

STRATEGIC PARTNER



GOLD SPONSOR



Future Processing

SILVER SPONSOR





BRONZE SPONSOR





Azure Orchestration – Applying Data Factory in Production

Paul Andrew | Principal Consultant & Solution Architect



Microsoft®
Most Valuable
Professional

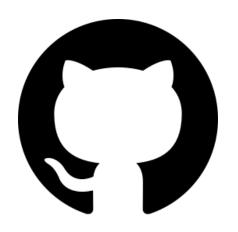












https://github.com/mrpaulandrew

CommunityEvents

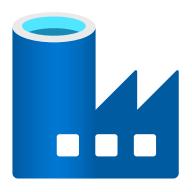
Demo code, content and slides from various community events.

C++

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

AGENDA ??





AGENDA - Short Stories

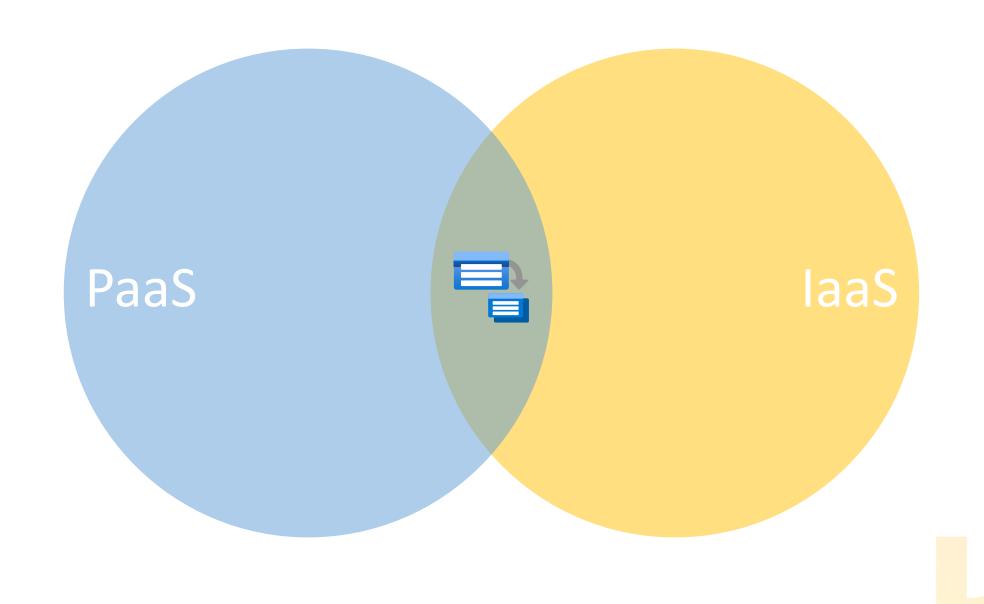
- **DD** Custom Activities
- (III) Controlling & Scaling Compute
- Scale Out Execution
- Metadata Driven Pipelines
- DD Deploying Data Factory
- DD Best Practices

Custom Activities

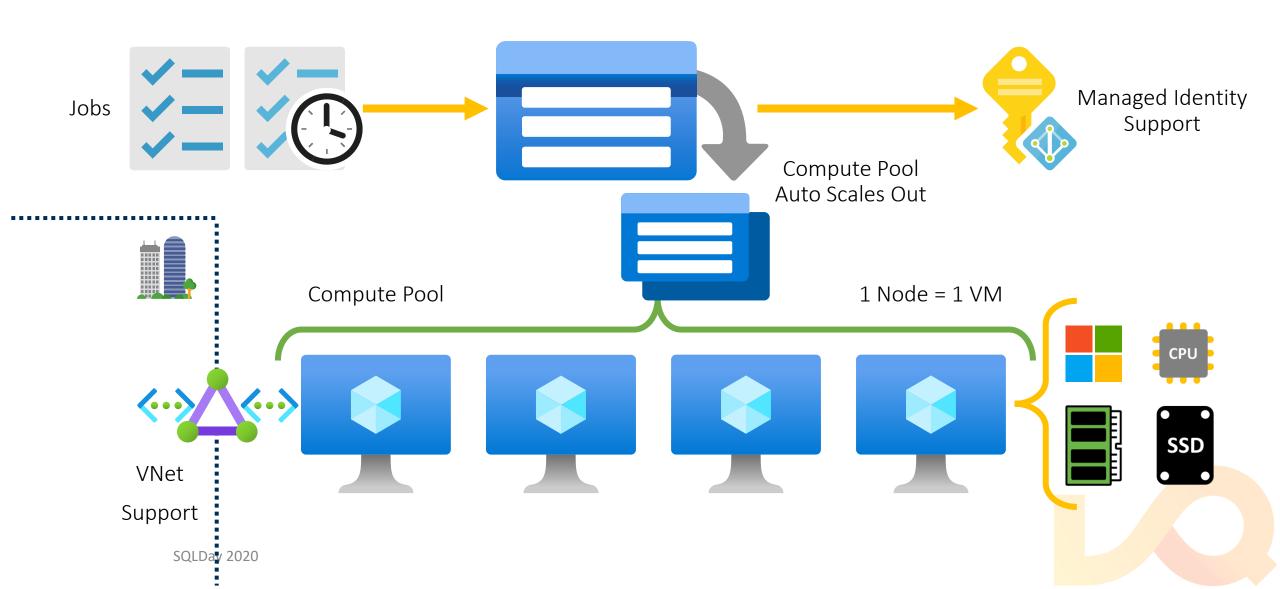




Azure Batch Service



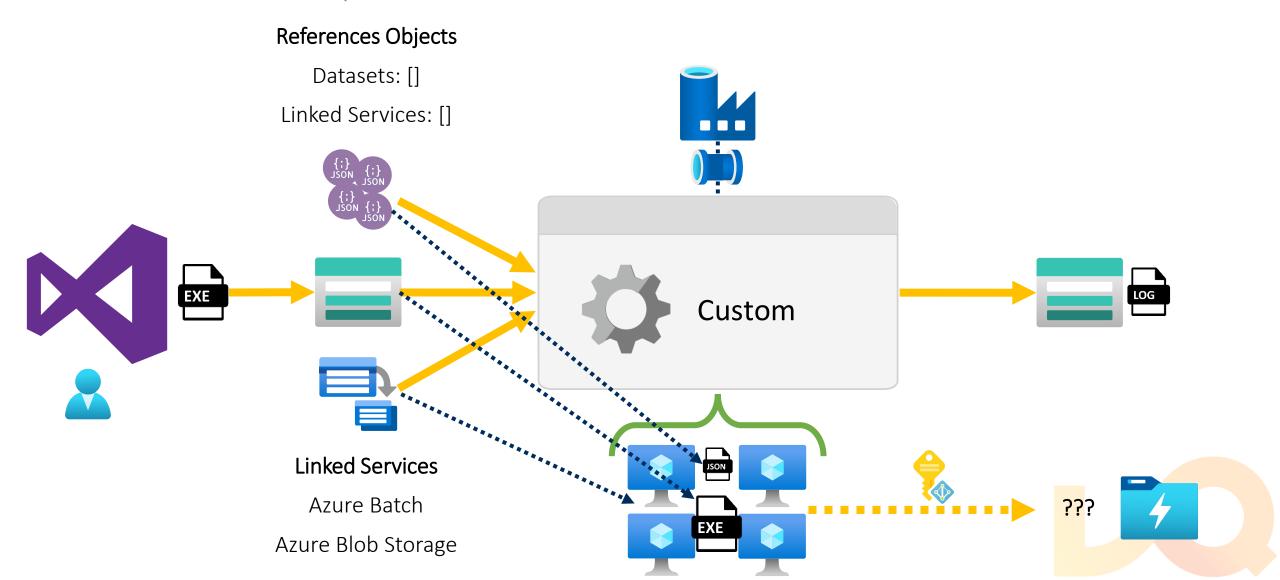
Azure Batch Service



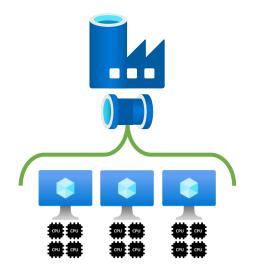
Custom Activity



Extend Data Factory with Custom Code



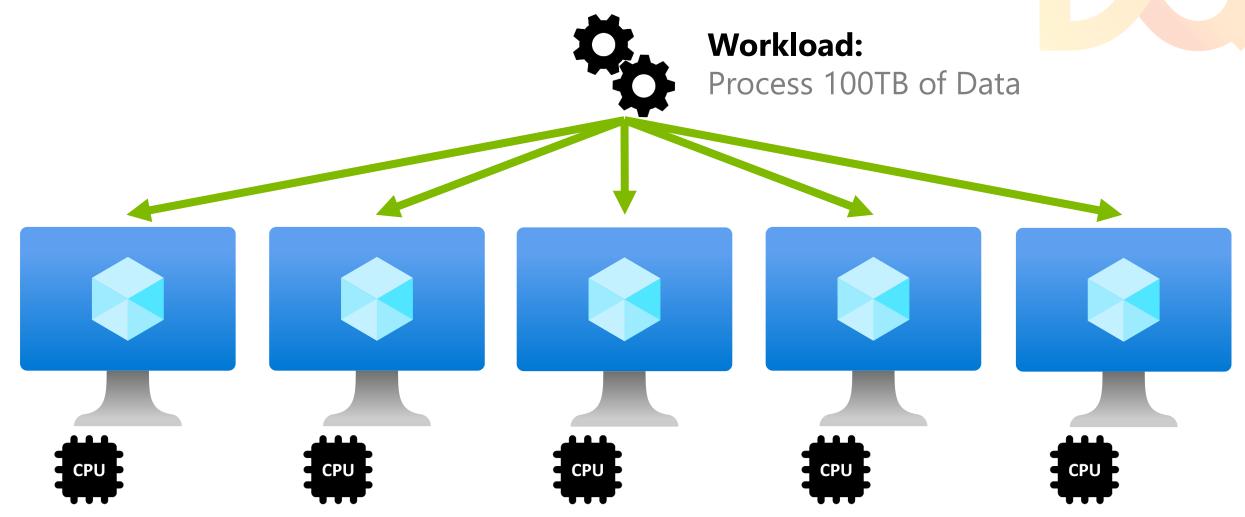
Controlling & Scaling Compute



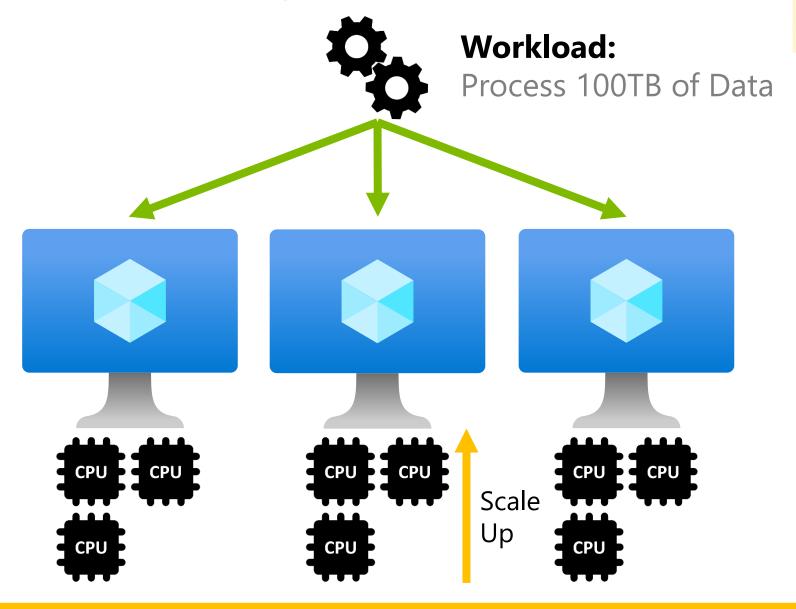


Scaling Up and/or Scaling Out Workload: Process 100TB of Data < CPU CPU -CPU CPU **CPU CPU CPU CPU CPU CPU** Scale Up **CPU** CPU **CPU CPU CPU CPU**

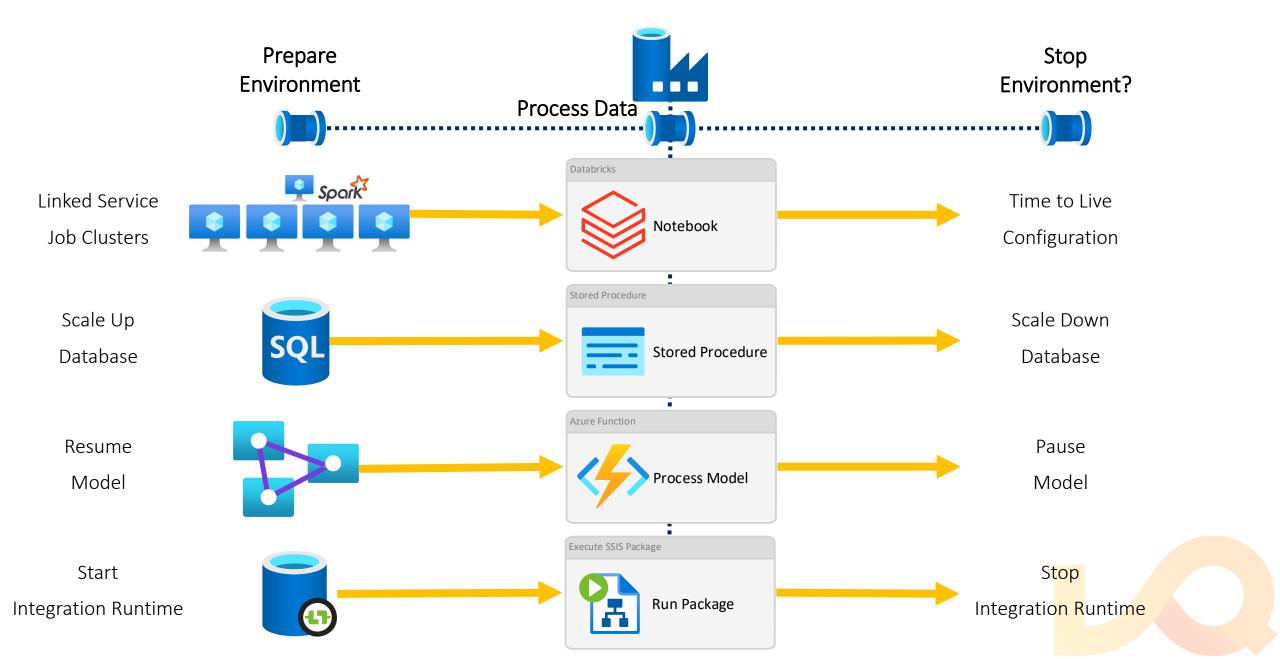
Scaling Up and/or Scaling Out



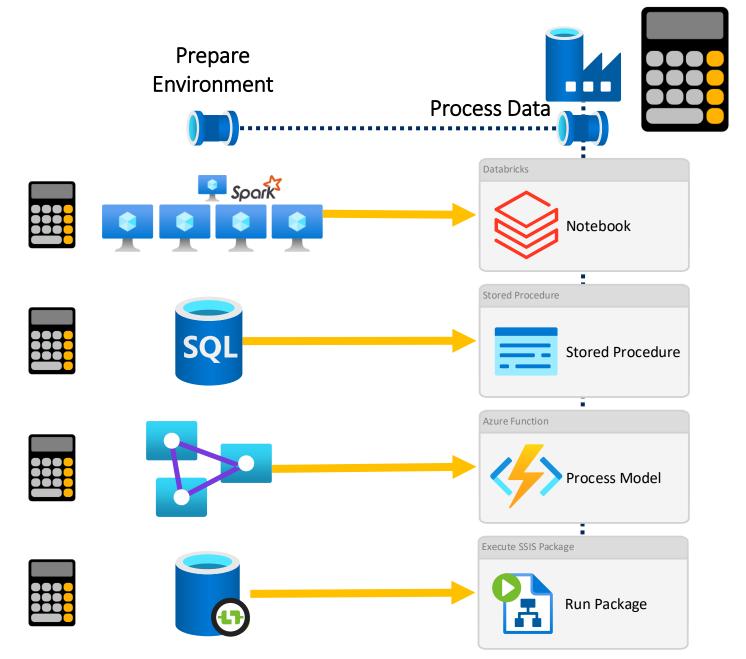
Scaling Up and/or Scaling Out



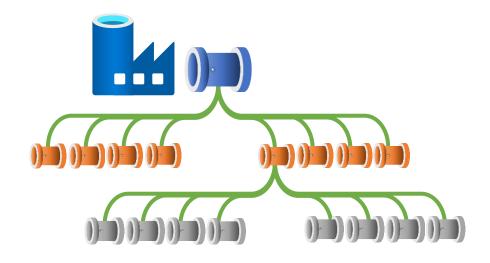
Resource Control



Resource Control - Cost



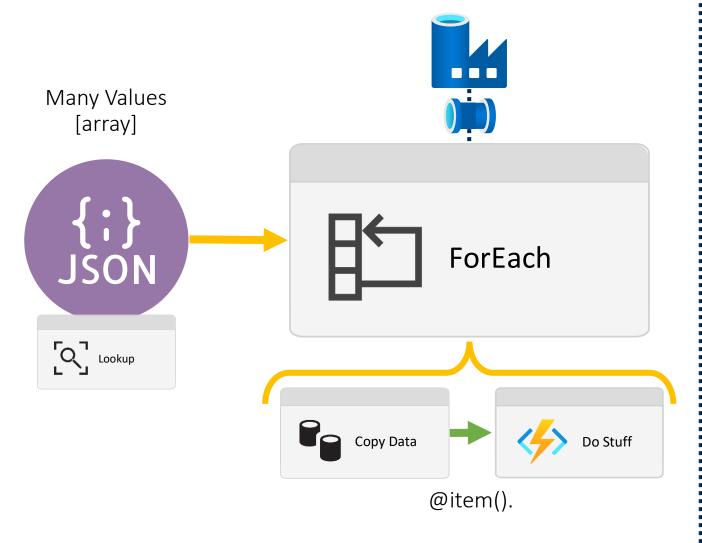
Scale Out Execution

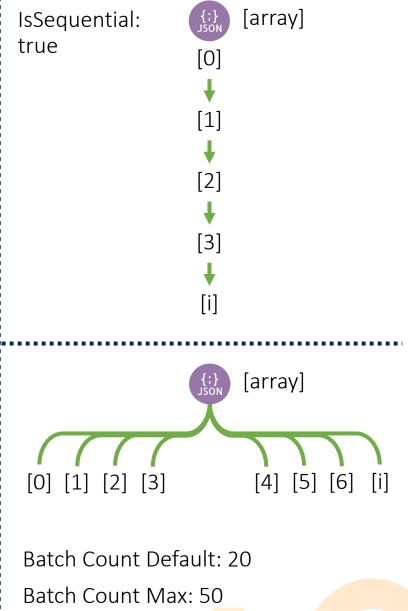




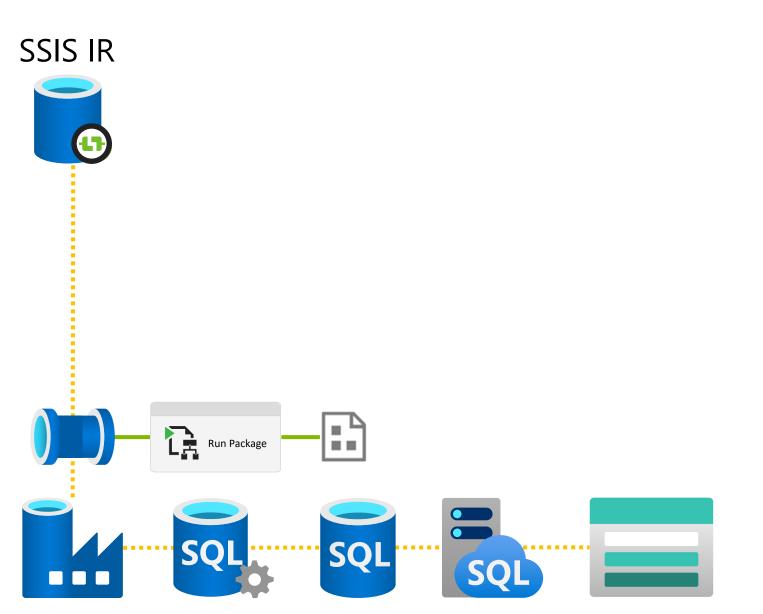
ForEach

Scaling Out Control Flow Activities

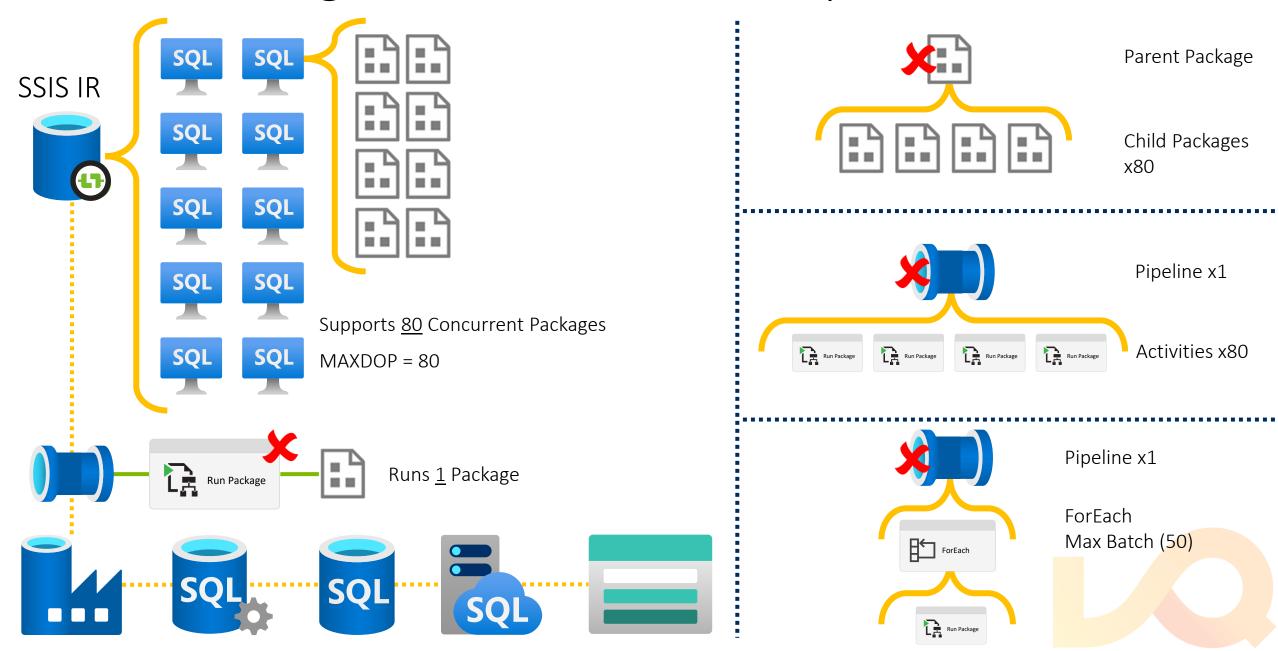




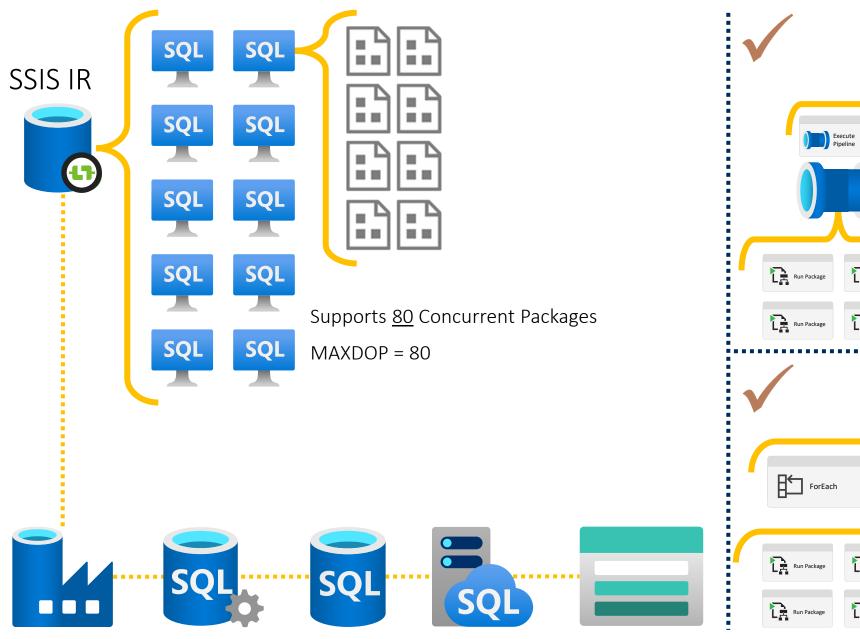
Problem: Using All Of The SSIS IR Compute

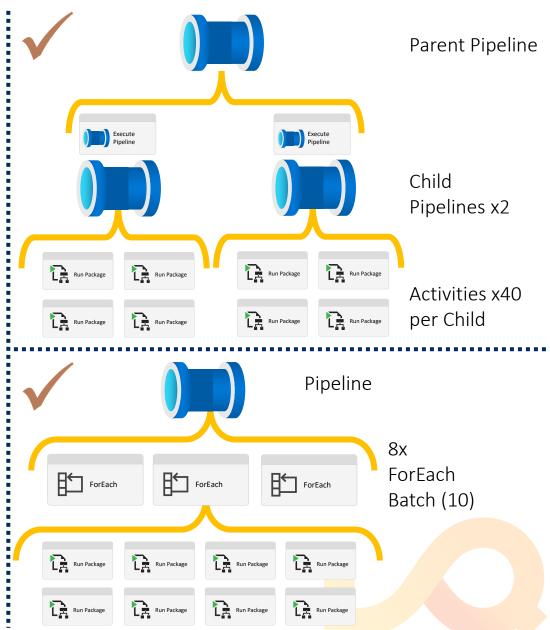


Problem: Using All Of The SSIS IR Compute

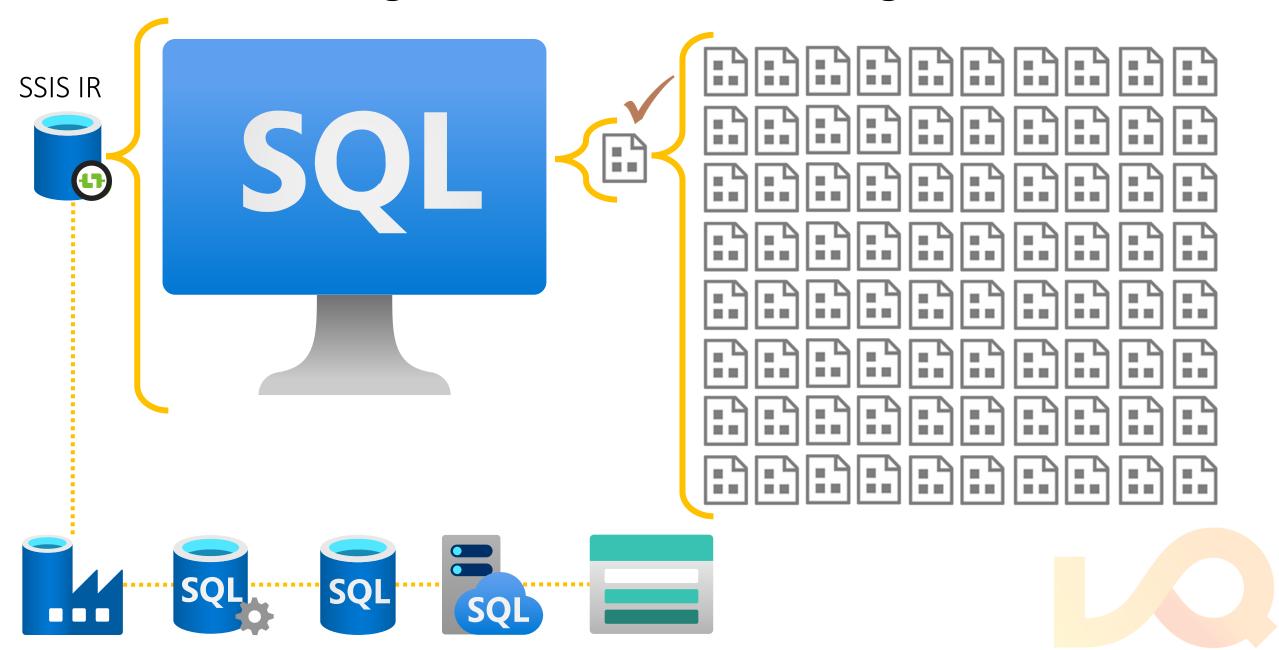


Solution 1 & 2: Static Pipelines

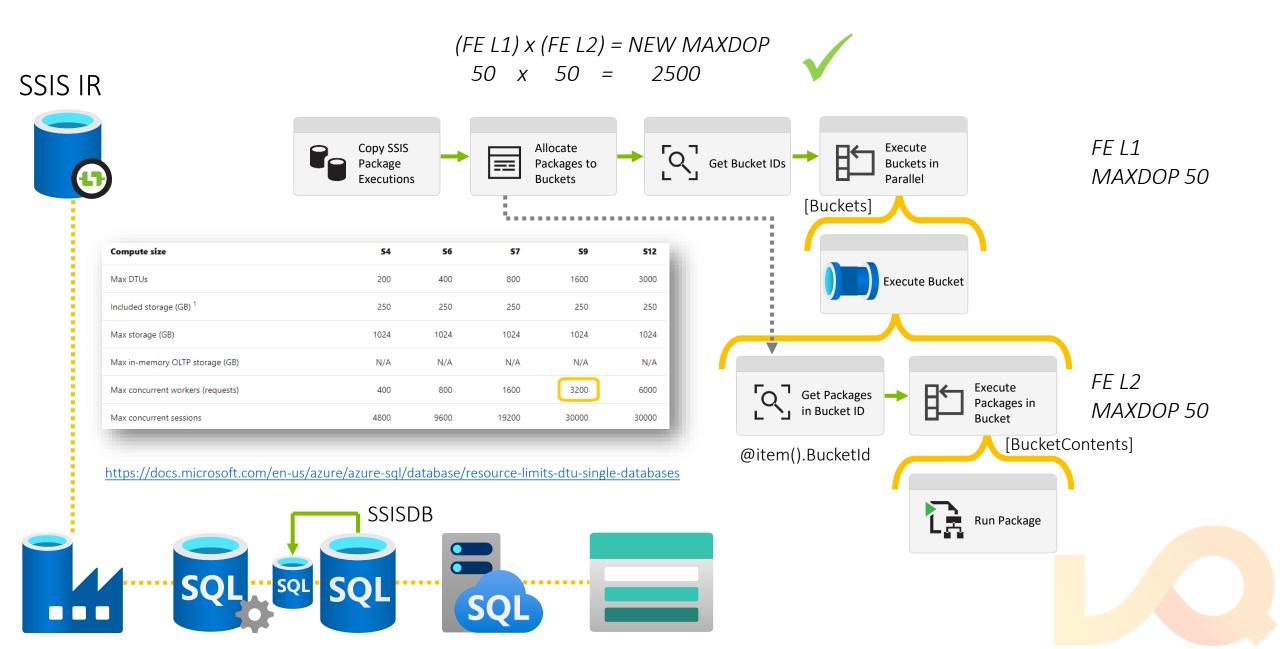




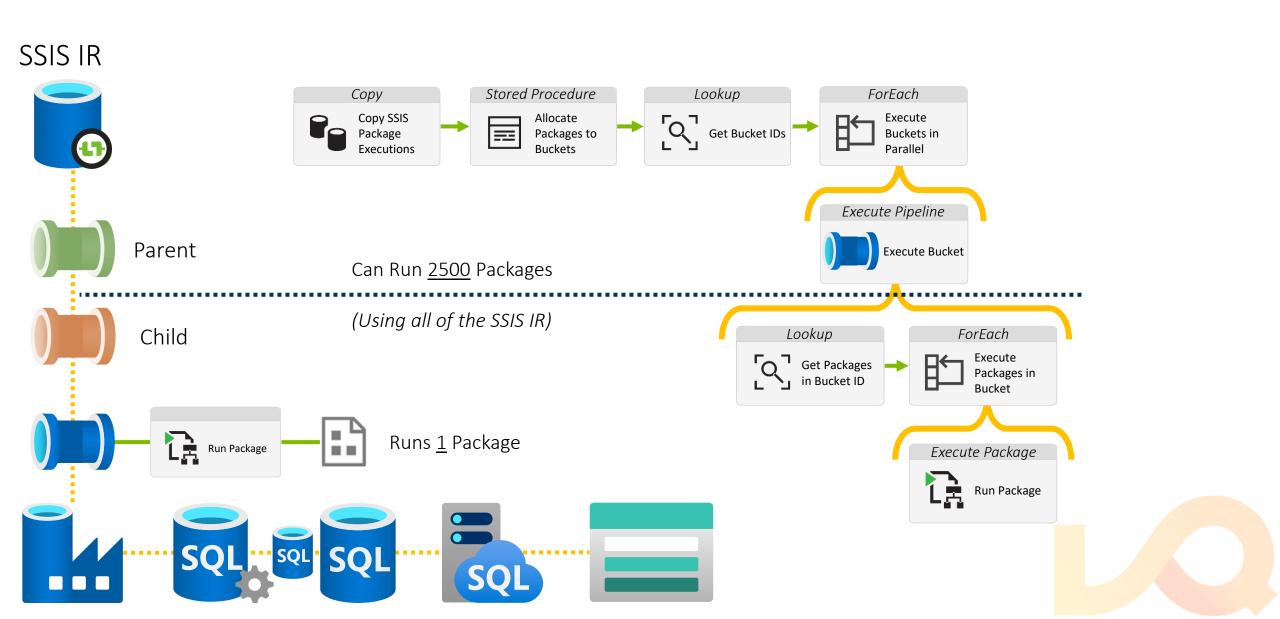
Solution 3: Packages Refactored on a Single Node IR

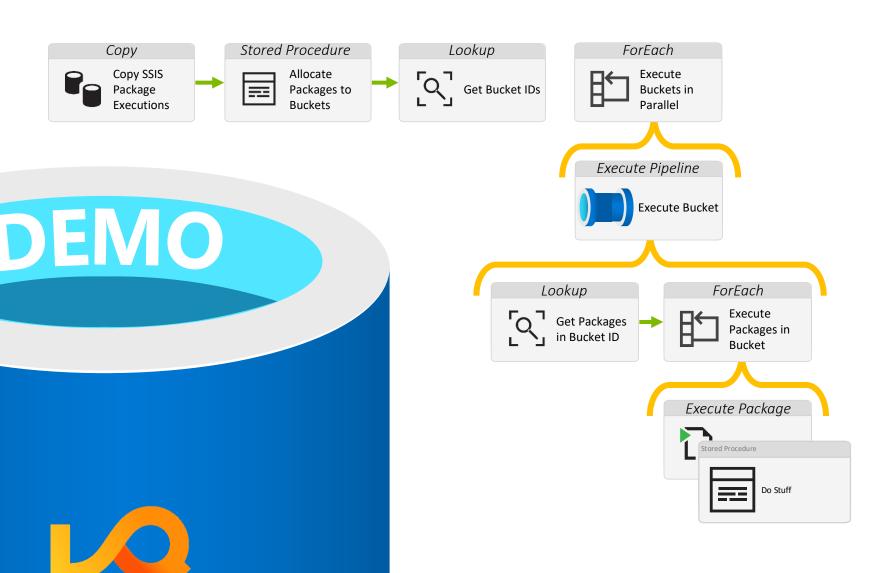


Solution 4: Nested ForEach Activities & Bucket Metadata



Solution 4: Nested ForEach Activities & Bucket Metadata



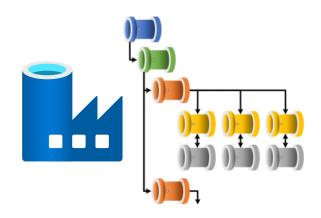




Demo Data Factory and code here:

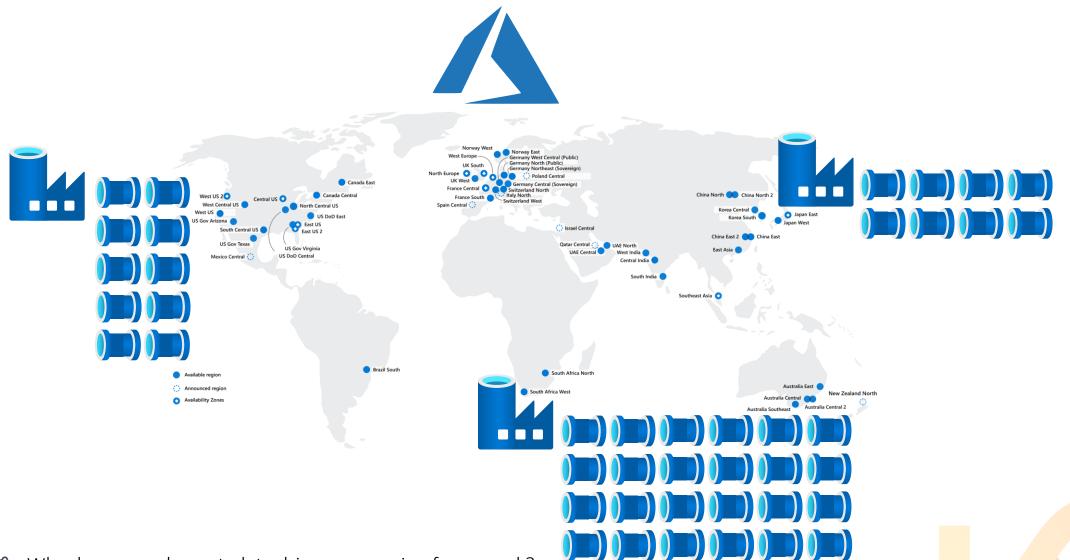
https://github.com/mrpaulandrew/A-Day-Full-of-Azure-Data-Factory

Metadata Driven Pipelines



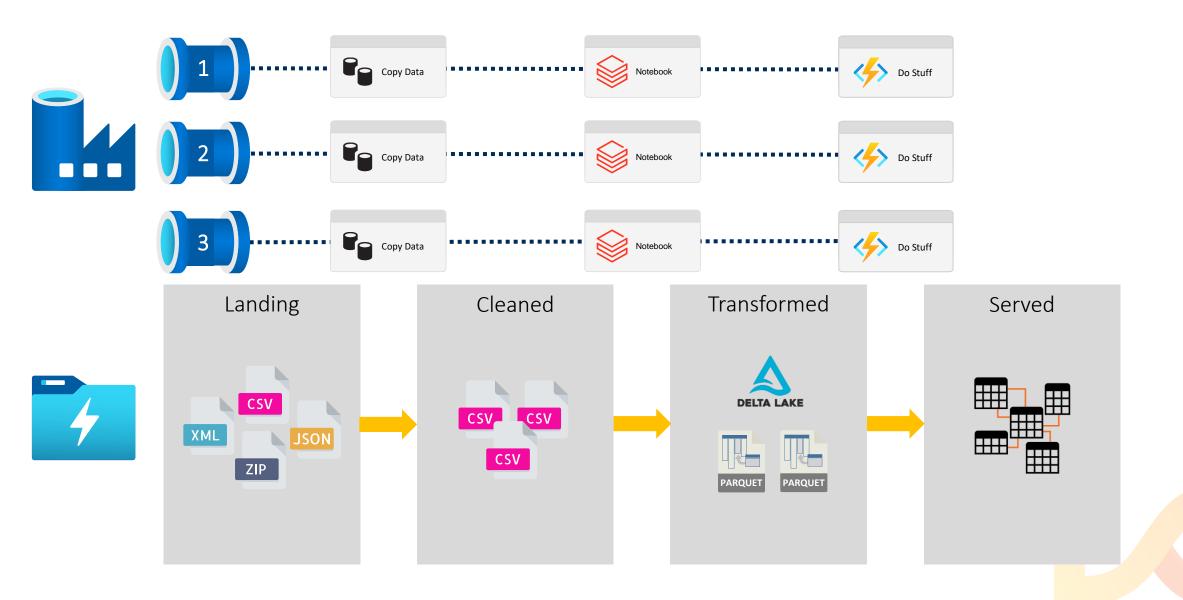


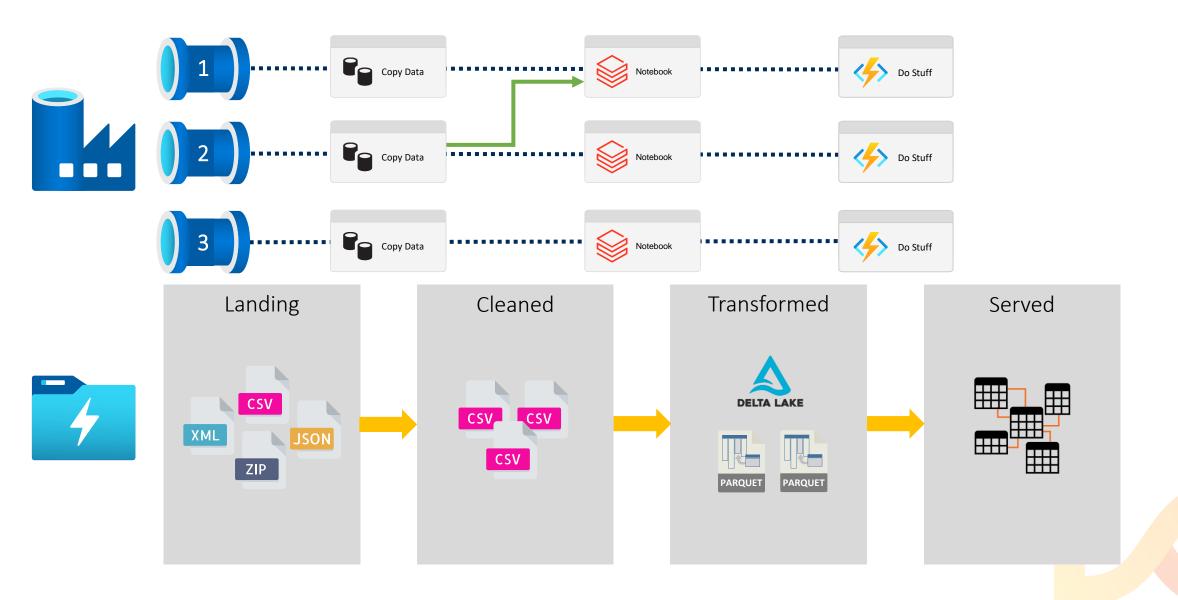
Problem: How should we structure our Data Factory Pipelines?

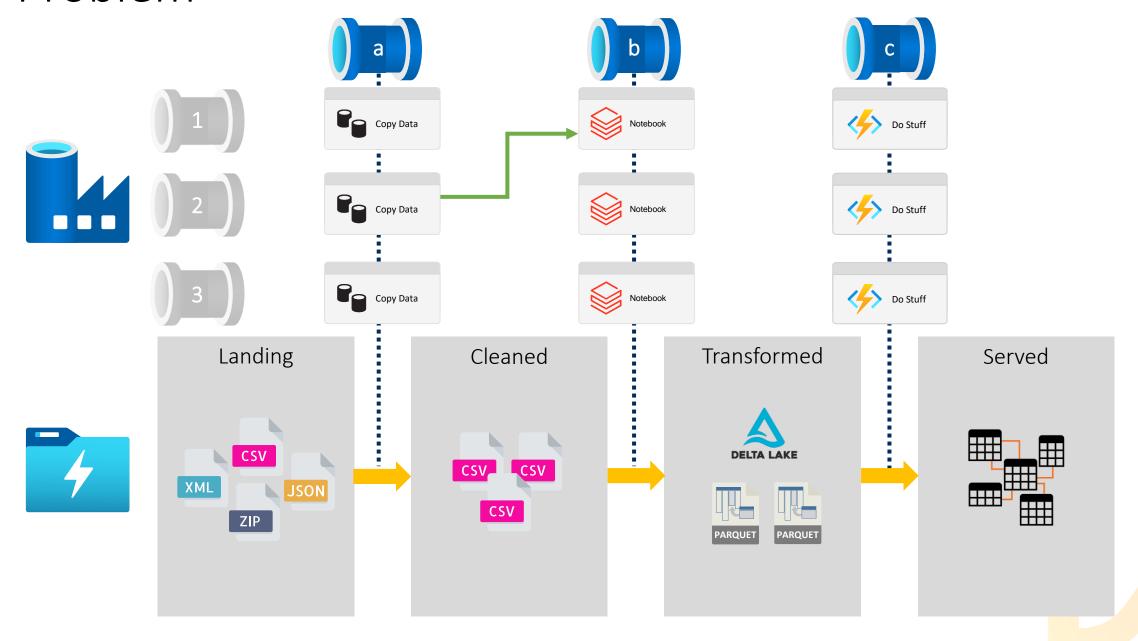


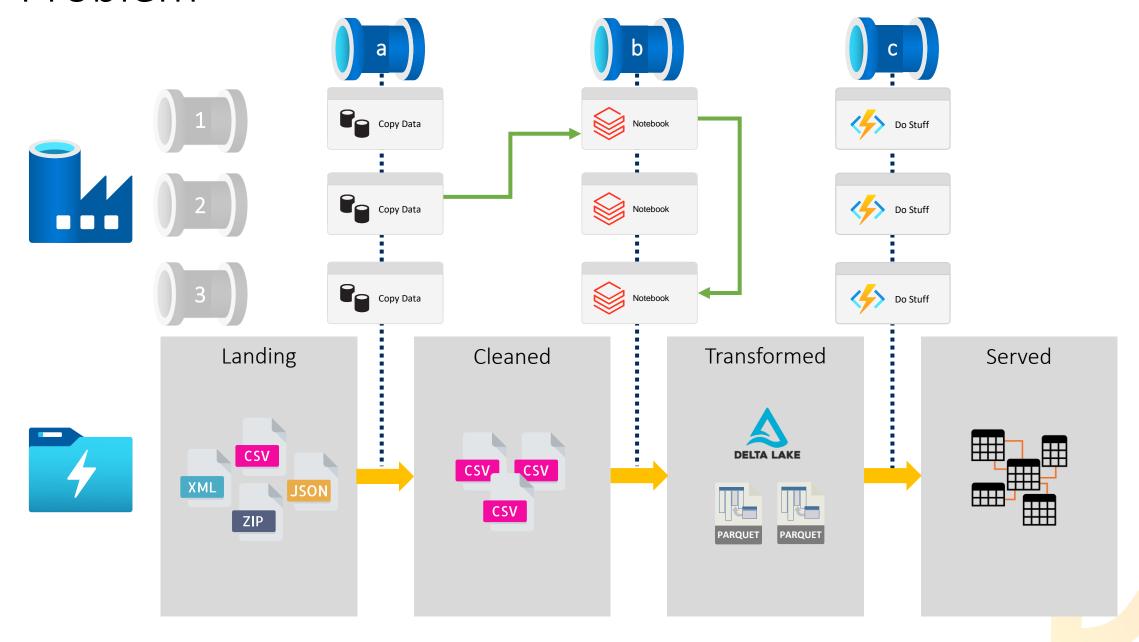


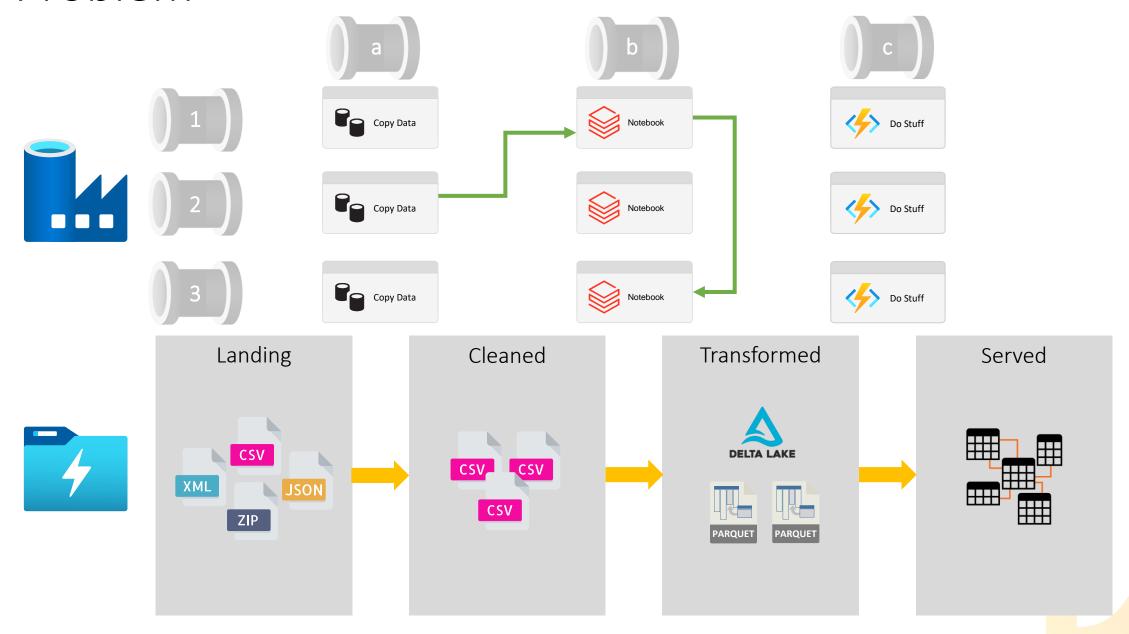
Why do we need a metadata driven processing framework?

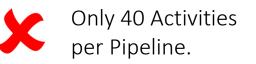


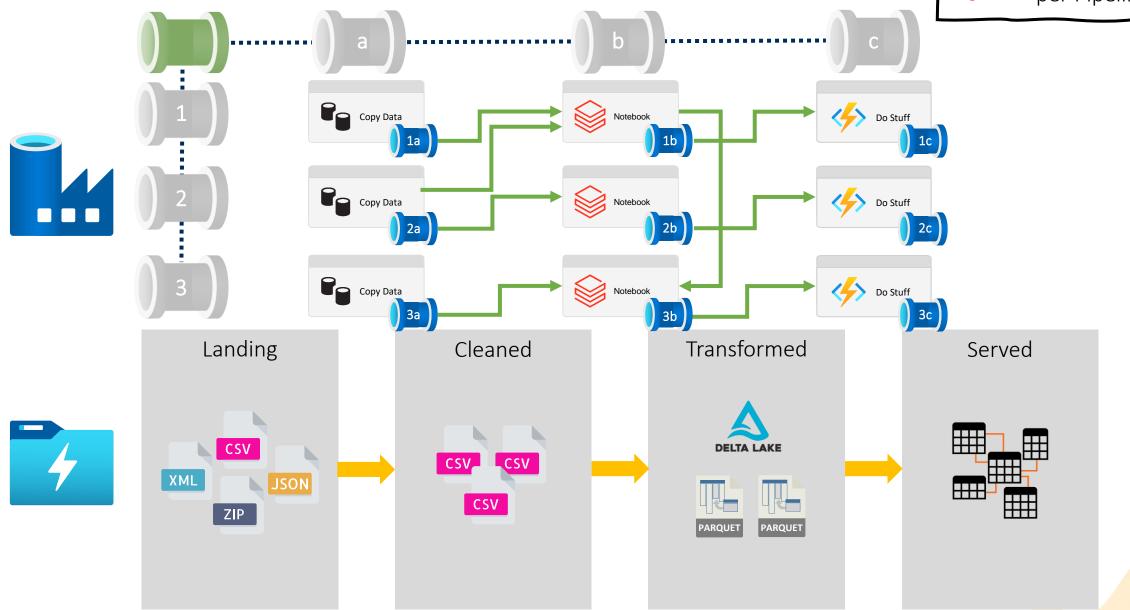


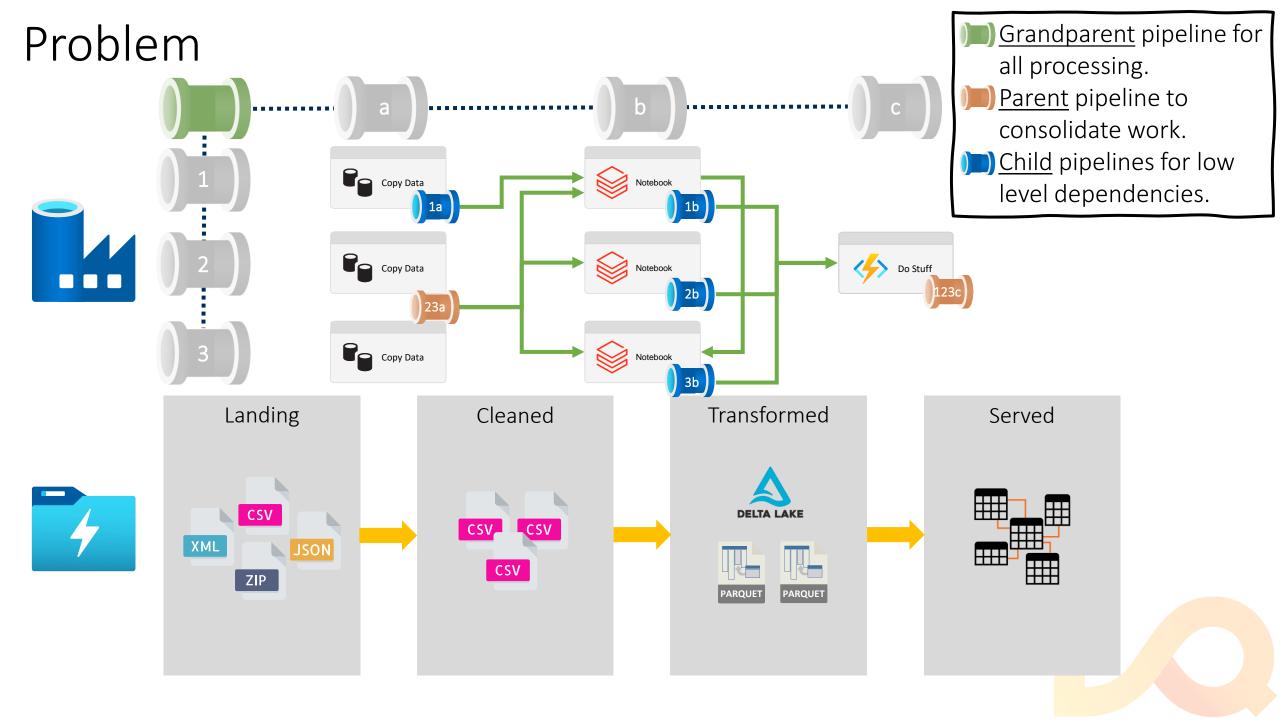




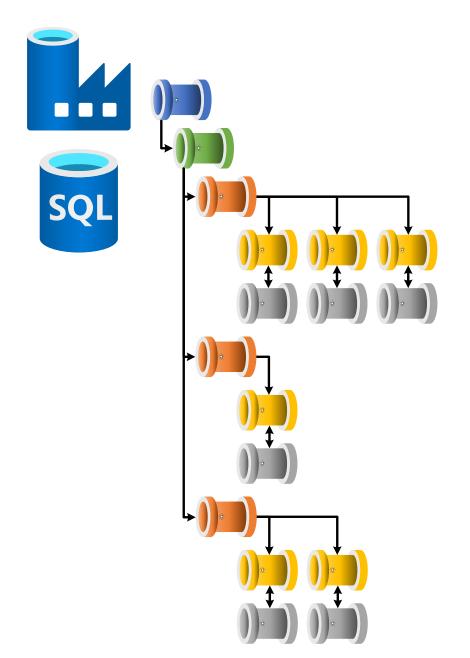






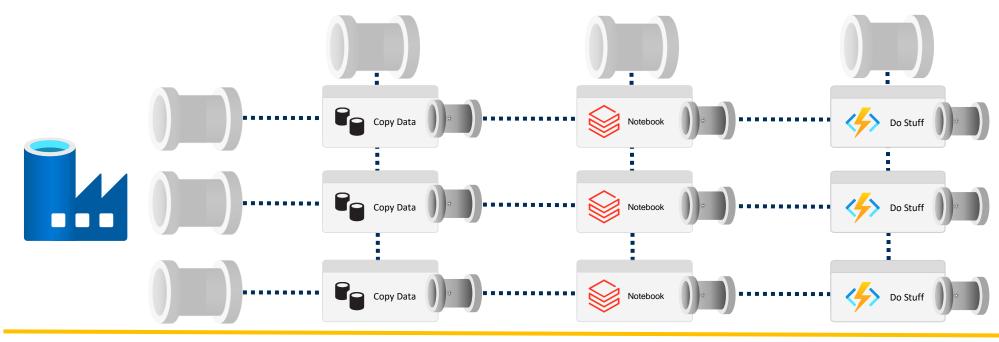


Solution: Use Metadata to Drive Data Factory Pipelines





Solution



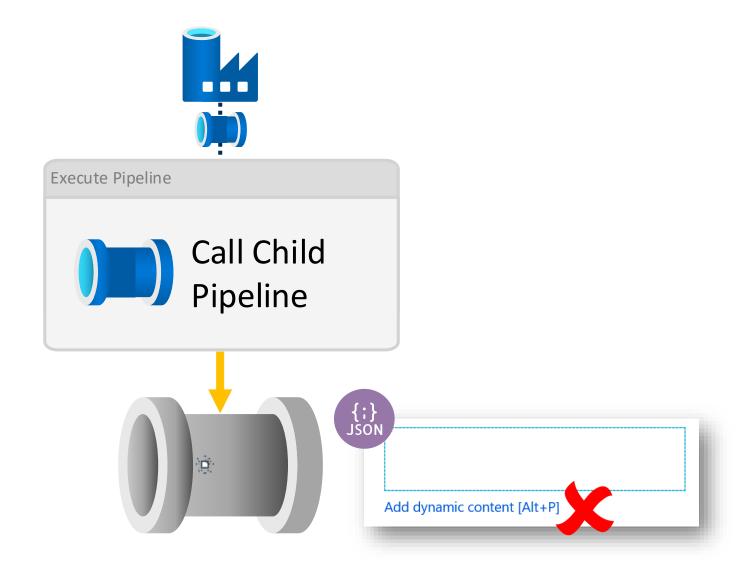


Stages	Pipelines
1	а
2	b
3	С
	d
	е
	f
	g
	h
	i

Stage	Pipeline
1	a
1	b
1	С
2	d
2	e
3	f
3	g
3	h
3	i



One More Problem



Calling Our Worker Pipelines



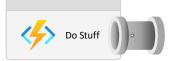






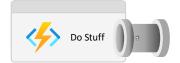














Stages	Pipelines
1	а
2	b
3	С
	d
	е
	f
	g
	h
	i

Stage	Pipeline
1	а
1	b
1	С
2	d
2	е
3	f
3	g
3	h
3	i

Option 1:



Option 2:



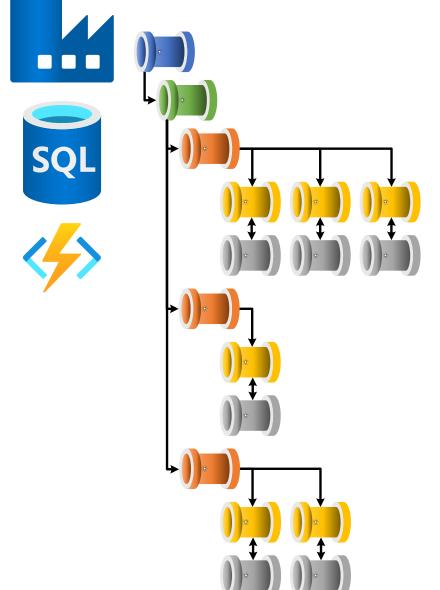
Option 3:



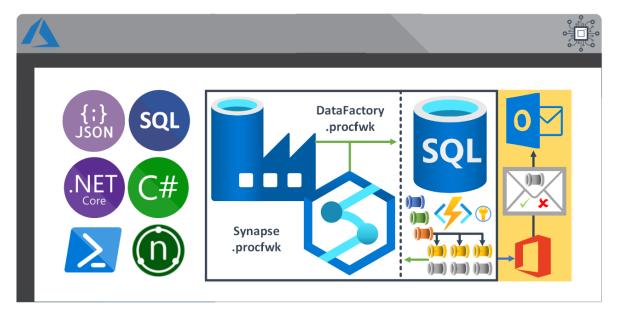


Solution: Use Metadata to Drive Data Factory Pipelines &

Functions to Handle the Worker Pipeline Interactions



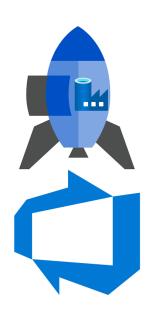
procfwk.com



github.com/mrpaulandrew/procfwk

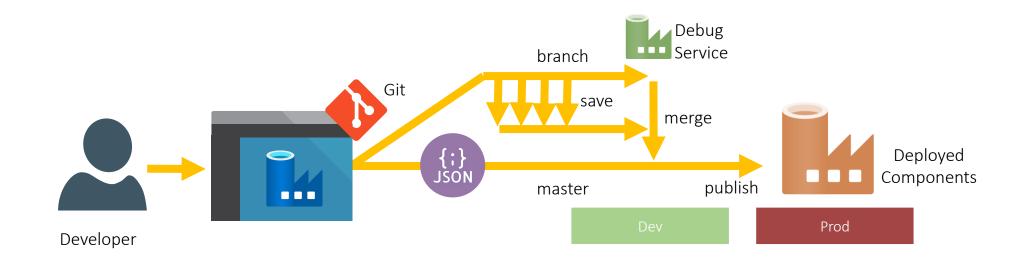


Deploying Data Factory





Deploying Data Factory

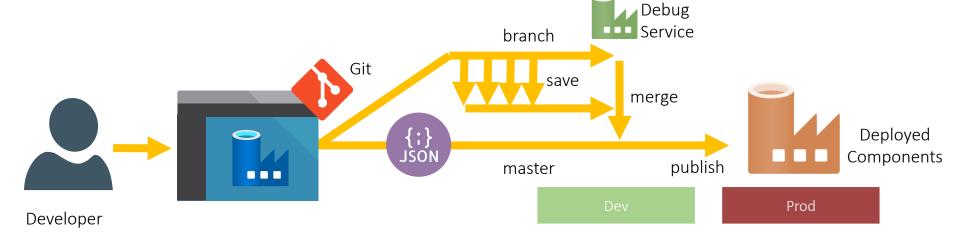




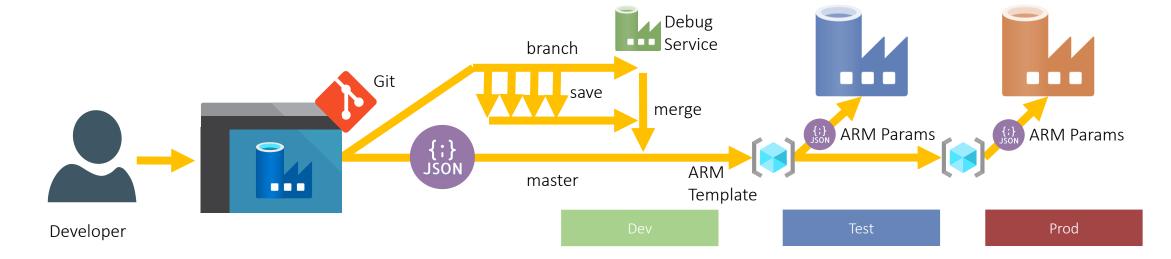
Deploying Data Factory

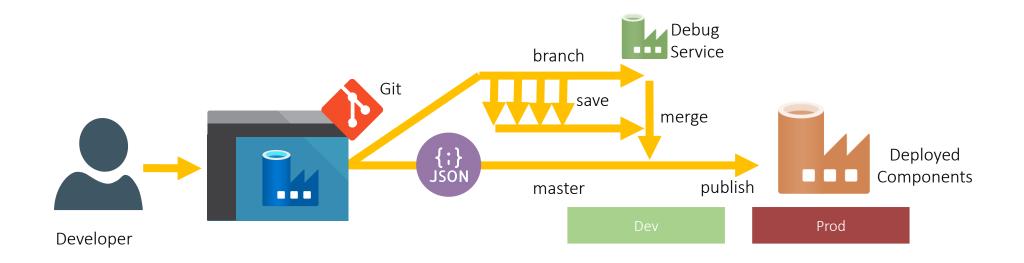
Option 1 – Single Data Factory Service



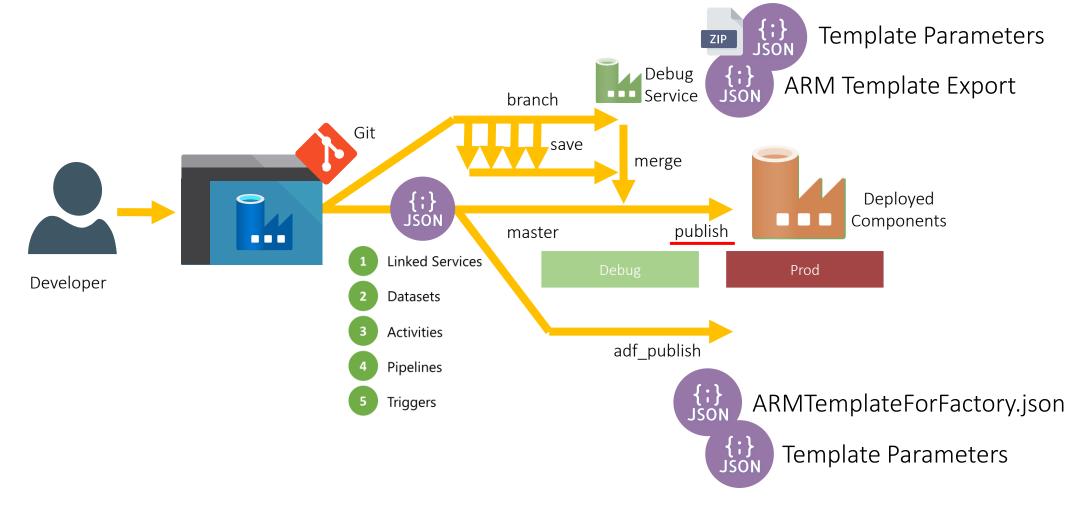


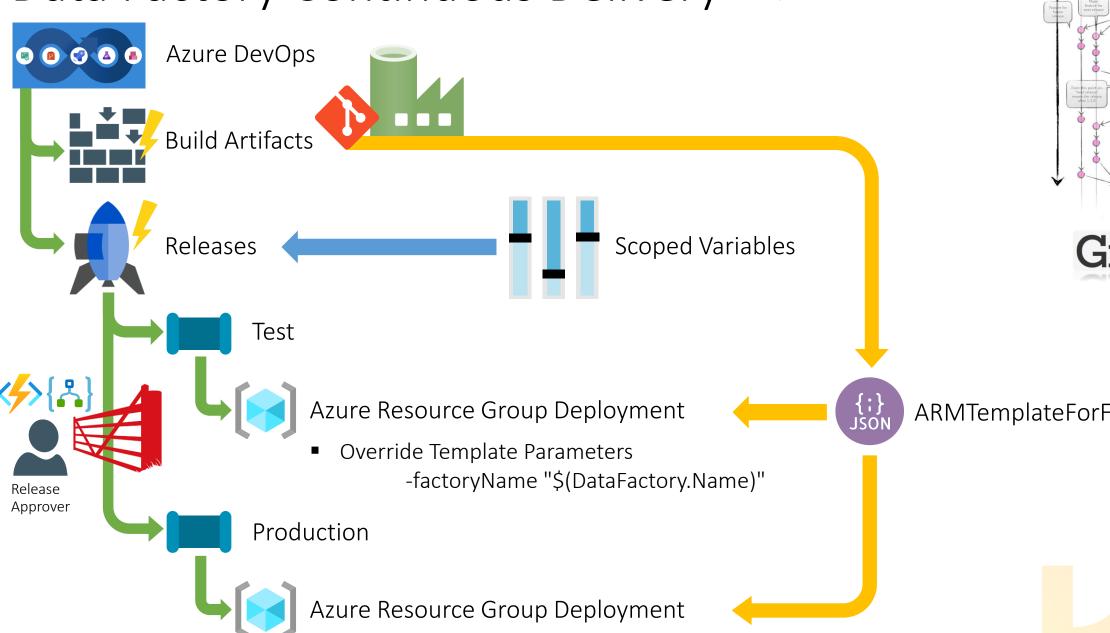
Option 2 – ARM Templates for Multiple Data Factory Services

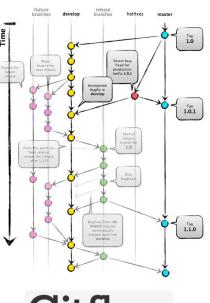




Getting Our ADF Source Code

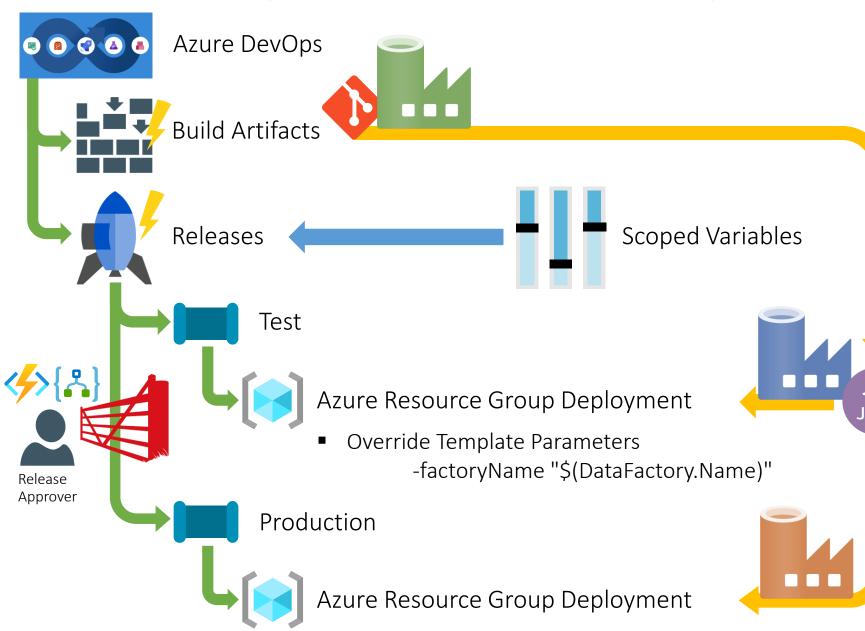










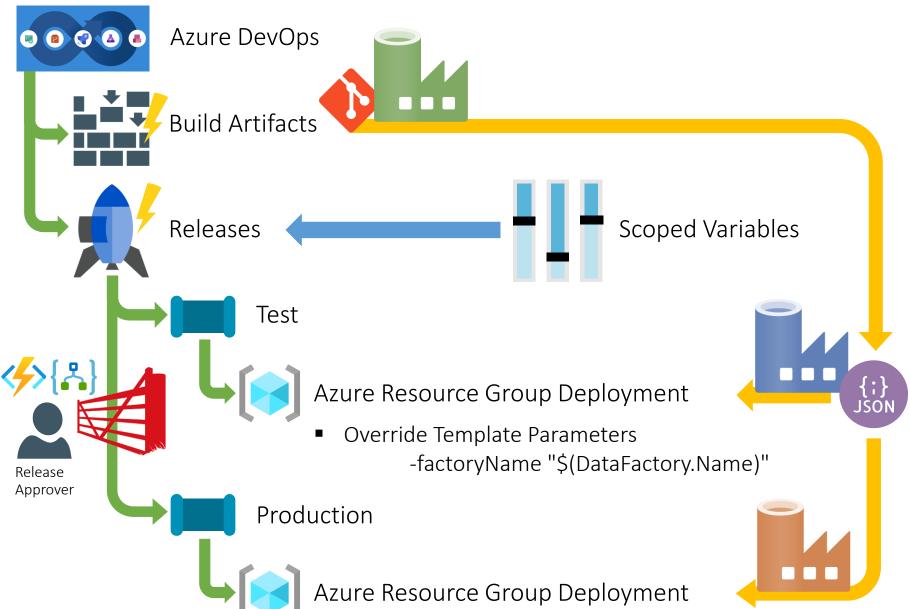


- 1 Linked Services
- 2 Datasets
- 3 Activities
- 4 Pipelines
- 5 Triggers

ARMTemplateForFactory.json



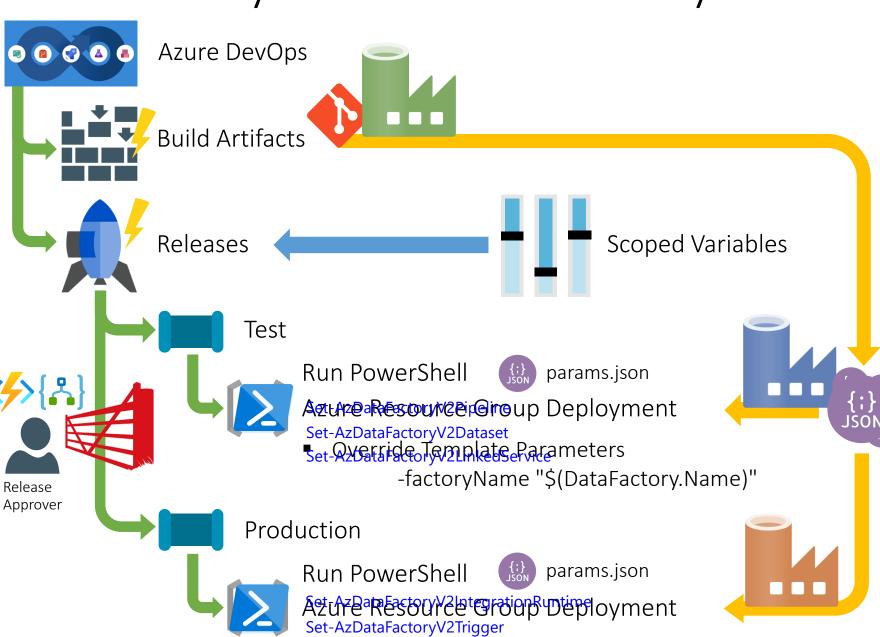






ARMTemplateForFactory.json

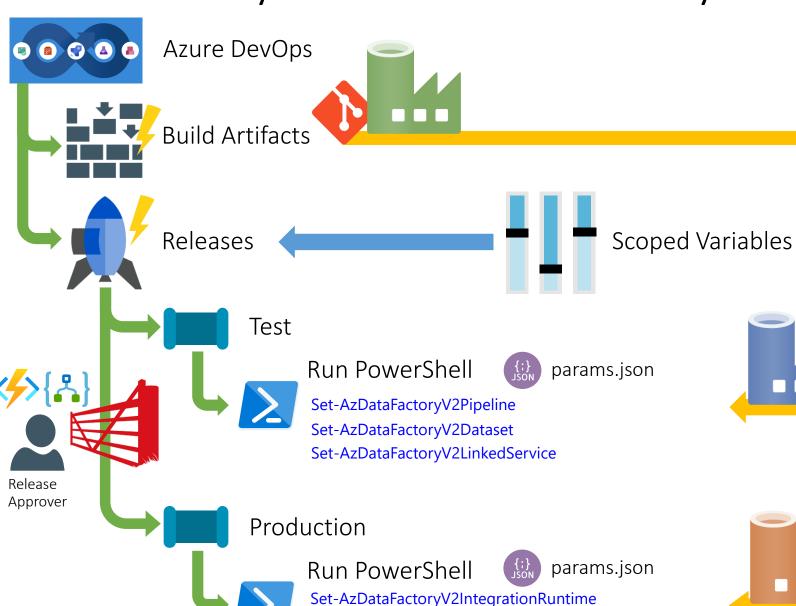




- 1 Linked Services
- 2 Datasets
- 3 Activities
- 4 Pipelines
- 5 Triggers

linkedservices.json
nipelines & activites.json
ARVITEMPIATEFORFACTORY.json
datasets.json
triggers.json
mdfstæublish





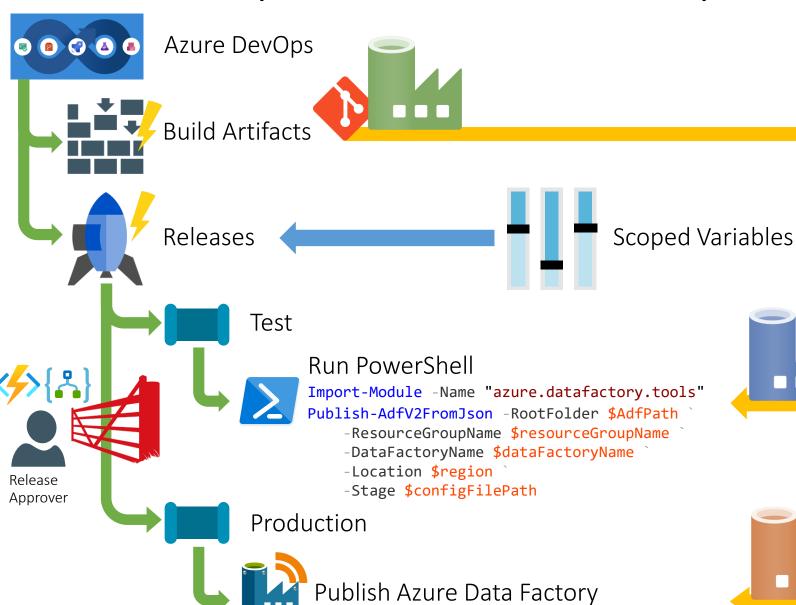
Set-AzDataFactoryV2Trigger

- **Linked Services**
- **Datasets**
- **Activities**
- **Pipelines**
- **Triggers**

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

1) Handle own dependencies.

2) Handle own removals.



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

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



Deployment Options Summary

Option 1 – Use a single Data Factory service.

Option 2 – ARM Templates for multiple Data Factory services (environments).

Option 3 – Use PowerShell cmdlets for each Data Factory JSON artifact.

Option 4 – Use a PowerShell module or custom Azure DevOps task.



Data Factory DevOps Story Summary

What is your code branching strategy?

Feature:

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?







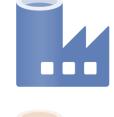




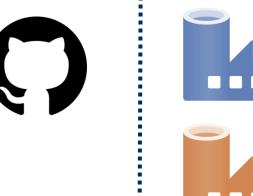
linkedservices.json activites.json datasets.json triggers.json



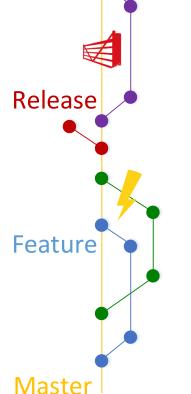












Best Practices





Key Points



Deployments

Mautomated <u>Testing</u>

Maming Conventions

© Pipeline <u>Hierarchies</u>

Impripeline & Activity Descriptions

DDFactory Component Folders

Inked Service Security via Azure Key Vault

Dynamic Linked Services

MGeneric Datasets

Metadata Driven Processing

© Parallel Execution

MHosted Integration Runtimes

MAzure Integration Runtimes

Wider <u>Platform Orchestration</u>

Custom Error Handler Paths

Monitoring via Log Analytics

Service <u>Limitations</u>

WUsing Templates

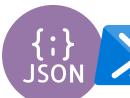
DDocumentation











```
unning checks for Data Factory ARM template:
D:\Stuff\arm_template2.json
Running check... Pipeline(s) without any triggers attached. Directly or indirectly.
Running check... Pipeline(s) with an impossible AND/OR activity execution chain.
Running check... Pipeline(s) without a description value.
Running check... Pipeline(s) not organised into folders.
Running check... Pipeline(s) without annotations.
Running check... Data Flow(s) without a description value.
Running check... Activitie(s) with timeout values still set to the service default value of 7 days.
Running check... Activitie(s) without a description value.
Running check... Activitie(s) ForEach iteration without a batch count value set.
Running check... Activitie(s) ForEach iteration with a batch count size that is less than the service maximum.
Running check... Linked Service(s) not using Azure Key Vault to store credentials.
Running check... Linked Service(s) not used by any other resource.
Running check... Linked Service(s) without a description value.
Running check... Linked Service(s) without annotations.
Running check... Dataset(s) not used by any other resource.
Running check... Dataset(s) without a description value.
Running check... Dataset(s) not organised into folders.
Running check... Dataset(s) without annotations.
Running check... Trigger(s) not used by any other resource.
Running check... Trigger(s) without a description value.
Running check... Trigger(s) without annotations.
Results Summary:
Checks ran against template: 21
Checks with issues found: 21
Total issue count: 264
Issue Count Check Details
                                                                                                         Severity
           Pipeline(s) without any triggers attached. Directly or indirectly.
           Pipeline(s) with an impossible AND/OR activity execution chain.
           Pipeline(s) without a description value.
           Pipeline(s) not organised into folders.
           Pipeline(s) without annotations.
           Data Flow(s) without a description value.
           Activitie(s) with timeout values still set to the service default value of 7 days.
           Activitie(s) without a description value.
           Activitie(s) ForEach iteration without a batch count value set.
           Activitie(s) ForEach iteration with a batch count size that is less than the service maximum. Medium
           Linked Service(s) not using Azure Key Vault to store credentials.
           Linked Service(s) not used by any other resource.
           Linked Service(s) without a description value.
           Linked Service(s) without annotations.
           Dataset(s) not used by any other resource.
                                                                                                         Medium
           Dataset(s) without a description value.
           Dataset(s) not organised into folders.
           Dataset(s) without annotations.
            Trigger(s) not used by any other resource.
            Trigger(s) without a description value.
           Trigger(s) without annotations.
```

Thank you for listening...

Paul Andrew





Blog: mrpaulandrew.com

YouTube: c/mrpaulandrew

Email: paul@mrpaulandrew.com

Twitter: @mrpaulandrew

LinkedIn: In/mrpaulandrew

/CommunityEvents github: github.com/mrpaulandrew /ContentCollateral

/procfwk



STRATEGIC PARTNER



GOLD SPONSOR



Future Processing

SILVER SPONSOR





BRONZE SPONSOR

