

# Building a Metadata driven ELT Framework in Azure Synapse Analytics



Erwin de Kreuk



@erwindekreuk

<https://erwindekreuk.com>



**Microsoft®**  
Most Valuable  
Professional



Erwin de Kreuk

Principal Consultant | Lead Data and AI  
| Data Platform MVP | Public Speaker ...



# Erwin de Kreuk

Principal Consultant – Lead Data & AI  
InSpark

Twitter icon

@erwindekreuk

LinkedIn icon

[linkedin.com/in/erwindekreuk](https://www.linkedin.com/in/erwindekreuk)

Globe icon

[erwindekreuk.com](http://erwindekreuk.com) [github.com/edkreuk](https://github.com/edkreuk)





**REBTECH**



**dbWatch**



**redgate**

**devart**



# We Are InSpark

We help organizations  
accelerating their digital  
transformation with impactful  
Microsoft solutions & expertise

# Azure Synapse Analytics



Metadata driven ELT Framework

# Objectives of today

- Why
- Out-of-the-box Framework
- Custom-made Framework
- Recap



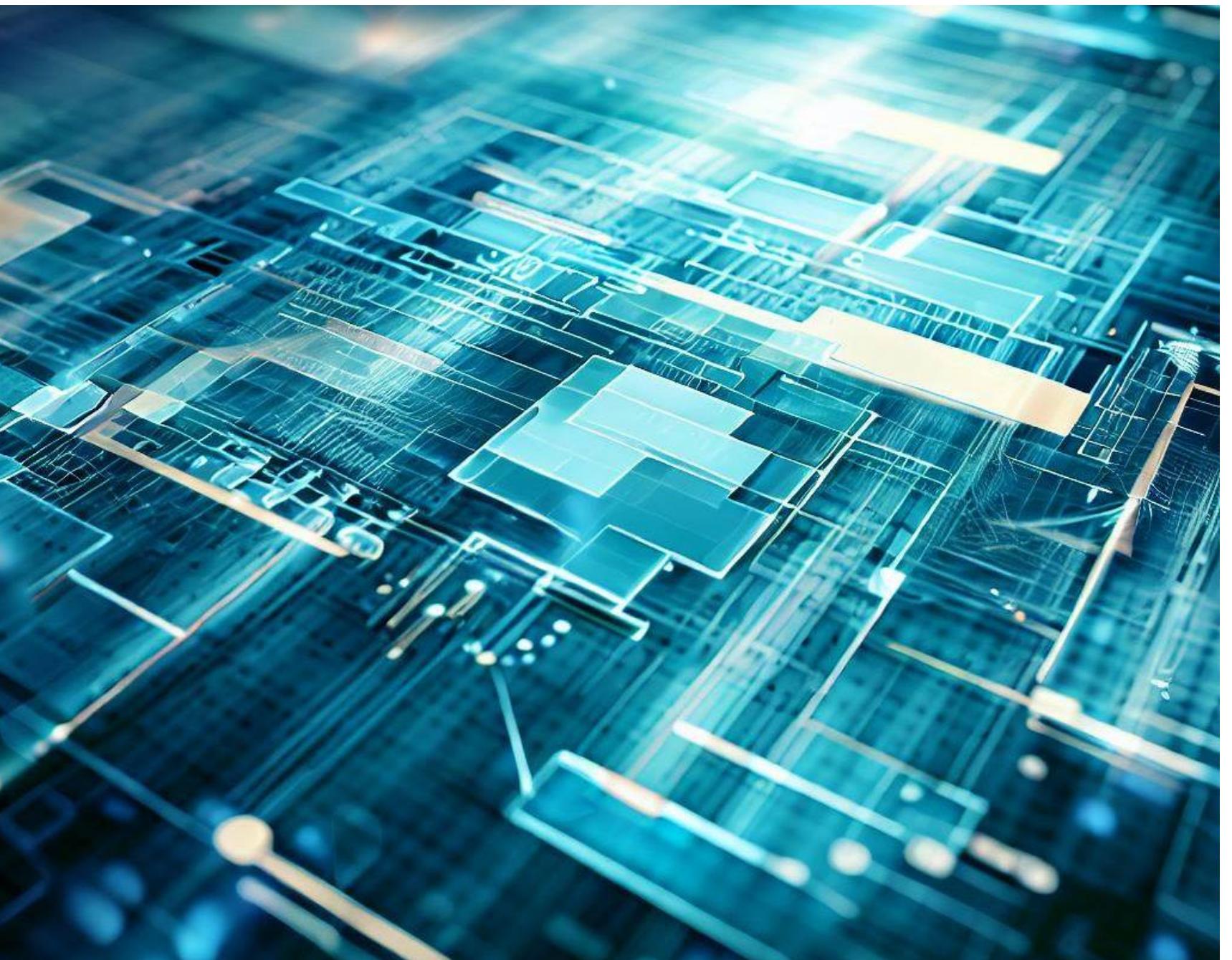
# Why

- ★ **Centralized** metadata management for better governance and maintenance.
- ★ **Flexibility** and **agility** by separating metadata from code.
- ★ **Scalability** and performance leveraging Synapse's distributed processing capabilities.
- ★ **Reusability** and standardization of data integration components.
- ★ **Simplified** maintenance through metadata modifications instead of manual code changes.
- ★ **Collaboration** and documentation for alignment between business and technical teams.
- ★ **Extensibility** and customization using Synapse's integration capabilities.

# Out-of-the-Box Framework

## Out-of-the-Box

- Ready-to-use.
- Rapid implementation.
- Limited customization.
- Lower development effort.
- Lower upfront costs.
- Ongoing support and updates.



# Custom-Made Framework

## Custom-Made

- Tailored to specific needs.
- Full control over design and features.
- Higher development effort.
- Flexibility and extensibility.
- Higher upfront costs.



# Out-of-the-Box Framework

- ★ Available by default in Azure Synapse Analytics and Azure Data Factory
- ★ Azure SQL Database as requirement (Linked Service)
- ★ Almost next, next, finish



## Metadata-driven copy task

You will get parameterized pipelines which can read metadata from an external store to load data at a large scale.

Metadata-driven Framework

## Out-of-the-Box Framework

- ★ Connect to table data store
- ★ Define Schema and Table name
- ★ Define Schedule
- ★ Connect to Source
- ★ Select Tables
- ★ Define full or Incremental
- ★ Define Destination
- ★ Create tables in datastore

Please run the following SQL script in your SQL server to create a control table. Then view your pipeline to execute a debug run.

[Download SQL script](#)

Generated SQL script for control table

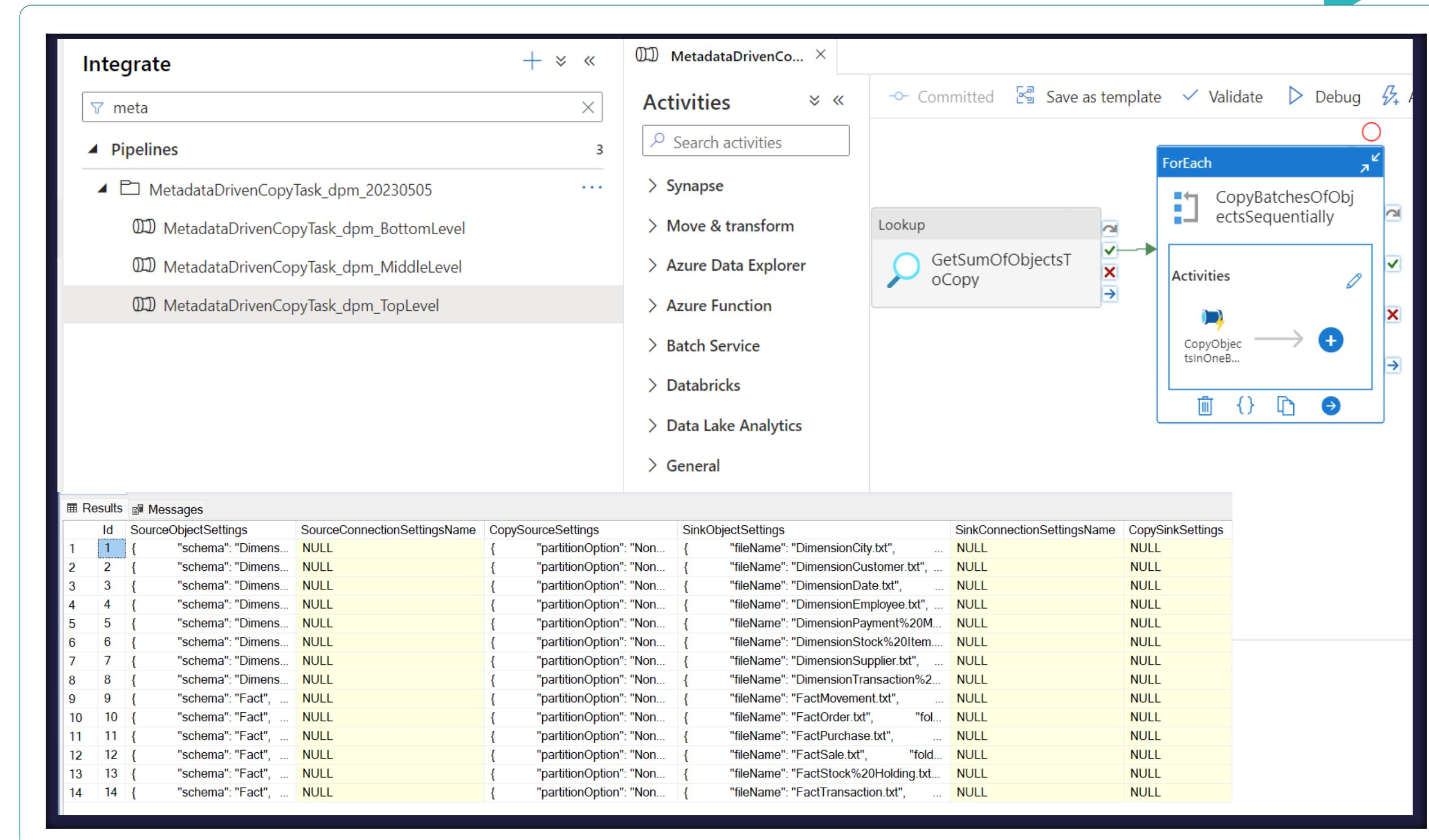
[Copy to clipboard](#)

```
***** Object: Table [MainControlTable_k4j] *****/
CREATE TABLE [MainControlTable_k4j](
    [Id] [int] IDENTITY(1,1) NOT NULL PRIMARY KEY,
    [SourceObjectSettings] [nvarchar](max) NULL,
    [SourceConnectionSettingsName] [varchar](max) NULL,
    [CopySourceSettings] [nvarchar](max) NULL,
    [SinkObjectSettings] [nvarchar](max) NULL,
    [SinkConnectionSettingsName] [varchar](max) NULL,
```

## Metadata-driven Framework

# Out-of-the-Box Framework

- ★ Execute pipeline
- ★ Top
  - ★ This pipeline will calculate the total number of objects (tables etc.) required to be copied in this run
- ★ Middle
  - ★ This pipeline will copy one batch of objects.
- ★ Bottom
  - ★ This pipeline will copy objects from one group



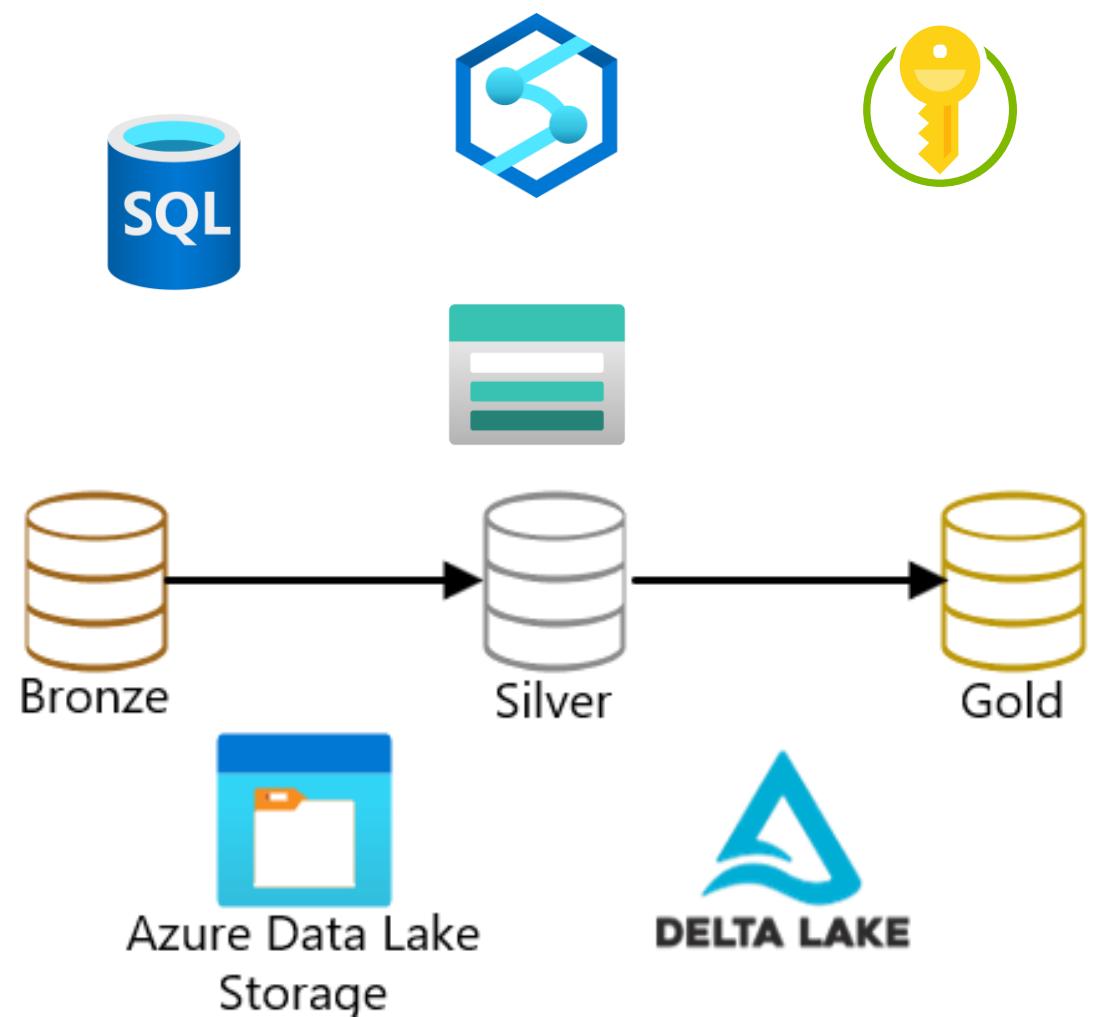
Innovate to accelerate

# DEMO

Metadata-driven Framework

# Custom-Made Framework

- ★ Based on parameters
- ★ Meta data => Azure SQL Database
- ★ Azure Synapse Analytics and Azure Data Factory
- ★ Azure Key Vault
- ★ Azure Data lake
- ★ Based on the Medaillon Architecture



# Who is already using Parameters?



# Custom-Made Framework

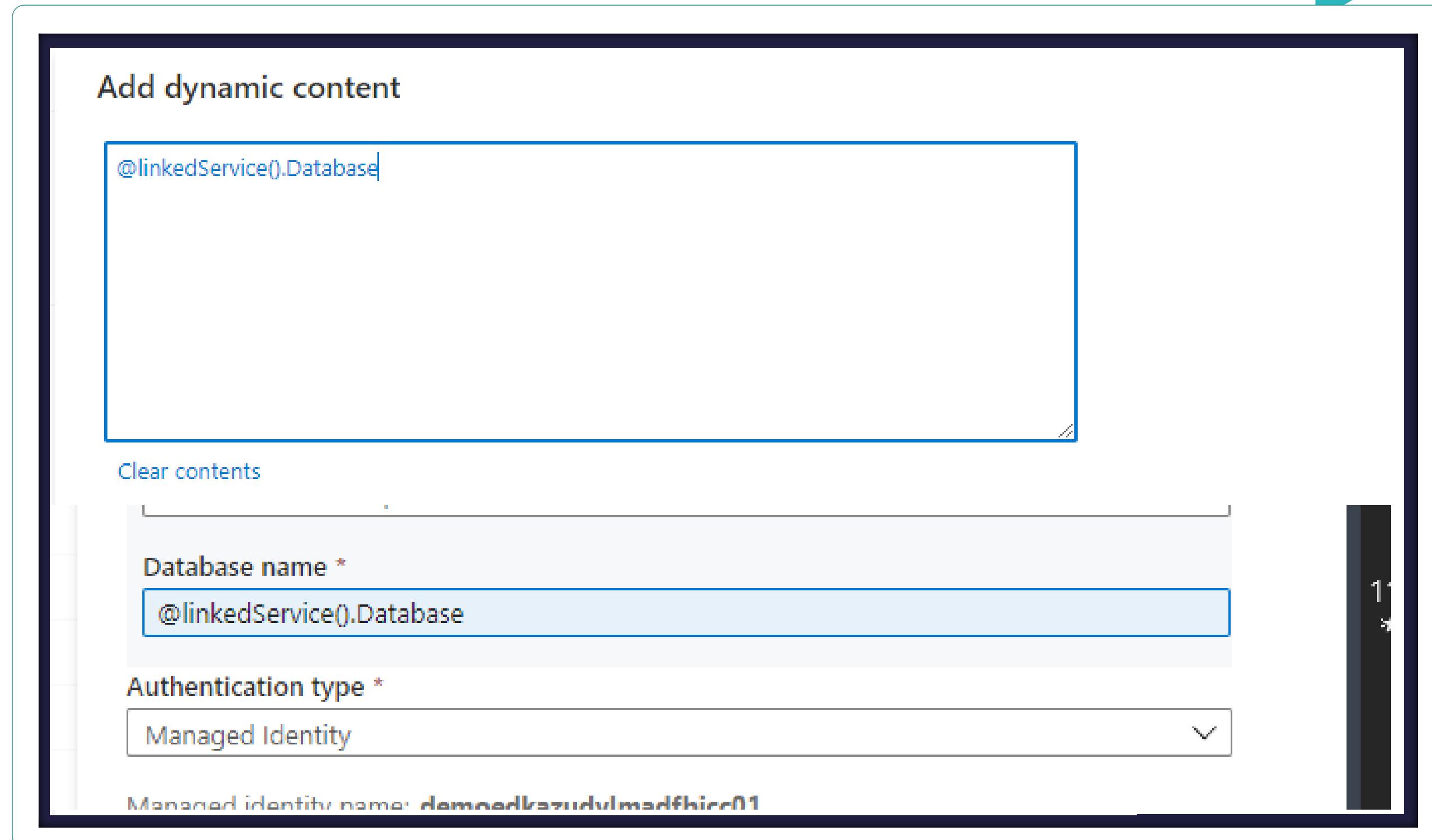
- ★ Linked Services Parameter
- ★ Dataset Parameter
- ★ Dataflow Parameter
- ★ Notebook Parameter
- ★ Pipeline Parameter

IMPLEMENTING  
**DEFAULT PARAMETERS**  
THAT DEPEND ON  
OTHER PARAMETERS

## Custom-Made Framework

# Linked Services Parameters

- ★ Connect to different Databases on same Server
- ★ Connect to different Logical Servers based on the same Linked Service



# Custom-Made Framework

## Dataset Parameters

- ★ Create 1 dataset for all your Linked Services activities

Add dynamic content

@dataset().FilePath

Parquet  
DS\_ADLS\_RAW

Connection Schema Parameters

Linked service \* LS\_A\_LS\_DLS2 Test connection Edit New Learn more

File path \* raw / @dataset().FilePath / @dataset().FileName Browse

Compression type snappy

Parameters

+ New | Delete

Name	Type	Default value
FilePath	String	Value
Filename	String	Value

Clear contents

## Custom-Made Framework

# Dataflow Parameters

- ★ Create a dynamically Dataflow
- ★ Easy to re-use

**General**   **Settings**   **Parameters**   **User properties**

Data flow parameters ⓘ

Name	Value	Type	Expression ⓘ
TargetFilePath	@pipeline().parameters.FilePath	string	<input type="text"/>
TargetFileName	@pipeline().parameters.FileName	string	<input type="text"/>
ProcessType	@pipeline().parameters.ProcessType	string	<input type="text"/>
DataStoreExcludeColum...	@pipeline().parameters.DataStoreExcl...	string	<input type="text"/>

The screenshot shows the Azure Data Factory Data Flow designer interface. At the top, there's a visual representation of a data flow pipeline with several stages: 'DataStoreExcludeColu...', 'IsCurrent', 'RemoveAuditColumnsA...', and 'MarkAsInsert'. Below this, there's a 'Parameters' tab where a parameter named 'TargetFilePath' is defined with a string type and an expression field containing the value 'DeltaLake/' + \$TargetFilePath + '/' + \$TargetFileName + '/Original'. The 'Parameters' table also includes columns for 'Name', 'Type', and 'Default value'. The 'Default value' column for 'TargetFilePath' contains the expression 'Enter expression... ANY'. The 'Expression' column for 'TargetFilePath' contains the expression '@pipeline().parameters.FilePath'. Arrows point from the 'Value' and 'Expression' columns of the 'Parameters' table to the corresponding fields in the 'Default value' and 'Expression' columns of the 'Parameters' table.

# Notebook Parameters

- ★ Pass Parameters from Synapse/ADF Pipeline to Databricks
  - ★ Using Widgets
- ★ Pass Parameters from Synapse Pipeline to Synapse Notebook
  - ★ Toggle parameter cell

The screenshot shows a Databricks notebook interface. On the left, there's a sidebar with a 'Parameters' section containing a list of variables and their values. On the right, there's a code editor window with some Python code.

Code in the editor:

```
1 # Creating widgets for leveraging parameters, and printing the parameters
2
3 datalake_connection = 'labseuwvd1mdloxgn01.dfs.core.windows.net'
4 target_file = ''
5 target_PKColumns = 'ProductId'
6 source_file = 'AW/SalesLT_Product/2021/06/09/AW_SalesLT_Product_202106090916.parquet'
7 target_file_path = ''
8 HashExcludeColumns = ''
9 delta_db = 'delta'
10 IsIncremental = ''
11 ReturnOutput = True
```

Below the code, there's a green checkmark icon.

<input type="checkbox"/> datalake_connection	@activity('Get EnvironmentSettings').o...
<input type="checkbox"/> target_file	@{pipeline().parameters.TargetFileName}
<input type="checkbox"/> target_PKColumns	@{pipeline().parameters.PKColumns}
<input type="checkbox"/> source_file	@{pipeline().parameters.FilePath}/@{pi...
<input type="checkbox"/> target_file_path	@pipeline().parameters.TargetFilePath
<input type="checkbox"/> HashExcludeColumns	@pipeline().parameters.HashExcludeC...
<input type="checkbox"/> delta_db	@{pipeline().parameters.DeltaDB}
<input type="checkbox"/> IsIncremental	@{pipeline().parameters.IsIncremental}
<input type="checkbox"/> ReturnOutput	@{pipeline().parameters.ReturnOutput}

## Custom-Made Framework

# Pipeline Parameters

- ★ Can be used across all your Pipelines
- ★ Can easily control dataset, linked services and dataflows etc

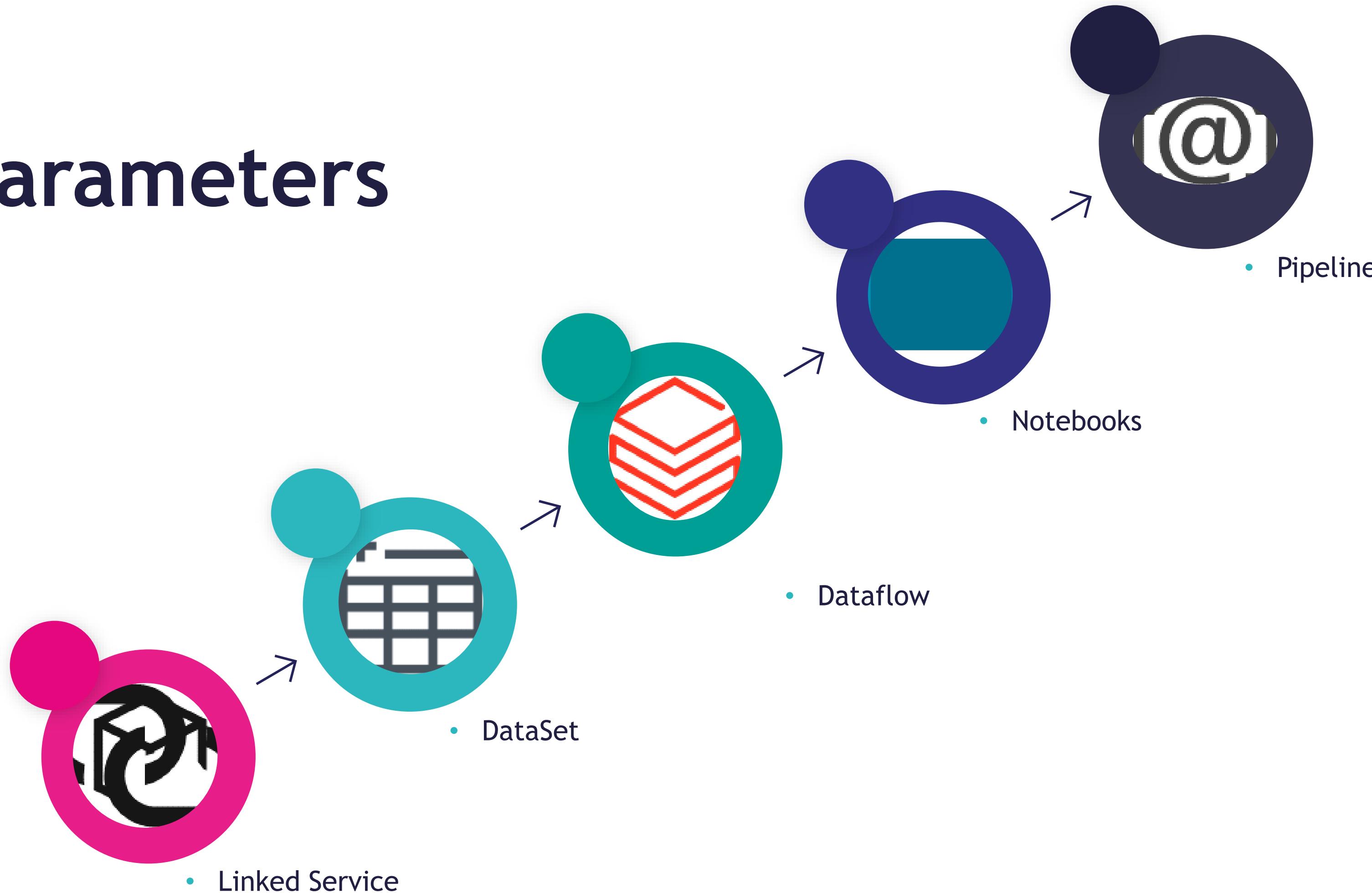
Add dynamic content

```
@pipeline().parameters.FilePath
```

The screenshot shows the 'Sink' tab of the 'Dataset properties' configuration for a 'DS\_ADLS\_RAW' dataset. The 'NAME' column lists 'FilePath' and 'FileName'. The 'VALUE' column for both contains the placeholder 'Value'. The 'TYPE' column indicates both are of type 'string'. Below this, there is a table for defining parameters:

Name	Type	Default value
FilePath	String	Value
FileName	String	Value

# Parameters



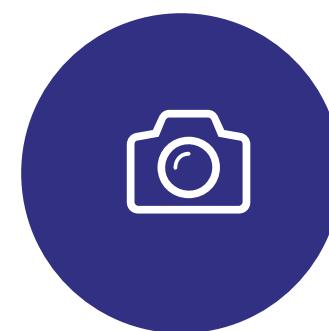
Innovate to accelerate

# DEMO

## Custom-Made Framework



Can we build Synapse Pipelines dynamically?



Can we load the active(current) or historical records to a DataStore?



Can we extract data from my sources based on MetaData?



Can we build history from extracted data based on MetaData?



Can we log the execution of the Pipelines?

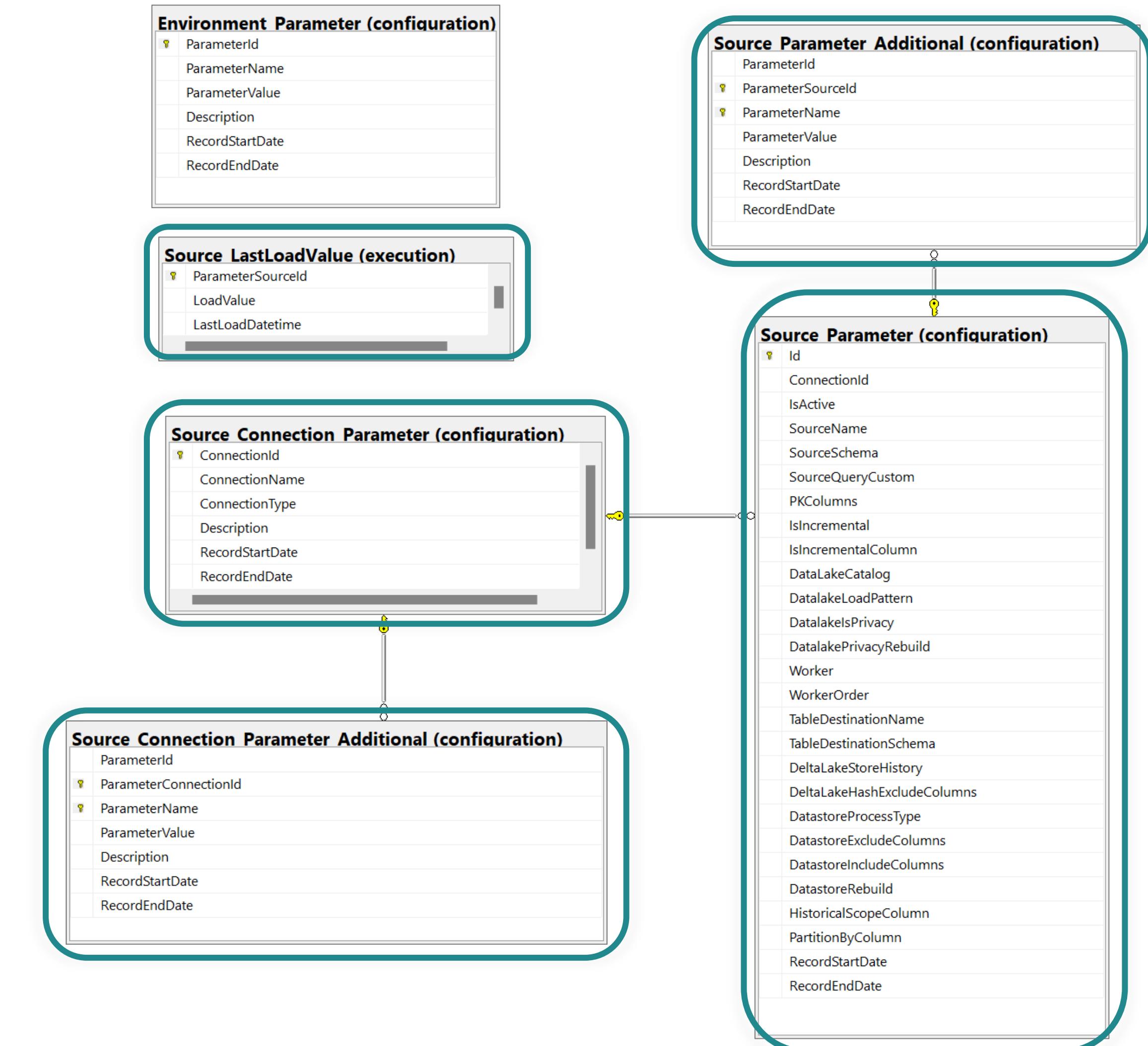


YES WE  
CAN.

Custom-Made Framework

# Configuration

- ★ Source Parameters
  - ★ Source Parameters Additional
- ★ Source Connection Parameters
  - ★ Source Connection Parameters additional
- ★ Source Last LoadValue
- ★ Versioning enabled (track changes)



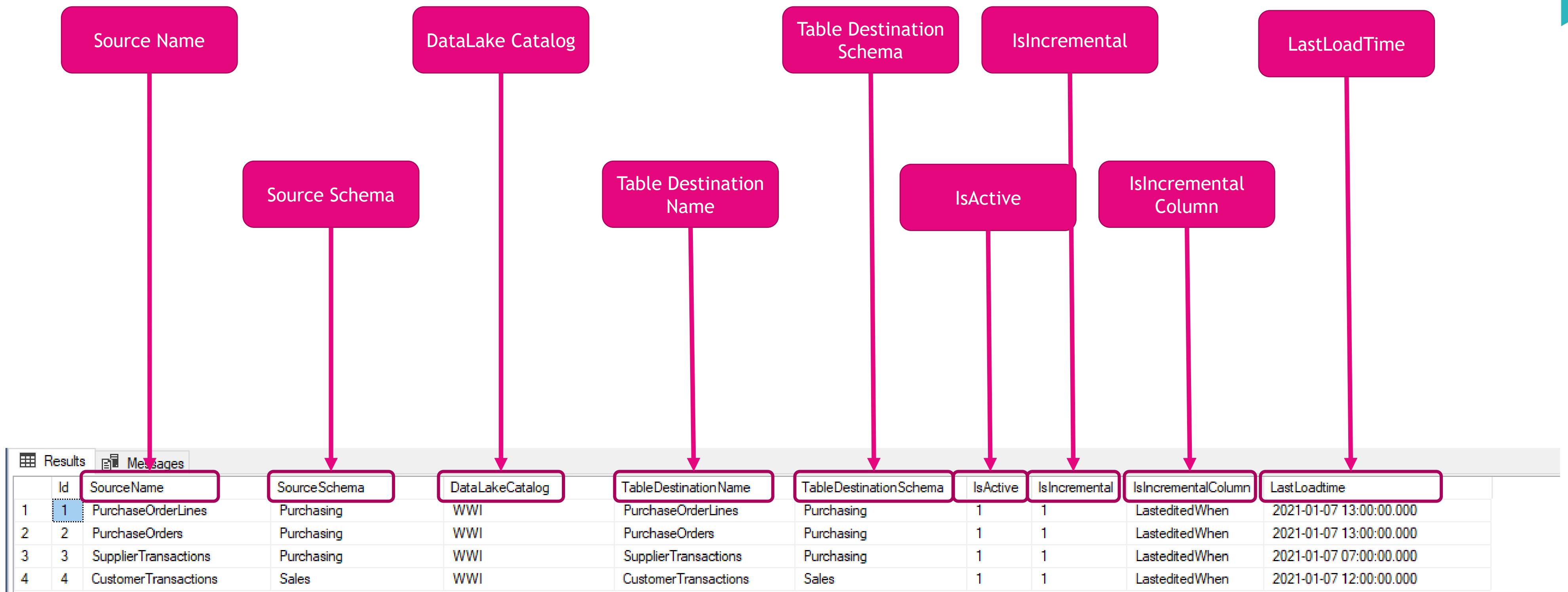
## Custom-Made Framework

# Environment

- ★ Overarching settings that can be constantly reused.
- ★ Global settings

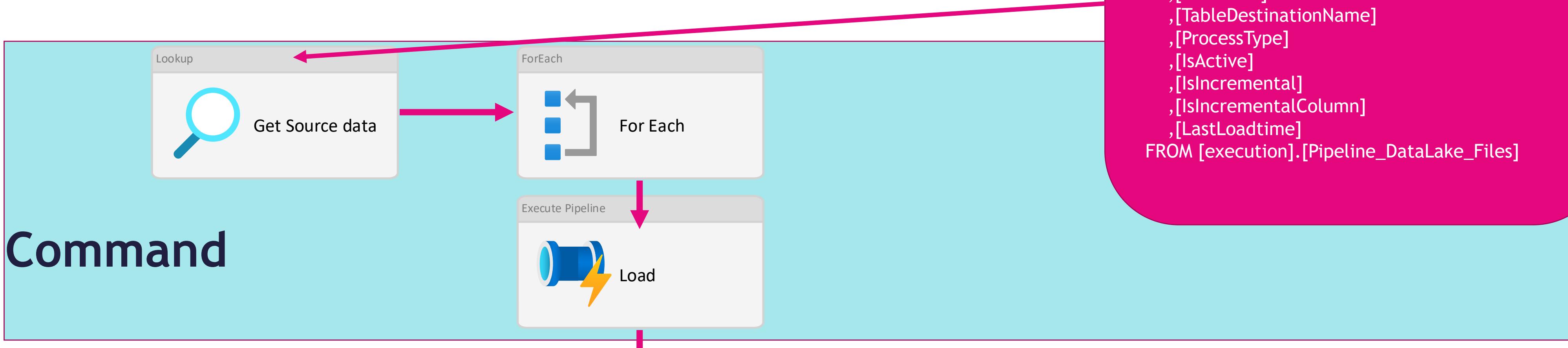
Environment Parameter (configuration)					
	ParameterId	ParameterName	ParameterValue	Description	RecordStartDate
	ParameterId	ParameterName	ParameterValue	Description	RecordEndDate
1	1	Environment	DVLM	Enviroment short DVLM-TEST-ACPT-PROD	
2	2	EnvironmentLong	development	Enviroment long develop-test-accpetance-production	
3	3	AzureDataBricksClusterId	XXXX	ClusterId for Databricks Environment	
4	4	ResourceGroupName	XXXX	ResourceGroupName of your Resources	
5	5	KeyVaultConnection	https://ededeudvilmvaultxgn01.vault.azure.net/	Connection to keyvault of your current environment	
6	6	DataLakeConnection	https://ededeudvldmloxgn01.dfs.core.windows.net/	DataLake Connection	
7	7	AAS_ServerName	XXXX	Name of your Azure Analysis Service Server	
8	8	SubscriptionId	xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxx	SubscriptionId	
9	9	TenantId	xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxx	TenantId	
10	10	KeyVaultConfigDB	ASQL-METADATAFRAMEWORK	Keyvault secret with Connection Details voor Config...	
11	11	KeyVaultDatastore	ASQL-OXGN01-DATASTORE	Keyvault secret with Connection Details voor DataS...	
12	12	IsDataStoreCSVEnabled	0	Do you want to create a CSV file from your DataSt...	
13	13	IsDataStoreParquetEnabled	0	Do you want to create a Parquet file from your Data...	
14	14	NumberOfDIUDataLake	4	DIU capacity for processing Source to DataLake av...	
15	15	NumberOfDIUDataStore	0	DIU capacity for processing Deltalake to DataStore...	
16	16	DataStoreSchema	Stg	Name of the schema which is used to Store your D...	
17	17	IsSourceLoggingEnabled	1	Feature functionality to log ADF/Synapse	
18	18	IsDataStoreLoggingEnabled	1	Feature functionality to log ADF/Synapse	
19	19	IsDataStoreHistoricalEnabled	0	Are you building a Histroical Datastore model	
20	20	IsDataVaultLoadEnabled	0	Are you building a DataVault model	
21	21	IsHistoricalDatastoreLoadEnabled	0	Create Historical layer in SQL DataStore, loading D...	
22	22	MaxRowsPerFile	1000000	Number of rows to split files	
23	23	SynapseDedicatedSQLPool	DWH	Name of the Dedicated SQL Pool in Synapse	
24	24	UseSynapseServerlessDataStore	1	Use Serverless Engine to read from DeltaLake to D...	
25	25	DeltaLakeModule	2	0=DataBricks 1=DataFlow 2=SynapseSpark	
26	26	DatastoreModule	1	1=Use datastore 0=No datastore	
27	27	LakehouseModule	1	1=Use lakehouse 0=No lakehouse	
28	28	TeamsWebhook		Webhook to sent general pipeline messages to	

## Custom-Made Framework



Custom-Made Framework

# Pipelines

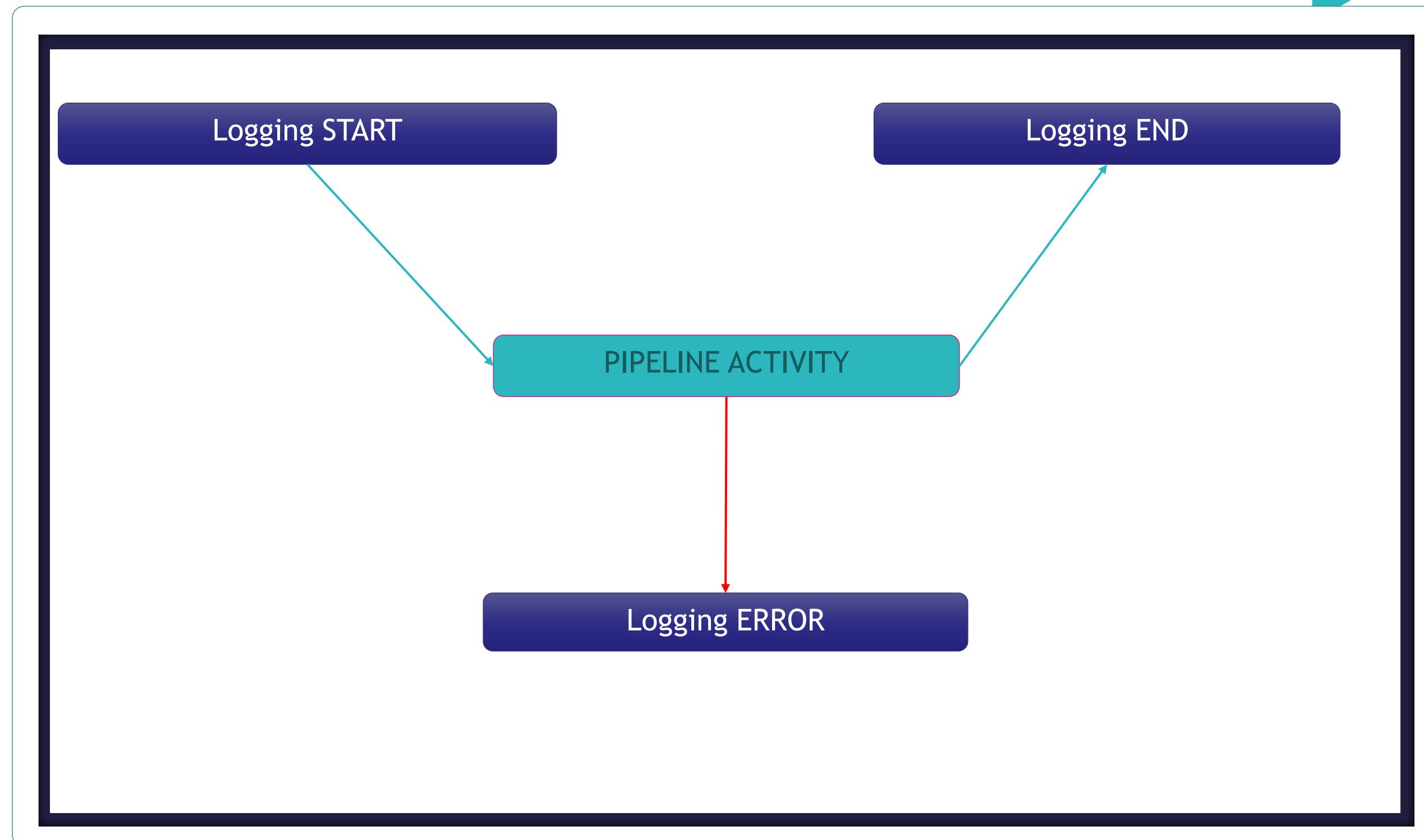


## Execute

Custom-Made Framework

# Logging

- ★ Log Start and End Time of records
- ★ Log Extracted Records
- ★ Log Execution Failure
- ★ Create Pipeline\_ExecutionLog table



Custom-Made Framework

# Logging

- ★ Log Start and End Time of records
- ★ Log Extracted Records
- ★ Log Execution Failure
- ★ Create Pipeline\_ExecutionLog table

BEGIN

Insert new Record  
Insert Metadata  
Insert Start time

END

End Time  
Status(1)  
Row Counts  
Pipeline Details

ERROR

End Time  
Status(2)  
Failure Message

Pipeline\_ExecutionLog

	Results	Messages											
LogId	StartTime	EndTime	Status	RowsRead	RowsNew	RowsUnchanged	RowsDeleted	RowsSkipped	Reference	NumberOfParallelCopies	IntegrationRuntime	DataConsistencyVerification	ReportLineageToCatalog
1	5	2021-02-08 19:55:35.517	2021-02-08 19:55:57.687	1	84	84	0	0	DataMovement	1	DefaultIntegrationRuntime (West Europe)	Verified	N/A
2	4	2021-02-08 19:55:04.670	2021-02-08 19:55:28.327	1	2	2	0	0	DataMovement	1	DefaultIntegrationRuntime (West Europe)	Verified	N/A
3	3	2021-02-08 19:54:35.920	2021-02-08 19:54:58.937	1	2	2	0	0	DataMovement	1	DefaultIntegrationRuntime (West Europe)	Verified	N/A
4	2	2021-02-08 19:54:09.117	2021-02-08 19:54:29.653	1	9	9	0	0	DataMovement	1	DefaultIntegrationRuntime (West Europe)	Verified	N/A
5	1	2021-02-08 19:53:57.763	NULL	0	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL	NULL

LogId	StartTime	En...	Status	FailureMessage
1	15	2021-02-0...	20...	1
2	14	2021-02-0...	20...	1
3	13	2021-02-0...	20...	2
4	12	2021-02-0...	20...	1
5	11	2021-02-0...	20...	2

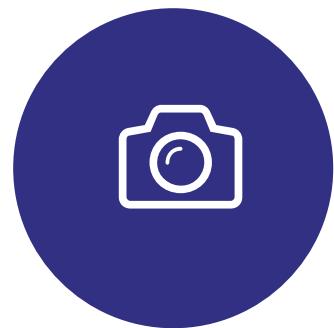
Innovate to accelerate

# DEMO

## Custom-Made Framework



Can we build Synapse Pipelines dynamically?



Can we load the active(current) or historical records to a Lakehouse?



Can we extract data from my sources based on MetaData?



Can we build history from extracted data based on MetaData?

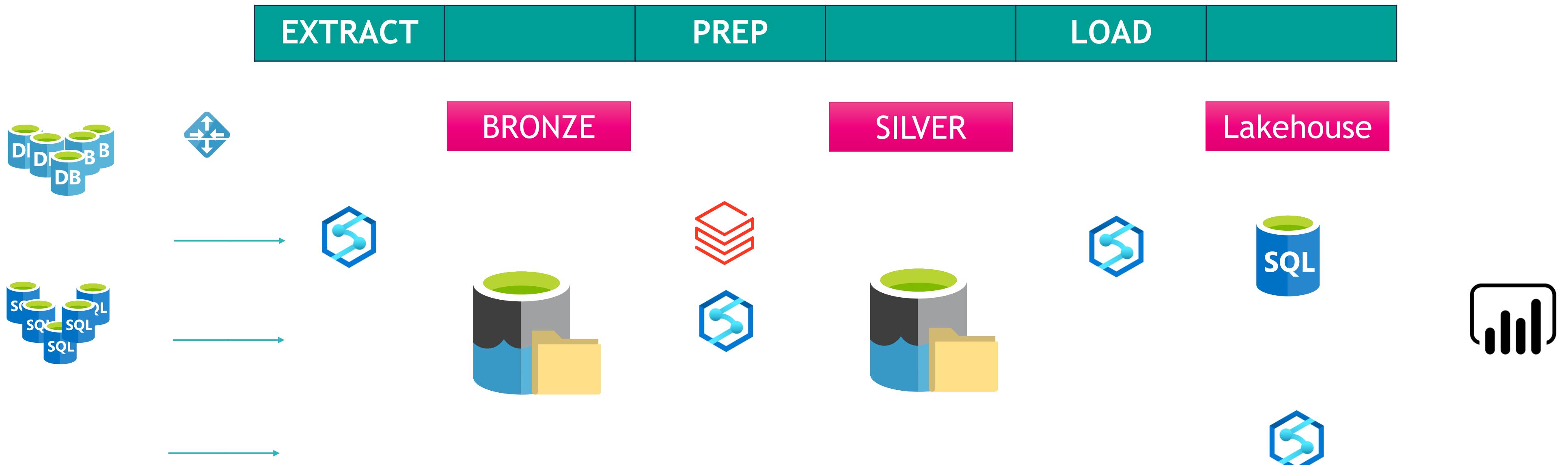


Can we log the execution of the Pipelines?



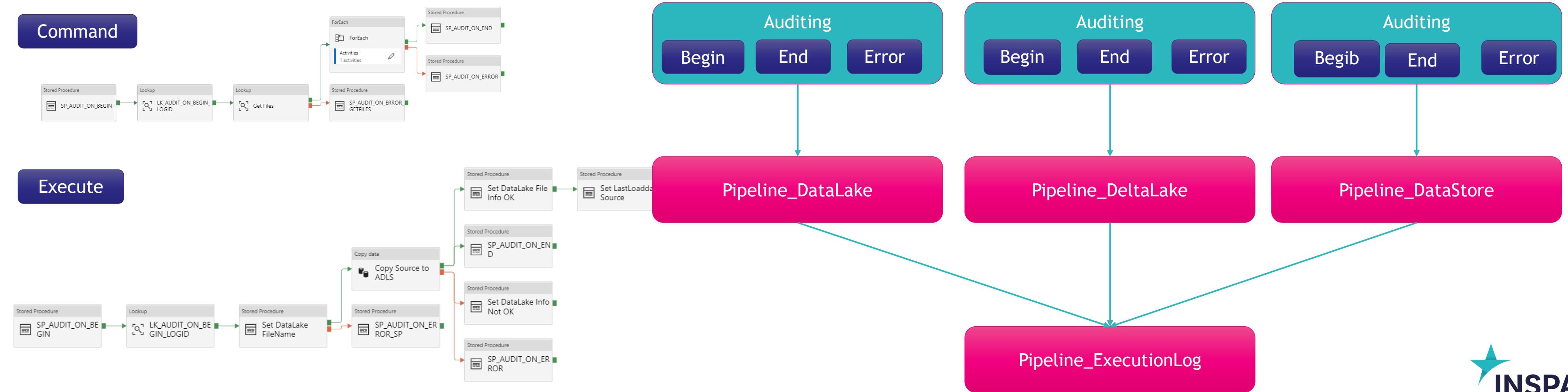
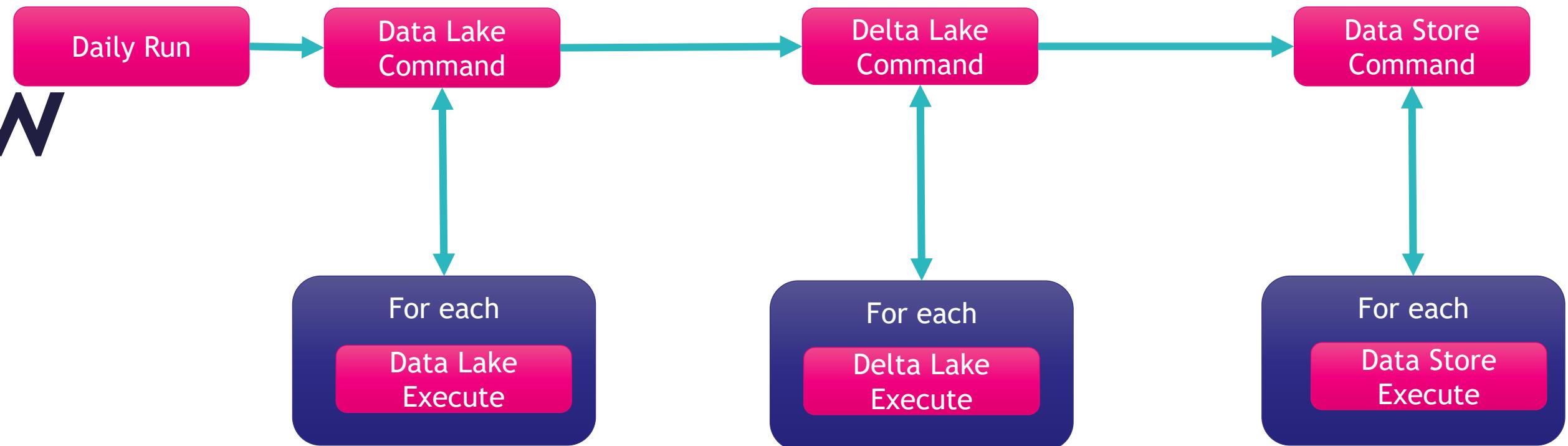
Custom-Made Framework

# Loading Pattern



# Custom-Made Framework

# Process Flow



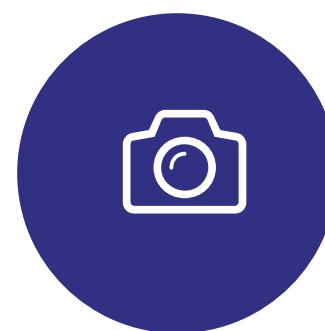
Innovate to accelerate

# DEMO

# Recap



Can we build Synapse Pipelines dynamically?



Can we load the active(current) or historical records to a Lakehouse?



Can we extract data from my sources based on MetaData?



Can we build history from extracted data based on MetaData?



Can we log the execution of the Pipelines?



# Questions?





Erwin de Kreuk

Principal Consultant | Lead Data and AI  
| Data Platform MVP | Public Speaker ...



# Erwin de Kreuk

Principal Consultant – Lead Data & AI  
InSpark

Twitter icon

@erwindekreuk

LinkedIn icon

[linkedin.com/in/erwindekreuk](https://www.linkedin.com/in/erwindekreuk)

Globe icon

[erwindekreuk.com](http://erwindekreuk.com) [github.com/edkreuk](https://github.com/edkreuk)

