# Module 6 – Execution Parallelism
## Control Flow Scale Out

Cloud Formations

# Data Factory Core Components

XML
CSV
JSON
ZIP

SQL

1 Linked Services

2 Datasets

3 Activities

4 Pipelines

5 Triggers

6 IR Compute

Azure IR

Hosted IR

SSIS IR

# For Each

## Iterating over other 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

# For Each Activity

## Iterating over other control flow activities

Many Values
[array]

ForEach

Lookup

* Running in Debug mode.

Execute Pipeline

Call Child Pipeline

@item().

IsSequential: true

[array]

[0]

[1]

[2]

[3]

[i]

[array]

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

Batch Count Default: 20

Batch Count Max: 50

# Problem: Using All Of The SSIS IR Compute

SSIS IR

Supports <u>80</u> Concurrent Packages

MAXDOP = 80

Run Package

Runs <u>1</u> Package

Parent Package

Child Packages
x80

Pipeline x1

ForEach

ForEach
Max Batch
(50)

Run Package

Pipeline x1

Run Package   Run Package   Run Package   Run Package

Activities x80

*April 2024

Cloud Formations - Knowledge Transfer & Training

© 2024 Cloud Formations Ltd

# Nested ForEach Activities & Bucket Metadata

SSIS IR

Copy SSIS Package Executions → Allocate Packages to Buckets → Get Bucket IDs → Execute Buckets in Parallel

[Buckets]

Execute Bucket

Get Packages in Bucket ID → Execute Packages in Bucket

@item().BucketId

[BucketContents]

Run Package

# A General Pattern for Scaling Out

Azure IR

Cloud Formations - Knowledge Transfer & Training

Copy SSIS Package Executions → Allocate Packages to Buckets → Get Bucket IDs → Execute Buckets in Parallel

[Buckets]

Execute Bucket

Get Packages in Bucket ID → Execute Packages in Bucket

@item().BucketId

[BucketContents]

Stored Procedure

Do Stuff

SQL

# Module 6 – Execution Parallelism
## Concurrency Limits vs
## Internal vs External Activities

Cloud Formations

# Data Factory Limitations

https://github.com/MicrosoftDocs/azure-docs/blob/main/includes/azure-data-factory-limits.md

800 Data Factory Instances per Subscription

80 Activities per Data Factory Pipeline

3 Active Data Flow Debug Sessions per Data Factory

5,000 Rows or 4MB of Data Returned per Lookup (No Error if More)

Minimum Tumbling Window Trigger – 15mins

4min Client Response Timeout Using Azure Functions Activity

5,000 Entities (Components) per Data Factory Instance

Cloud Formations - Knowledge Transfer & Training

© 2024 Cloud Formations Ltd

# Data Factory Core Components

1  Linked Services

2  Datasets

3  Activities

4  Pipelines

5  Triggers

6  IR Compute

Azure IR

Hosted IR

SSIS IR

# Compute Concurrency

Orchestrator

Runtime

Fixed Region

Flexible Region

*AutoResolveIntegrationRuntime*

Runtime 3

Runtime 1

Runtime 2

## Edit linked service

Azure SQL Database  Learn more

ⓘ  To avoid publishing immediately to Data Factory, please use Azure Key Vault to retrieve secrets securely. Learn more here

Name *

AzureSqlDatabase1

Description

Connect via integration runtime * ⓘ

AutoResolveIntegrationRuntime

**Connection string**    Azure Key Vault

Source LS

Sink LS

Copy Files

Source

Sink

© 2024 Cloud Formations Ltd

# Compute Concurrency

Orchestrator

Runtime

Fixed Region

Flexible Region

AutoResolve
IntegrationRuntime

Runtime 1

Runtime 2

Runtime 3

Available region

Announced region

Availability Zones

Internal

External

Cloud Formations - Knowledge Transfer & Training

# Concurrency – Pipelines vs Activities

10,000

Internal

1,000

External

3,000

Cloud Formations - Knowledge Transfer & Training

# Concurrency – Pipelines vs Activities

Example

ForEach

Execute 50

Wait — Wait 1 Hour (repeated across grid)

$$1 + (50 \times 39) = 1951$$

1 IR | 1 Pipeline | 1 ForEach + ( 50 Batches X 39 Wait Activities ) = 1951 Concurrent Activities

# Module 6 – Execution Parallelism

## Metadata Driven Frameworks

Cloud Formations

# Problem

How should we structure and trigger our Integration Pipelines?



*Pipeline-asaurus!*

# Problem



Landing · Cleaned · Transformed · Served

# Problem

**1** ---- Copy Data ---- Notebook ---- Do Stuff

**2** ---- Copy Data ---- Notebook ---- Do Stuff

**3** ---- Copy Data ---- Notebook ---- Do Stuff

## Landing
CSV
XML
JSON
ZIP

## Cleaned
CSV  CSV
CSV

## Transformed
DELTA LAKE
PARQUET  PARQUET

## Served

# Problem

a

b

c

1

2

3

Copy Data

Copy Data

Copy Data

Notebook

Notebook

Notebook

Do Stuff

Do Stuff

Do Stuff

## Landing

CSV

XML

JSON

ZIP

## Cleaned

CSV    CSV

CSV

## Transformed

DELTA LAKE

PARQUET    PARQUET

## Served

# Problem

**a**

| 1 | | Copy Data |
| 2 | | Copy Data |
| 3 | | Copy Data |

**b**

Notebook

Notebook

Notebook

**c**

Do Stuff

Do Stuff

Do Stuff

## Landing

CSV
XML
JSON
ZIP

## Cleaned

CSV  CSV
CSV

## Transformed

DELTA LAKE

PARQUET  PARQUET

## Served

# Problem



Landing

Cleaned

Transformed

Served

# Problem

Cloud Formations - Knowledge Transfer & Training

| a | b | c |

1 — Copy Data — 1a → Notebook — 1b → Do Stuff — 1c

2 — Copy Data — 2a → Notebook — 2b → Do Stuff — 2c

3 — Copy Data — 3a → Notebook — 3b → Do Stuff — 3c

## Landing
CSV
XML
JSON
ZIP

## Cleaned
CSV CSV
CSV

## Transformed
DELTA LAKE
PARQUET PARQUET

## Served

# Problem

Landing

Cleaned

Transformed

Served

Grandparent pipeline for all processing.
Parent pipeline to consolidate work.
Child pipelines for low level dependencies.

Copy Data  1a
Copy Data  23a
Copy Data

Notebook  1b
Notebook  2b
Notebook  3b

Do Stuff  123c

CSV
XML
ZIP
JSON

CSV  CSV
CSV

DELTA LAKE
PARQUET  PARQUET

# Solution

Use metadata to drive all pipeline executions

# Solution



| Stages | Pipelines |
|--------|-----------|
| 1 | a |
| 2 | b |
| 3 | c |
| | d |
| | e |
| | f |
| | g |
| | h |
| | i |

| Stage | Pipeline |
|-------|----------|
| 1 | a |
| 1 | b |
| 1 | c |
| 2 | d |
| 2 | e |
| 3 | f |
| 3 | g |
| 3 | h |
| 3 | i |

Cloud Formations - Knowledge Transfer & Training

# Calling Our Worker Pipelines

One More Problem to Consider



Execute Pipeline

Call Child
Pipeline

{;}
JSON

Add dynamic content [Alt+P]

Cloud Formations - Knowledge Transfer & Training

# Calling Our Worker Pipelines

Copy Data

Notebook

Do Stuff

Copy Data

Notebook

Do Stuff

Copy Data

Notebook

Do Stuff

**Option 1:**

Web
HTTP Run Worker

**Option 2:**

Web
HTTP Run Worker

**Option 3:**

Azure Function
Run Worker

| Stages | Pipelines |
|--------|-----------|
| 1 | a |
| 2 | b |
| 3 | c |
| | d |
| | e |
| | f |
| | g |
| | h |
| | i |

| Stage | Pipeline |
|-------|----------|
| 1 | a |
| 1 | b |
| 1 | c |
| 2 | d |
| 2 | e |
| 3 | f |
| 3 | g |
| 3 | h |
| 3 | i |

SQL

# Solution

Use Metadata to Drive Data Factory Pipelines & Functions as the middleware to handle worker execution

# Framework Pipeline Hierarchy

- Solution Executor

**Role:** Optional level platform setup, for example, scale up/out compute services ready for the framework to run.

- Batch Executor

**Role:** Execution run wrapper for batches and execution stage iterator.

- Stage Executor

**Role:** Scale out triggering of worker pipelines within the execution stage(s).

- Pipeline Executor

**Role:** Worker validator, executor, monitor and reporting of the outcome for the single worker pipeline.

- Workers & Utilities

**Role:** Anything specific to the process needing to be performed, Ingest and Transform pipelines included.

# Introducing CF.Cumulus

cloudformations.org/cumulus

A cloud data platform product to accelerate time to insights. Our open-source framework is designed for the real world. Stripping away the complexity, giving you the power to build, scale, and manage your dataflows with ease, accelerating data delivery.
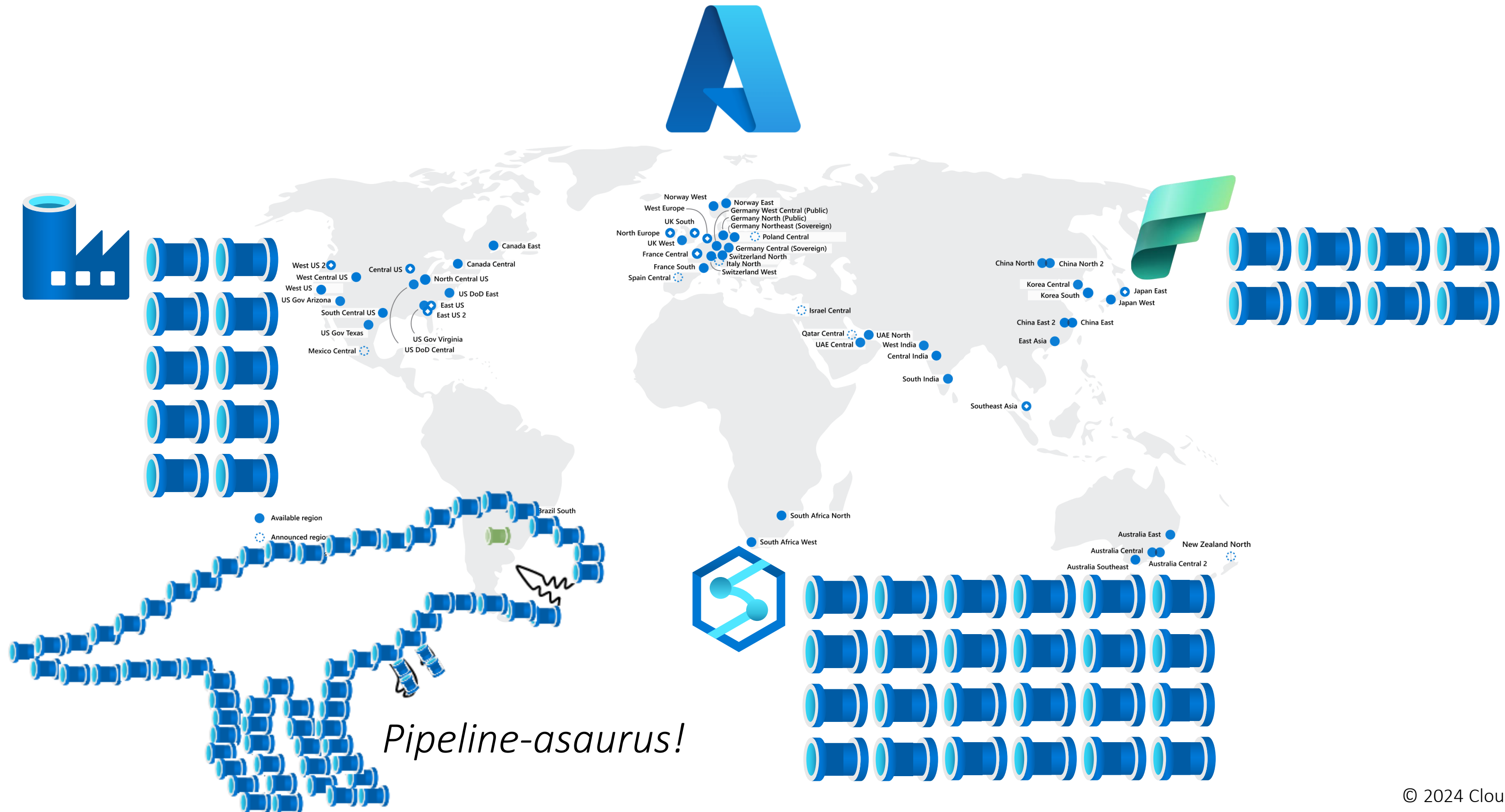


Cloud Formations - Knowledge Transfer & Training

**Sources**
- Files
- Databases
- APIs

**Cumulus**
Cloud Formations

- Ingest
- Control
- Transform

**Lakehouse**
- Raw
- Cleansed
- Curated

**Platform & Plugins**
- Secret Management
- Monitoring & Alerting
- Cost Optimisation
- Identity Management
- Deployment
- Data Catalogue & Governance

**Model**
- Datasets
- Data Marts
- Data Science

**Consume**
- Dashboards
- Data Products
- APIs & Integration

Underlying Technology: Spark, DELTA LAKE, C# {;} JSON

Works With:

# Sources

- Files
- Databases
- APIs

# Cumulus
Cloud Formations

## Ingest

## Control

## Transform

# Lakehouse

- Raw
- Cleansed
- Curated

## Platform & Plugins

- Secret Management
- Monitoring & Alerting
- Cost Optimisation
- Identity Management
- Deployment
- Data Catalogue & Governance

# Model

- Datasets
- Data Marts
- Data Science

# Consume

- Dashboards
- Data Products
- APIs & Integration

DataFactory .procfwk

SQL

Synapse .procfwk

ProcFwk.com

An evolution from metadata driven orchestration to complete data delivery. Accelerating your time to insight. And still open-source.

Cumulus

Cloud Formations

# CF.Cumulus Framework

An Azure product agnostic view of compute, storage and orchestration resources as used by CF.Cumulus.

A primary design focus for CF.Cumulus is to align with open-standards and open technologies. Offering a flexible, plug and play approach to the core components.

In this context, using Apache Spark compute and Delta Lake structures for the storage layer.

# Flexible Technology Choices

Building on the generic visuals above, but now applying a pure Microsoft product perspective to the logical delivery of Cumulus. The following diagrams could apply to a solution depending on the choice of tooling for Orchestration, Compute and Storage with each offering different advantages and coupling of resources.

## Using Azure Synapse Analytics

This implementation represents a <u>partly decoupled</u> architecture, in terms of cloud resources used to delivery Cumulus.

- **Orchestration** and **Compute** delivered using Azure Synapse Analytics.

- **Storage** delivered using Azure Data Lake.

## Using Microsoft Fabric

This implementation represents a <u>coupled architecture</u>, in terms of cloud resources used to delivery Cumulus.

- **Orchestration, Compute** and **Storage** delivered using Microsoft Fabric experiences.

## Using Azure Data Factory & Azure Databricks

**Recommended**

This implementation represents a <u>fully decoupled</u> architecture, in terms of cloud resources used to delivery Cumulus.

- **Orchestration** delivered using Azure Data Factory.
- **Compute** delivered using Azure Databricks.
- **Storage** delivered using Azure Data Lake.

© 2024 Cloud Formations Ltd

# Lakehouse Implementations - Simplified Through Metadata

New source datasets delivered from producers to consumers with ease using only a few lines of metadata and open cloud native tooling.



INSERT INTO [ingest].[Datasets] …
INSERT INTO [transform].[Datasets] …
INSERT INTO [control].[Pipelines] …

# Framework Database

**AlertOutcomes (procfwk)**
- OutcomeBitPosition
- PipelineOutcomeStatus
- BitValue

**PipelineAlertLink (procfwk)**
- AlertId
- PipelineId
- RecipientId
- OutcomesBitValue
- Enabled

**Recipients (procfwk)**
- RecipientId
- Name
- EmailAddress
- MessagePreference
- Enabled

**Batches (procfwk)**
- BatchId
- BatchName
- BatchDescription
- Enabled

**BatchStageLink (procfwk)**
- BatchId
- StageId

**Stages (procfwk)**
- StageId
- StageName
- StageDescription
- Enabled

**Subscriptions (procfwk)**
- SubscriptionId
- Name
- Description
- TenantId

**Tenants (procfwk)**
- TenantId
- Name
- Description

**Pipelines (procfwk)**
- PipelineId
- OrchestratorId
- StageId
- PipelineName
- LogicalPredecessorId
- Enabled

**PipelineParameters (procfwk)**
- ParameterId
- PipelineId
- ParameterName
- ParameterValue
- ParameterValueLastUsed

**PipelineDependencies (procfwk)**
- DependencyId
- PipelineId
- DependantPipelineId

**Orchestrators (procfwk)**
- OrchestratorId
- OrchestratorName
- OrchestratorType
- IsFrameworkOrchestrator
- ResourceGroupName
- SubscriptionId
- Description

**PipelineAuthLink (procfwk)**
- AuthId
- PipelineId
- OrchestratorId
- CredentialId

**ServicePrincipals (dbo)**
- CredentialId
- PrincipalName
- PrincipalId
- PrincipalSecret
- PrincipalIdUrl
- PrincipalSecretUrl

**Properties (procfwk)**
- PropertyId
- PropertyName
- PropertyValue
- Description
- ValidFrom
- ValidTo

**BatchExecution (procfwk)**
- BatchId
- ExecutionId
- BatchName
- BatchStatus
- StartDateTime
- EndDateTime

**CurrentExecution (procfwk)**
- LocalExecutionId
- StageId
- PipelineId
- CallingDataFactoryName
- ResourceGroupName
- DataFactoryName
- PipelineName
- StartDateTime
- PipelineStatus
- LastStatusCheckDateTime
- EndDateTime
- IsBlocked
- AdfPipelineRunId
- PipelineParamsUsed

**ExecutionLog (procfwk)**
- LogId
- LocalExecutionId
- StageId
- PipelineId
- CallingDataFactoryName
- ResourceGroupName
- DataFactoryName
- PipelineName
- StartDateTime
- PipelineStatus
- EndDateTime
- AdfPipelineRunId
- PipelineParamsUsed

**ErrorLog (procfwk)**
- LogId
- LocalExecutionId
- AdfPipelineRunId
- ActivityRunId
- ActivityName
- ActivityType
- ErrorCode
- ErrorType
- ErrorMessage

# Framework Database

**AlertOutcomes (procfwk)**
- 🔑 OutcomeBitPosition
- PipelineOutcomeStatus
- BitValue

**PipelineAlertLink (procfwk)**
- 🔑 AlertId
- PipelineId
- RecipientId
- OutcomesBitValue
- Enabled

**Recipients (procfwk)**
- 🔑 RecipientId
- Name
- EmailAddress
- MessagePreference
- Enabled

**Batches (procfwk)**
- 🔑 BatchId
- BatchName
- BatchDescription
- Enabled

**PipelineParameters (procfwk)**
- 🔑 ParameterId
- PipelineId
- ParameterName
- ParameterValue
- ParameterValueLastUsed

**BatchStageLink (procfwk)**
- 🔑 BatchId
- 🔑 StageId

**Pipelines (procfwk)**
- 🔑 PipelineId
- OrchestratorId
- StageId
- PipelineName
- LogicalPredecessorId
- Enabled

**PipelineDependencies (procfwk)**
- 🔑 DependencyId
- PipelineId
- DependantPipelineId

**Stages (procfwk)**
- 🔑 StageId
- StageName
- StageDescription
- Enabled

**Orchestrators (procfwk)**
- 🔑 OrchestratorId
- OrchestratorName
- OrchestratorType
- IsFrameworkOrchestrator
- ResourceGroupName
- SubscriptionId
- Description

**PipelineAuthLink (procfwk)**
- 🔑 AuthId
- PipelineId
- OrchestratorId
- CredentialId

**Subscriptions (procfwk)**
- 🔑 SubscriptionId
- Name
- Description
- TenantId

**ServicePrincipals (dbo)**
- 🔑 CredentialId
- PrincipalName
- PrincipalId
- PrincipalSecret
- PrincipalIdUrl
- PrincipalSecretUrl

**Tenants (procfwk)**
- 🔑 TenantId
- Name
- Description

**Properties (procfwk)**
- 🔑 PropertyId
- 🔑 PropertyName
- PropertyValue
- Description
- ValidFrom
- ValidTo

**BatchExecution (procfwk)**
- 🔑 BatchId
- 🔑 ExecutionId
- BatchName
- BatchStatus
- StartDateTime
- EndDateTime

**ExecutionLog (procfwk)**
- 🔑 LogId
- LocalExecutionId
- StageId
- PipelineId
- CallingDataFactoryName
- ResourceGroupName
- DataFactoryName
- PipelineName
- StartDateTime
- PipelineStatus
- EndDateTime
- AdfPipelineRunId
- PipelineParamsUsed

**CurrentExecution (procfwk)**
- 🔑 LocalExecutionId
- 🔑 StageId
- 🔑 PipelineId
- CallingDataFactoryName
- ResourceGroupName
- DataFactoryName
- PipelineName
- StartDateTime
- PipelineStatus
- LastStatusCheckDateTime
- EndDateTime
- IsBlocked
- AdfPipelineRunId
- PipelineParamsUsed

**ErrorLog (procfwk)**
- 🔑 LogId
- LocalExecutionId
- AdfPipelineRunId
- ActivityRunId
- ActivityName
- ActivityType
- ErrorCode
- ErrorType
- ErrorMessage

Configuration & Behaviour

Core Metadata

Execution Handling

Location & Authentication

Email Alerting

Runtime & Logging

# CF.Cumulus.Control Activity Chain

## Orchestration Framework

### Grandparent *(Optional)*

- Scale Up Compute
- Environment Setup
- Framework Bootstrap
- Sync Serving Layers
- Scale Down Compute
- Archive Cold Data

### Parent

- Is Parent Already Running
- Execute Precursor
- Check Metadata Integrity
- Check Previous Execution
- Clean Up Previous Run
- Execution Wrapper
- Set Execution Id
- Get Stages
- Execute Stages
- Check Outcome and Update Logs

- Get SPN Details
- Log Pipeline Checking
- Get Pipeline Status
- Set Pipeline Status
- Set Last Check DateTime
- Pipeline Status Unknown
- Pipeline Status Queued
- Pipeline Status Failed
- Pipeline Status InProgress - Running
- Pipeline Status Succeeded
- Pipeline Status Cancelled

- Check and Update Blockers
- Log Stage Preparing
- Stage Executor

### Child

- Get Pipelines
- Execute Pipelines
- Worker Executor

### Infant

- Get Worker Core Details
- Capture Worker Core Detail as an Array
- Validate Pipeline
- Log Pipeline Validating
- Is Target Worker Validate
- Update Execution With Invalid Worker
- Throw Exception – Invalid Worker
- Log Validate Function Activity Failure
- Get Pipeline Params
- Log Pipeline Running
- Execute Pipeline
- Log Execute Function Activity Failure
- Set Run Id
- Update Run Id
- Get Wait Duration
- Wait Until Pipeline Completes
- Set Pipeline Result
- Get Worker Pipeline Status
- Set Last Check DateTime
- Set Worker State
- Wait if Running
- Wait for Pipeline
- Log Check Function Activity Failure
- Pipeline Status Unknown
- Pipeline Status Cancelled
- Pipeline Status Succeded
- Get Error Details
- Log Error Details
- Pipeline Status Failed
- Check For Alerts
- Send Alerts
- Get Email Parts
- Call Email Sender

### Utilities - Check For Running Pipeline

- Get Framework Orchestrator Details
- Set Subscription Id
- Set Resource Group Name
- Set Orchestrator Type
- Switch For Orchestrator Type
- Get Query Run Days Value
- Set Orchestrator Type
- Throw Exception – Invalid Worker
- Check for Valid ADF Pipeline Name
- Get ADF Pipeline Runs
- Set ADF Runs Output
- Check for Valid SYN Pipeline Name
- Get SYN Pipeline Runs
- Set SYN Runs Output
- Filter Running Pipelines
- If Using Batch Executions
- Get Execution Batch Status
- Filter for Batch Name
- Set Run Count for Batch
- Throw Exception – Invalid Worker
- Set Run Count Without Batch
- If Pipeline Is Running

### Utilities – Throw Exception

- Raise Error

### Utilities – Email Sender

- Send Email

## Worker Pipelines

- Worker 1 - Extract
- Worker 2 - Clean
- Worker 3 - Transform
- Worker 4 - Load
- Worker 5 - Serve
- Worker n - ..........

# CF.Cumulus Framework

An Azure product agnostic view of compute, storage and orchestration resources as used by CF.Cumulus.

A primary design focus for CF.Cumulus is to align with open-standards and open technologies. Offering a flexible, plug and play approach to the core components.

In this context, using Apache Spark compute and Delta Lake structures for the storage layer.



Data Storage

Producers

Ingestion Compute

Raw

PARQUET
CSV
JSON
TXT

Cleansed

DELTA LAKE

Curated

DELTA LAKE

Consumers

Transformation Compute

Metadata

Worker Pipelines

Middleware Functions

Bootstrap Pipelines

**Ingest**

**Transform**

**Control**

# Problem

How should we structure and trigger our Integration Pipelines?



*Pipeline-asaurus!*

# Solution

## Use Metadata to Drive Integration Pipeline execution

Canada East

West US 2
West Central US
Central US
West US
US Gov Arizona
South Central US
East US
East US 2
US Gov Texas
US Gov Virginia
US DoD Central
Mexico Central

Canada Central
North Central US
US DoD East

Brazil South

Norway West
West Europe
UK South
North Europe
UK West
France Central
France South
Spain Central

Norway East
Germany West Central (Public)
Germany North (Public)
Germany Northeast (Sovereign)
Poland Central
Germany Central (Sovereign)
Switzerland North
Italy North
Switzerland West

Israel Central

Qatar Central
UAE Central
UAE North
West India
Central India

South India

Southeast Asia

China North
China North 2
Korea Central
Korea South
China East 2
China East
East Asia
Japan East
Japan West

South Africa North

South Africa West

Australia East
Australia Central
Australia Southeast
Australia Central 2
New Zealand North

● Available region
○ Announced region
◈ Availability Zones

### Cumulus
Cloud Formations

Ingest    Transform

Control

# Solution

## Use Metadata to Drive Integration Pipeline execution

Tenant 1

Subscription 1

Resource Group 1

Tenant 3

Subscription 3

Resource Group 3

Tenant 2

Subscription 2

Resource Group 2

**Cloud Formations - Knowledge Transfer & Training**

### Cumulus
Cloud Formations

- Ingest
- Transform
- Control

West US 2
West Central US
West US
US Gov Arizona
South Central US
US Gov Texas
Mexico Central
Central US
North Central US
East US
East US 2
US Gov Virginia
US DoD Central
US DoD East
Canada East
Canada Central

Brazil South

Norway West
West Europe
UK South
North Europe
UK West
France Central
France South
Spain Central
Norway East
Germany West Central (Public)
Germany North (Public)
Germany Northeast (Sovereign)
Germany Central (Sovereign)
Switzerland North
Italy North
Switzerland West
Poland Central

Israel Central

Qatar Central
UAE Central
UAE North
West India
Central India
South India

China North
China North 2
Korea Central
Korea South
China East 2
China East
East Asia
Japan East
Japan West

Southeast Asia

South Africa North
South Africa West

Australia East
Australia Central
Australia Southeast
Australia Central 2
New Zealand North

- Available region
- Announced region
- Availability Zones

© 2024 Cloud Formations Ltd

# Hub & Spoke Integration Architecture
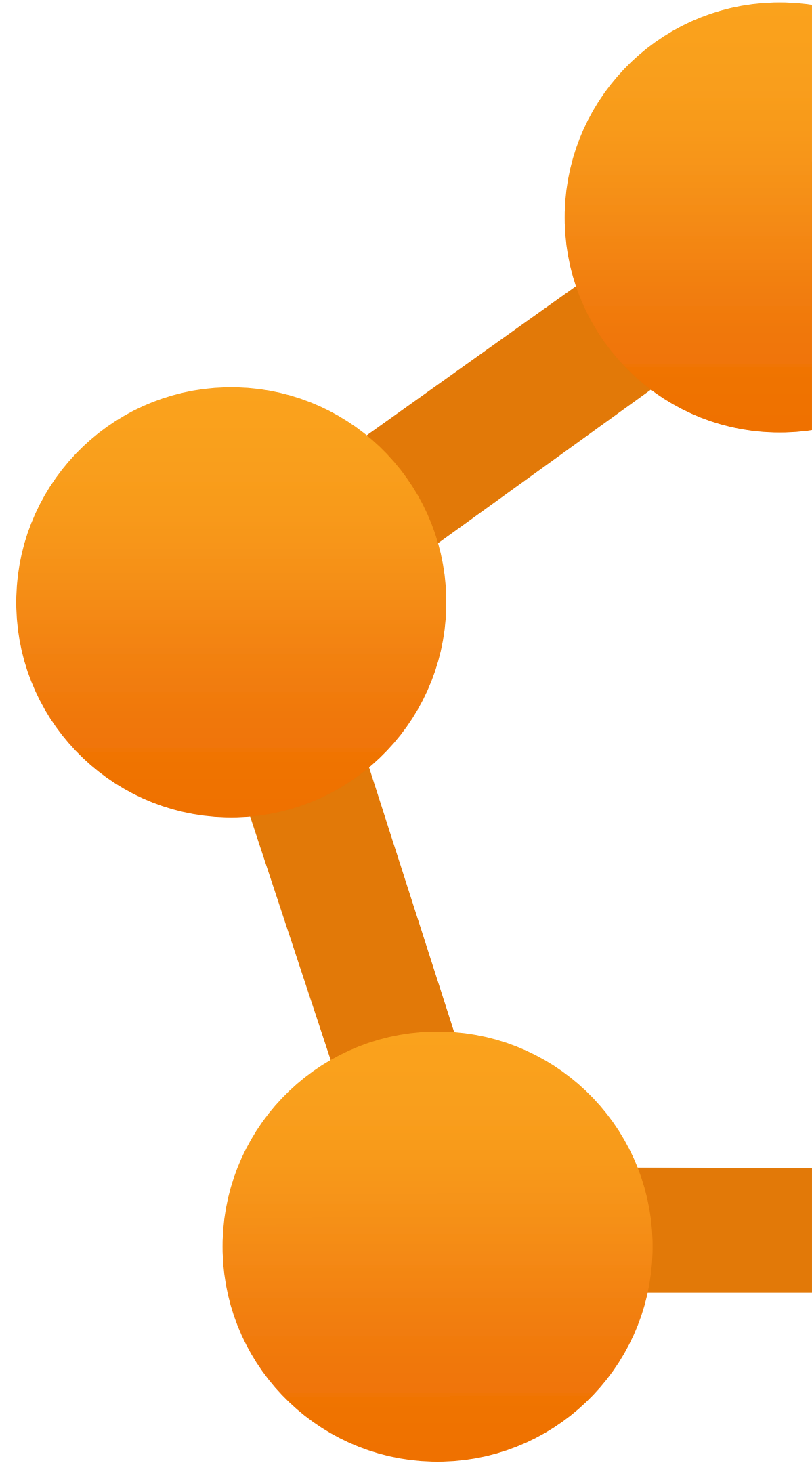
Cloud Formations - Knowledge Transfer & Training

# Module 6

Execution Parallelism ✔

Any questions?

Cloud Formations

# Agenda: Data Integration Pipelines

## Fundamentals to Level 300

**Cloud Formations - Knowledge Transfer & Training**

- **Module 1:** Pipeline Fundamentals
  - The History of Azure Orchestration
  - Synapse Analytics vs Data Factory vs Microsoft Fabric
  - Integration Components
  - Common Activities
  - Execution Dependencies
- **Module 2:** Integration Runtime Design Patterns
  - Compute Types
    - Azure
    - Hosted
    - SSIS
  - Patterns & Configuration
- **Module 3:** Data Transformation
  - Data Flows
  - Power Query Injection
  - Spark Configuration
  - Use Cases
- **Module 4:** Dynamic Pipelines
  - Expressions & Interpolation
  - Simple Metadata Driven Execution
  - Dynamic Content Chains
  - Reference Names

<< BREAK

- **Module 5:** Pipeline Extensibility
  - Azure Batch Service
  - Pipeline Custom Activities
  - Azure Management API
  - Azure Functions
- Labs
  - Create Azure resources
  - Build a copy pipeline
  - Create a reusable pipeline
  - Author a data flow
  - Monitor factory activity
  - Explore Synapse pipelines
  - Explore Fabric pipelines
  - Mini-project

<< LUNCH

- **Module 6:** Execution Parallelism
  - Control Flow Scale Out
  - Concurrency Limitations
  - Internal vs External Activities
  - Orchestration Framework
- **Module 7:** VNet Integration
  - Private Endpoints
  - Managed VNet's
  - Firewall Bypass

**Development**

- **Module 8:** Security
  - Service Principals
  - Managed Identities
  - Azure Key Vault Integration
  - Customer Managed Keys
  - Pipeline Access & Permissions
- **Module 9:** Monitoring & Alerting
  - Studio Monitoring
  - Log Analytics & Kusto Queries
  - Operational Dashboards
  - Advanced Alerting

<< BREAK

- **Module 10:** Solution Testing
  - Development Time Validation
  - Test Coverage
  - NUnit Tests
- **Module 11:** CI/CD
  - Source Control vs Developer UI
  - Basic ARM Template Deployments
  - Advanced Deployment Patterns
- **Module 12:** Final Thoughts
  - Costs & Conclusions
  - Best Practices

**Production**

© 2024 Cloud Formations Ltd