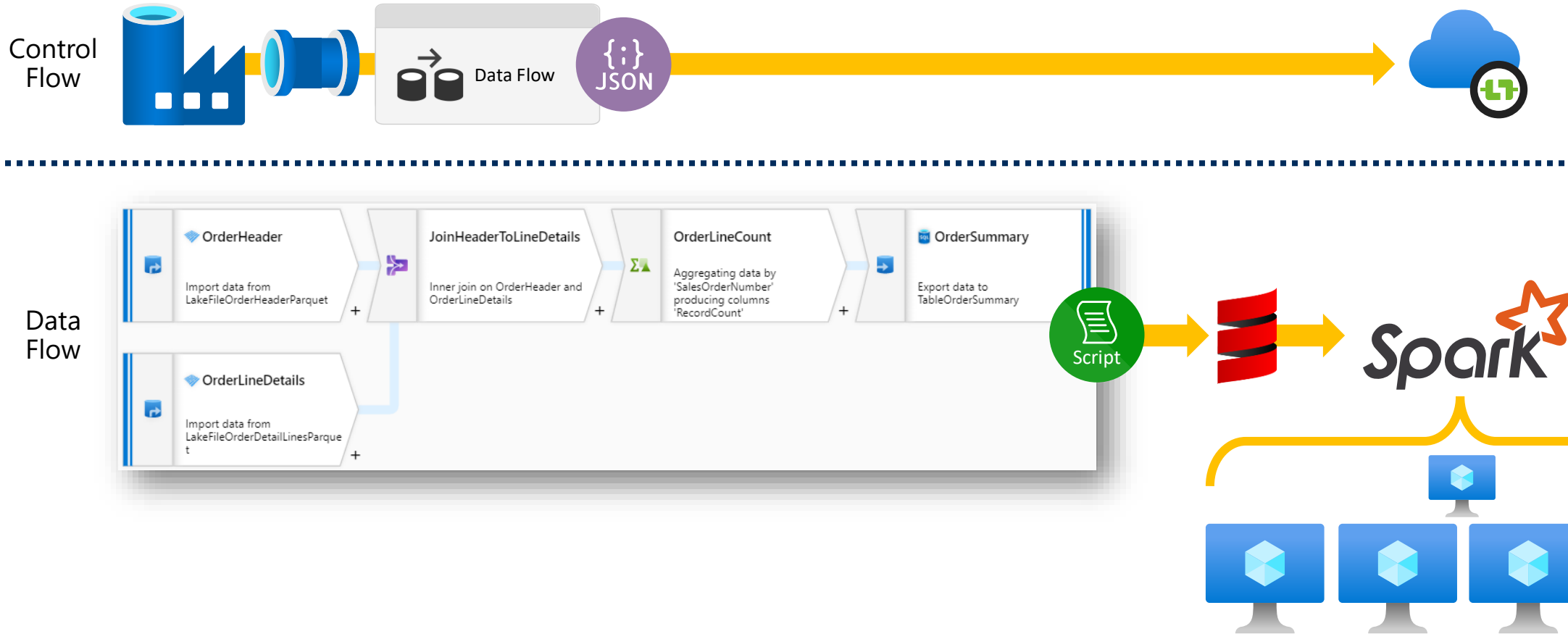
Wrangling

CSV
XML
JSON
ZIP

SQL

1 Linked Services

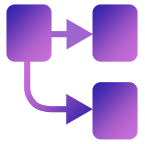2 Datasets

3 Activities

4 Pipelines

5 Triggers

# What is a Mapping Data Flow?



**A:** Graphic data transformation tool that sits on top of Apache Spark.

# What can a Mapping Data Flow do? - Transformations
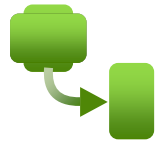
New Branch

Join

Conditional Split

Exists

Union

Lookup

Derived Column

Select

Aggregate

Surrogate Key

Pivot
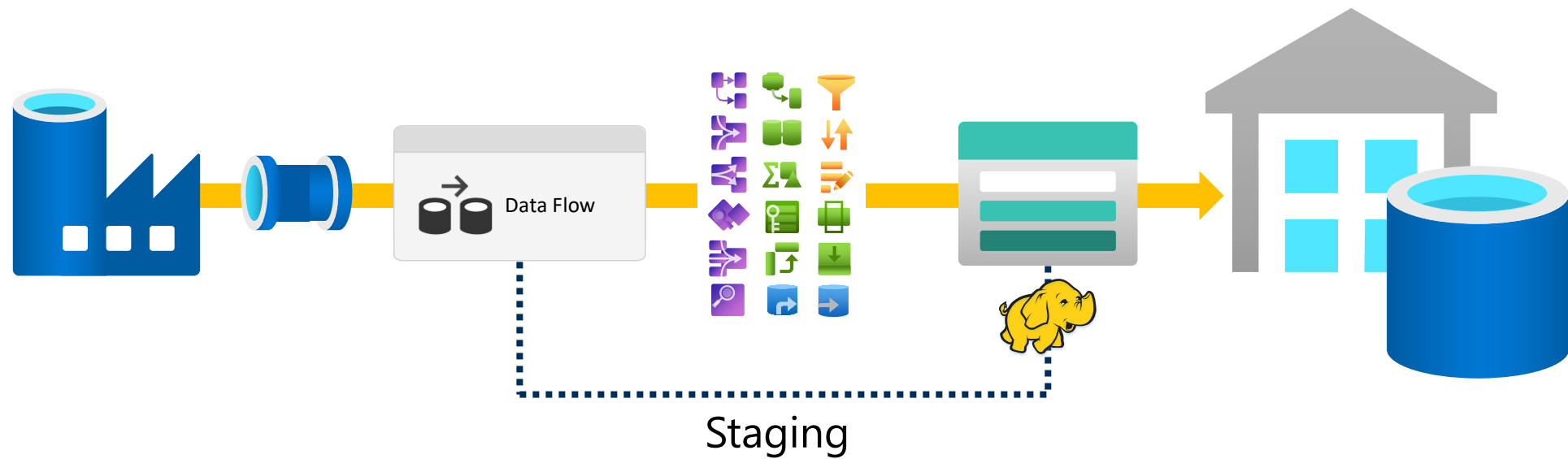
Unpivot

Window

Flatten

Filter

Sort

Alter Row

Key

**Input & Output Modifiers**

**Schema Modifiers**

**Row Modifiers**

# What can a Mapping Data Flow do? - PolyBase
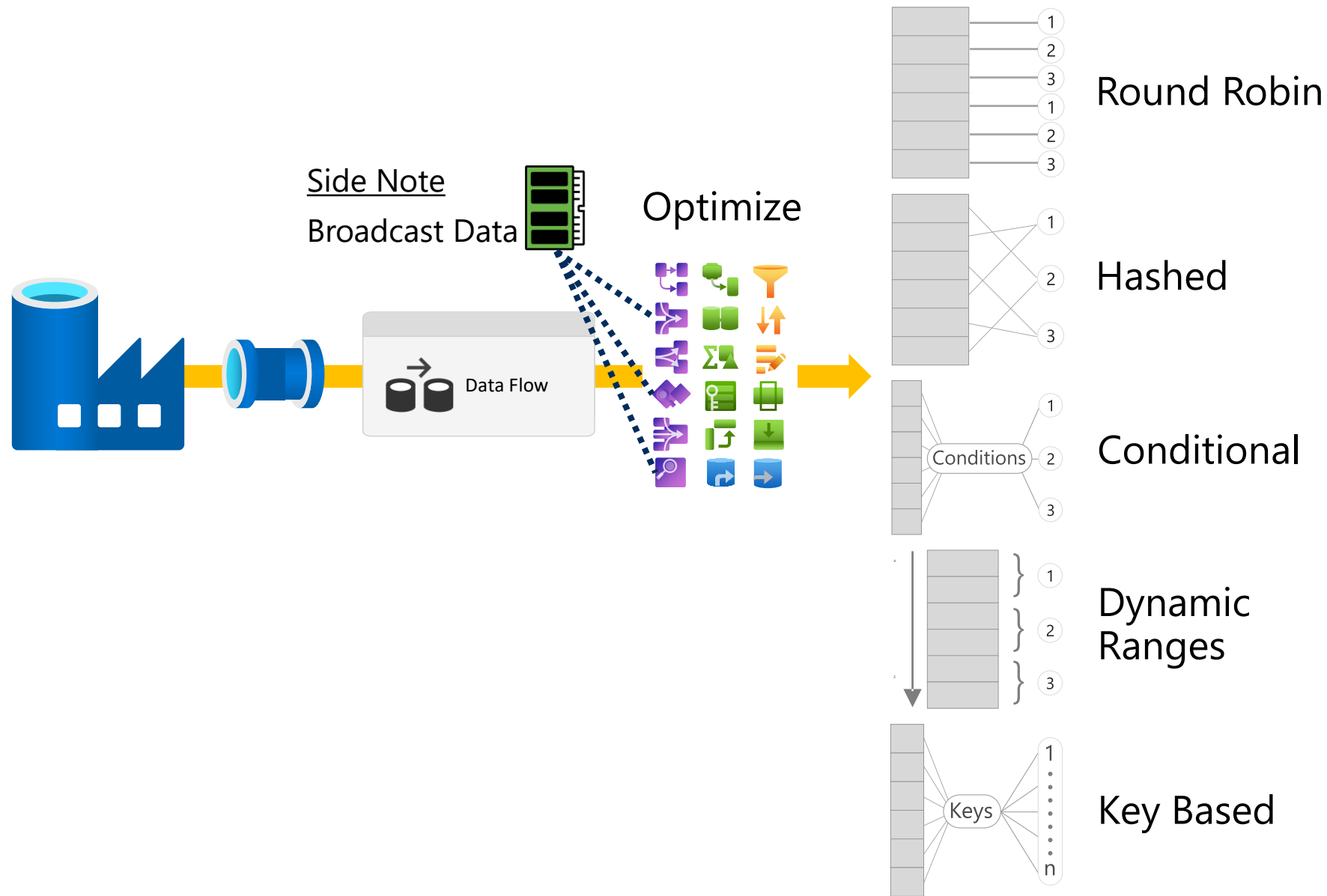
Data Flow

Staging

PolyBase ⓘ

Staging linked service     Select...    ⓘ   + New

Staging storage folder     Container   /   Directory     📁 Browse | ⌄

https://docs.microsoft.com/en-gb/azure/data-factory/data-flow-expression-functions

# What can a Mapping Data Flow do? - Partition Handling

Round Robin

Hashed

Conditional

Dynamic Ranges

Key Based

Side Note
Broadcast Data

Optimize

Data Flow

Conditions

Keys

# What is a Wrangling Data Flow?

# What can a Wrangling Data Flow do? - Home

# Data Flow Cluster Configuration



Default Azure IR

- General Purpose
- 4x Worker Nodes
- 0 Minutes

Control Flow

Data Flow

JSON

Data Flow

Mapping

Wrangling

Script

Spark

- Compute Type
- Number of Worker Nodes
- Cluster Time to Live

# Other Data Transformation Services in Azure



| SSIS Packages | HD Insight | Data Lake Analytics | Synapse SQLDW | SQL Database | Batch Service | Durable Functions | Synapse Spark | Databricks Spark | Analysis Services | Cosmos DB |

# When Should We Use Data Flows?



SSIS Packages · HD Insight · Data Lake Analytics · Synapse SQLDW · SQL Database · Batch Service · Durable Functions · Synapse Spark · Databricks Spark · Analysis Services · Cosmos DB

Mapping

Wrangling

Data Flows

# Data Transformations in Azure Comparisons

| Transformation Method | Graphical UI | Scales Out | Scales Up | Cloud Native Tech |
|---|:---:|:---:|:---:|:---:|
| T-SQL (SQLDB) | ✖ | ✖ | ✔ | ✖ |
| SSIS | ✔ | ✖ | ✔ | ✖ |
| Scala (Databricks) | ✖ | ✔ | ✔ | ✔ |
| Data Factory Data Flows | ✔ | ✔ | ✔ | ✔ |

# 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 power users who has grown out of Power BI following a series of Data Lake exploration sessions.

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

Simpler and quicker data engineering for data scientists that want to quickly prepare raw data for model training and testing, also with the ability to use large amounts of compute.

# Source Code &
ARM Deployments

# Getting Our ADF Source Code



Developer

Git

Debug Service

Template Parameters

ARM Template Export

branch

save

merge

master

publish

Debug

Prod

Deployed Components

1 Linked Services
2 Datasets
3 Activities
4 Pipelines
5 Triggers

adf_publish

ARMTemplateForFactory.json

Template Parameters

# 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

Gitflow

# 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 **Datasets**

3 **Activities**

4 **Pipelines**

5 **Triggers**

# Monitoring & Logging

# Diagnostic Settings



Log Analytics

Storage

Event Hub

# Diagnostic Settings

Log Analytics

# Using Log Analytics

```
ADFPipelineRunDurations
 | project
        TimeGenerated,
        Start,
        End,
        ['DataFactory'] = substring(ResourceId, 121, 100),
        Status,
        PipelineName,
        Parameters,
        ["RunDuration"] = datetime_diff('Minute', End, Start)
 | where

        TimeGenerated > ago(1h)
        and Status !in ('InProgress','Queued','Cancelling')
```

**Pipeline Durations**
procfwkloganalytics

RunDuration

20

15

10

5

1:10 PM  1:20 PM  1:30 PM  1:40 PM  1:50 PM  2:00 PM

TimeGenerated

— FRAMEWORKFACTORY    — FRAMEWORKFACTORYDEV    — WORKERSFACTORY

ADF.procfwk
SQL

**Resources and Content**   Edit

| | Blogs | mrpaulandrew.com/ADF.procfwk |
| | GitHub | github.com/mrpaulandrew/ADF.procfwk |
| | Twitter | #ADFprocfwk |

**FrameworkSupportF...**
Function App

Running

**Function Call Durations**
ProcFwkAppInsights

duration
50,000
25,000
0
1:10 PM  1:20 PM  1:30 PM  1:40 PM  1:50 PM  2:00 PM
timestamp

● CheckPipelineStatus   ◆ ExecutePipeline   ■ SendEmail

**ProcFwkAppInsights**
Application Insights

**procfwkloganalytics**
Workspace

**Pipeline Durations**
procfwkloganalytics

20
15
10
5
0
RunDuration

1:10 PM  1:20 PM  1:30 PM  1:40 PM  1:50 PM  2:00 PM
TimeGenerated

─●─ FRAMEWORKFACTORY   ─◆─ FRAMEWORKFACTORYDEV   ─■─ WORKERSFACTORY

**FrameworkFactor**
Data factory

**FrameworkFactor**
Data factory

**FrameworkFactor**
Data factory

**WorkersFactory**
Data factory

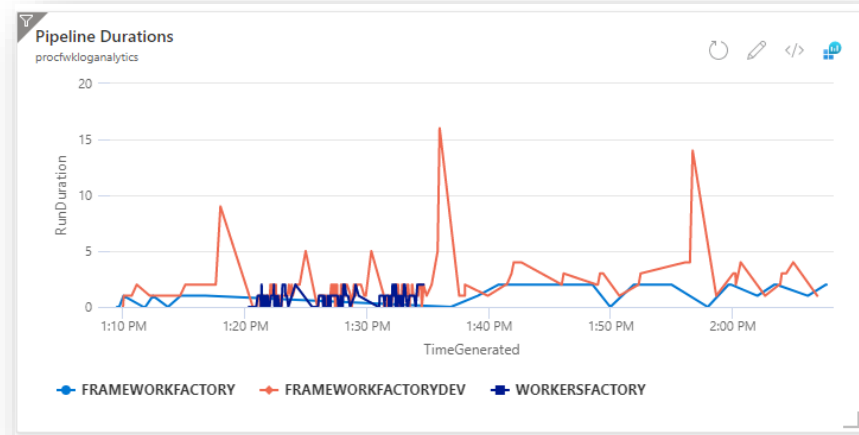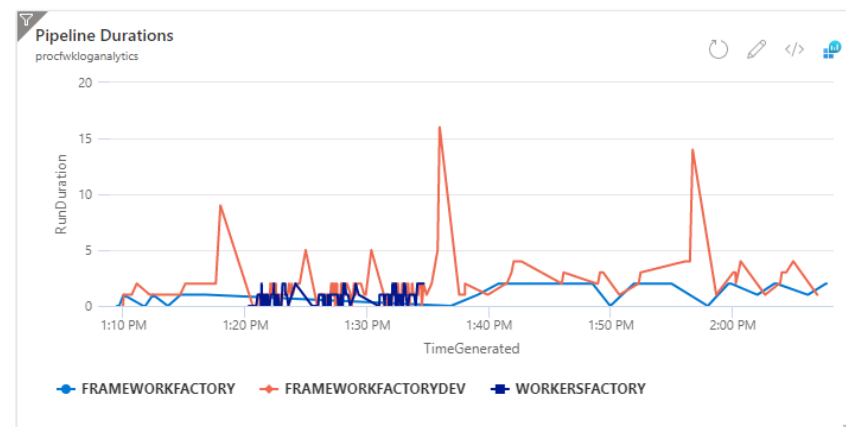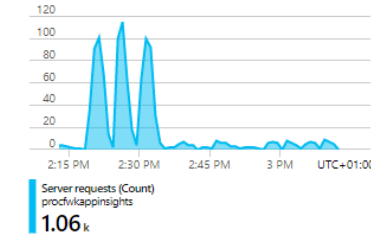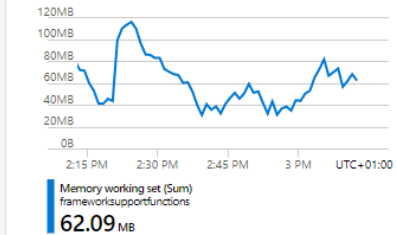**Server requests**

120
100
80
60
40
20
0
2:15 PM  2:30 PM  2:45 PM  3 PM  UTC+01:00

Server requests (Count)
procfwkappinsights
**1.06** k

**Memory working set**

120MB
100MB
80MB
60MB
40MB
20MB
0B
2:15 PM  2:30 PM  2:45 PM  3 PM  UTC+01:00

Memory working set (Sum)
frameworksupportfunctions
**62.09** MB

**Resources**
ADF.procfwk

- ProcFwkLogAnalytics
- FrameworkFactory
- FrameworkFactoryDev
- FrameworkKeys
- platformsupport01
- FrameworkMetadataDev (pl...
- frameworksupportstore
- frameworkstorage01
- FrameworkSupportFunctions
- FrameworkFactoryTest
- WorkersFactory
- frameworkonsynapse
- UKSouthPlan
- FrameworkMetadataTest (pl...
- ProcFwkAppInsights
- 9a4fe00e-39d9-4ec8-8f88-5...
- frameworkdatalake01
- sqlvaexht4i7t63enw

**FrameworkMetadat...**
SQL database
Online

**Compute utilization**

10%
0%
2:15 PM  2:30 PM  2:45 PM  3 PM  UTC+01:00

DTU percentage (Max)
platformsupport01/frameworkmetadatadev
**13** %

**FrameworkKeys**
Key vault

**Average latency**

500ms
400ms
300ms
200ms
100ms
0ms
6 PM    Sep 7    6 AM    12 PM    UTC+01:00

auxiliary          secret
frameworkkeys      frameworkkeys
**45.45** ms       **33.57** ms

# Using Data Explorer



Data Explorer
Operational Data Warehouse

Write KQL queries in Azure
Data Studio Notebooks

KQL

**Small/Medium**

**Large/Enterprise**

# Conclusions

# What is Azure Data Factory (ADF)?

# What is Azure Data Factory?



1. A complete Microsoft Azure integration tool.
2. Orchestrator of our Control Flow operations – with scale out Activities.
3. Orchestrator of our Data Flow 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**

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

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

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

# Thank you for listening...

Paul Andrew

**Microsoft MVP**
Most Valuable Professional

**altius**

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

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

/CommunityEvents

**GitHub:**    github.com/mrpaulandrew

/ContentCollateral