



Drive Efficiency and Reliability Using Metadata-Driven Frameworks

09-10-2024

Erwin de Kreuk

Marco Hansma

We Are InSpark

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

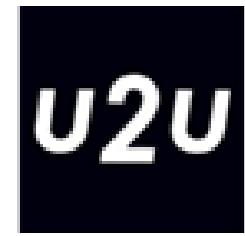
Thank you, partners



PLAINSIGHT



delaware



EpicData.

LACO/

tillit
data shapers

lytix

bmatix
Act informed

randstad
digital

KOHERA

MONIN
Database Managed Services

inetum.
realdolmen
Positive digital flow

TriFinance
BEYOND ADVISORY

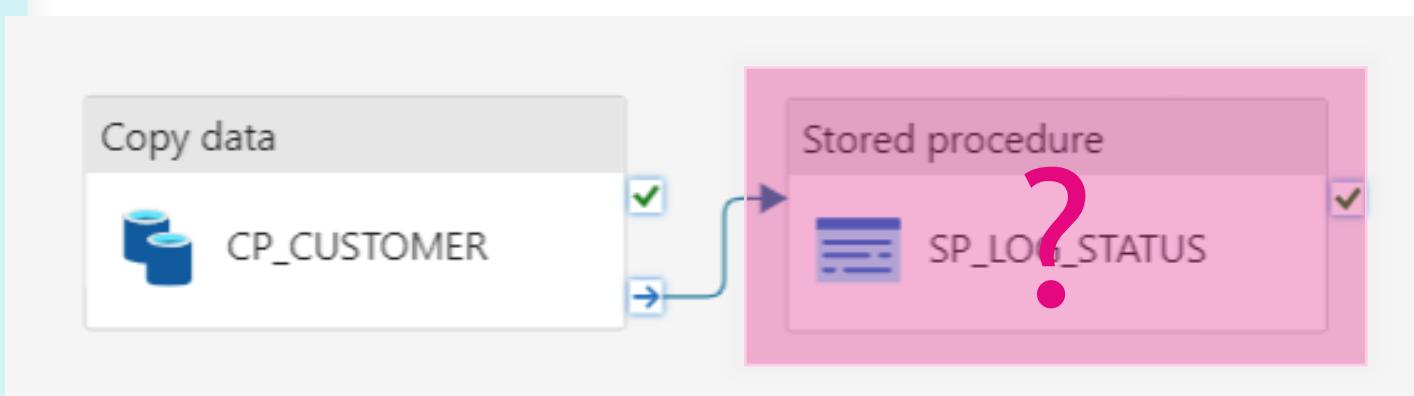
P

The screenshot shows the Power BI FabricDemo workspace. On the left, there's a navigation bar with icons for Home, Create, Browse, OneLake, Apps, Metrics, Monitor, Learn, Real-Time hub, Workspaces, and FabricDemo. The main area displays a list of data pipelines under the 'FabricDemo > Chaos' folder. The table has columns for Name and Type. The pipelines listed are:

Name	Type
PL_EXTRACT_AW_CUSTOMER	Data pipeline
PL_EXTRACT_AW_PRODUCT_CATEGORY	Data pipeline
PL_EXTRACT_AW_PRODUCT_DESCRIPTION	Data pipeline
PL_EXTRACT_AW_PRODUCT_DESCRIPTION_LINEAGE	Data pipeline
PL_EXTRACT_AW_PRODUCT_MODEL	Data pipeline
PL_EXTRACT_AW_PRODUCT_MODEL_PRODUCT_DESCRIPTION	Data pipeline
PL_EXTRACT_AW_PRODUCTS	Data pipeline
PL_EXTRACT_AW_SALES_ORDER_DETAIL	Data pipeline
PL_EXTRACT_AW_SALES_ORDER_HEADER	Data pipeline
PL_EXTRACT_AW_SALES_ORDER_LINE	Data pipeline
PL_EXTRACT_WWI_ADDRESS	Data pipeline
PL_EXTRACT_WWI_CITIES	Data pipeline
PL_EXTRACT_WWI_COUNTRIES	Data pipeline
PL_EXTRACT_WWI_DELIVERY_METHODS	Data pipeline
PL_EXTRACT_WWI_PAYMENT_METHODS	Data pipeline
PL_EXTRACT_WWI_PEOPLE	Data pipeline
PL_EXTRACT_WWI_STATE_PROVINCES	Data pipeline
PL_EXTRACT_WWI_TRANSACTION_TYPES	Data pipeline

lement

Content of each pipeline

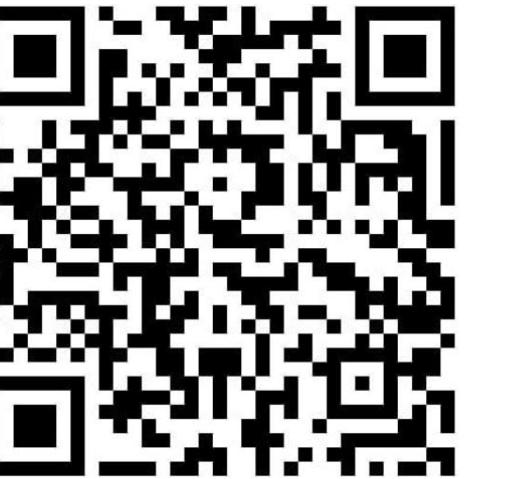


Erwin de Kreuk

Principal Consultant

Lead Data & AI InSpark

Let's connect



 @erwindekreuk

 linkedin.com/in/erwindekreuk

 erwindekreuk.com

 github.com/edkreuk

 <https://sessionize.com/erwin-de-kreuk/>



Marco Hansma

Architect InSpark

Let's connect

 linkedin.com/in/marcohansma



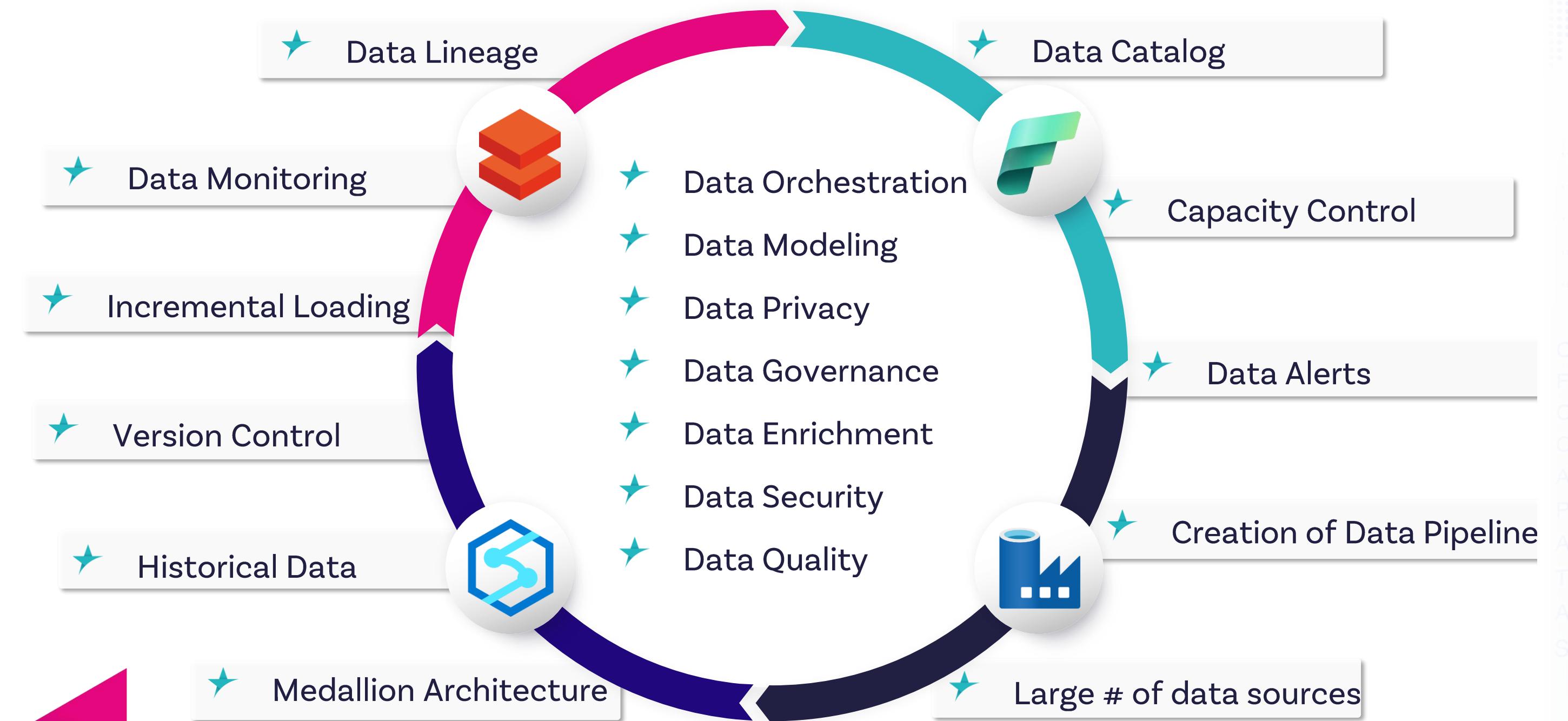
Objectives

Data platform challenges
Why
Medallion Architecture
Data Pipeline Parameters
Meta Data DrivenFramework
Data Pipeline Logging
Recap



Data platform Challenges

'From data source to data model' to report



**“Simplify,
I am the Chief Data Information
officer and don’t want to be the
Chief Integration Officer**

**Help me to simplify and automate my
Data Orchestration.”
Every CDO, Every organisation**



What would help

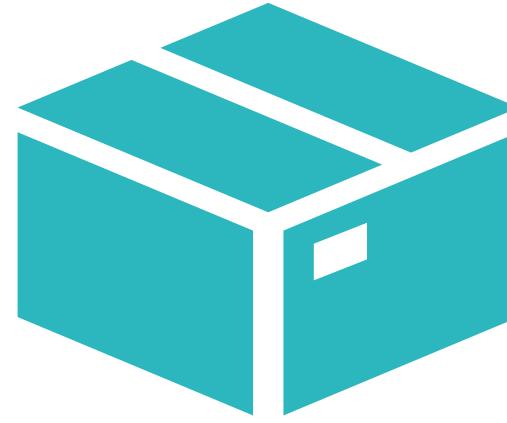
- ★ Simplicity in connecting data sources
- ★ First results within hours
- ★ Focus on business value instead of data integration
- ★ Meta Data Driven
 - Standardized data pipelines, Notebooks, orchestration and Way of Work
- ★ Overview of data process flows
- ★ Detailed logging information
- ★ Integration of solution like:
 - Data Privacy
 - Data Quality
 - Data Governance
- ★ From data “Spaghetti to Lasagna”
 - Building a uniform data architecture

*I am the Chief Data Information Officer
and don't want to be the Chief Integration Officer*

*Help me to simplify and automate my Data
Orchestration*



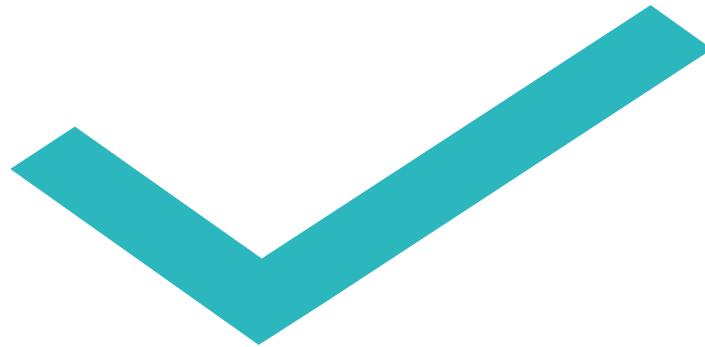
Why



Scalability



Automation



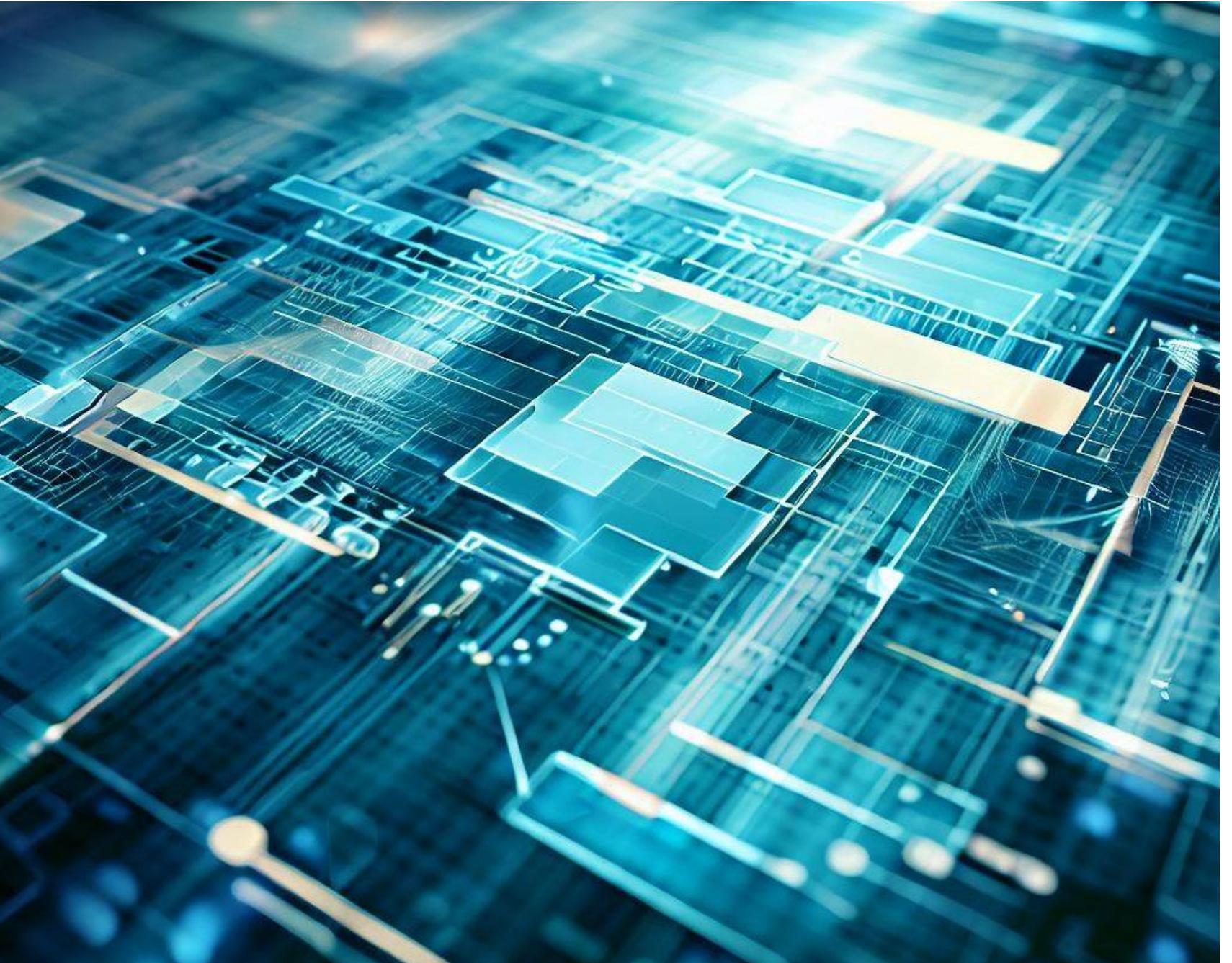
Traceability



Flexibility

Out-of-the-Box Framework

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



Custom-Made Framework

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



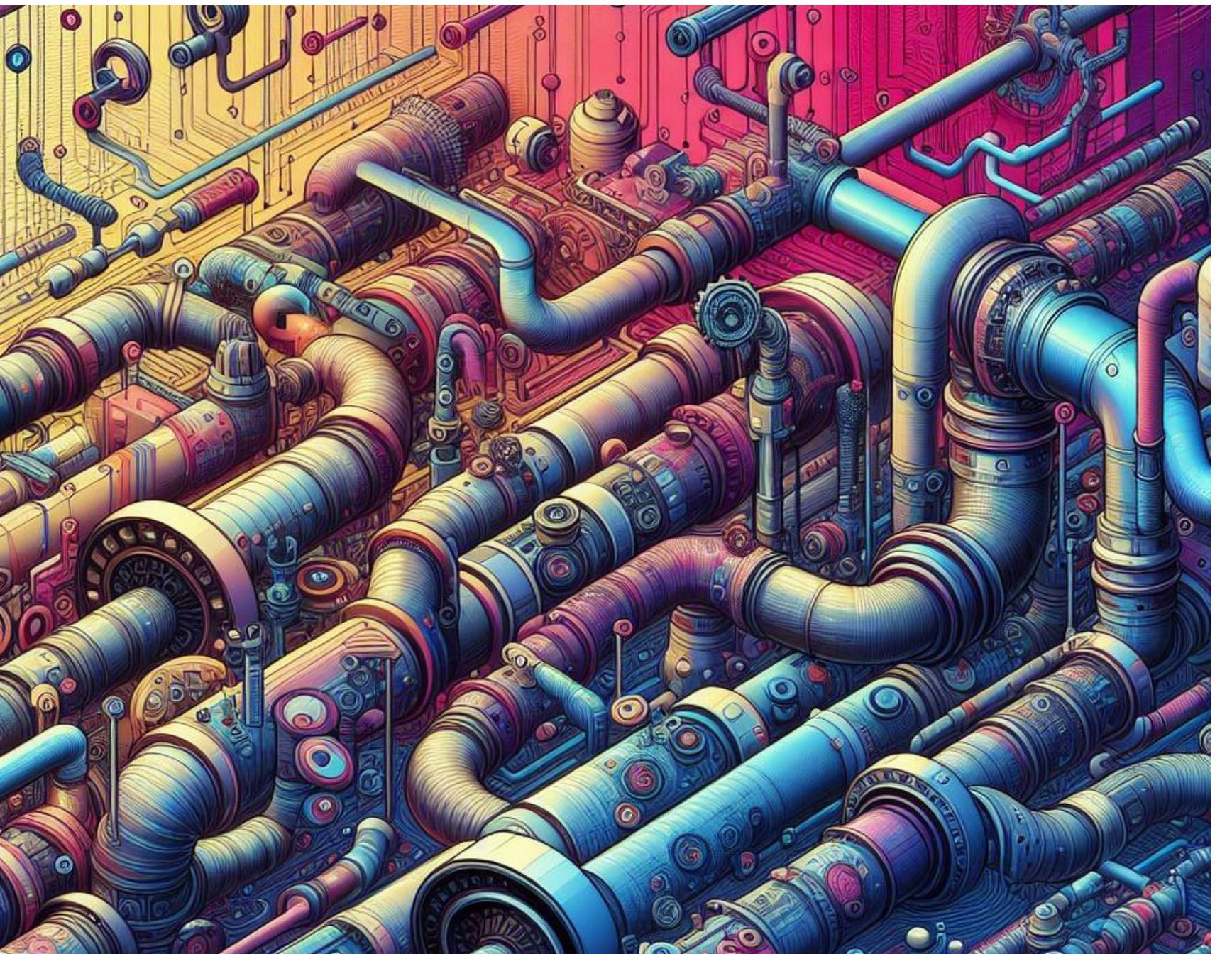
Custom-Made Framework

Based on parameters

Meta data => Azure SQL Database / Json /

Microsoft Fabric but also on Azure Synapse Analytics and
Azure Data Factory

Based on the Medallion Architecture





Who is already using Parameters?



Microsoft Fabric

The unified data platform for the era of AI



Data
Factory



Synapse Data
Engineering



Synapse Data
Science



Synapse Data
Warehousing



Synapse Real
Time Analytics



Power BI



Data
Activator



AI



OneLake



Purview

Unified
architecture

Unified
experience

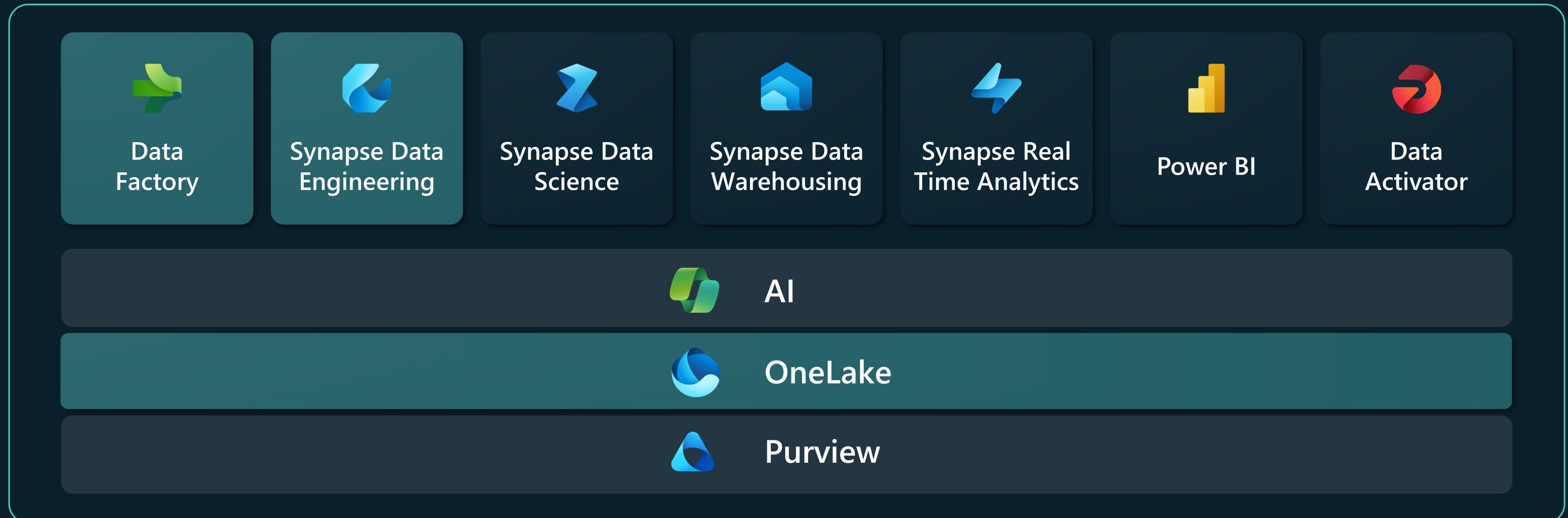
Unified
governance

Unified
business model



Microsoft Fabric

The unified data platform for the era of AI



Unified
architecture

Unified
experience

Unified
governance

Unified
business model

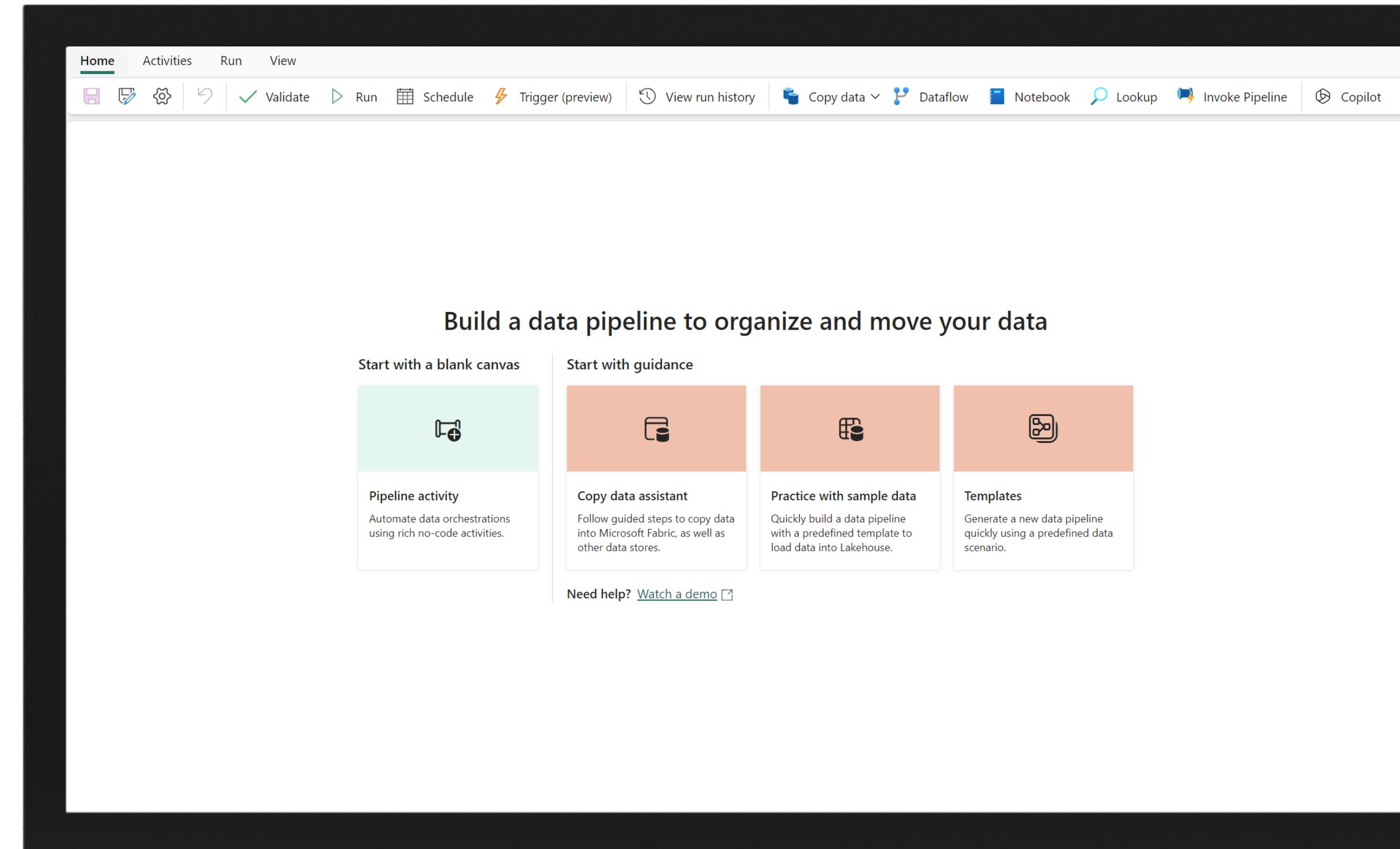
Parameters

- Pipeline Parameters
 - Copy Activity Parameters
 - Foreach Parameters
- Notebook Parameters

IMPLEMENTING
DEFAULT PARAMETERS
THAT DEPEND ON
OTHER PARAMETERS

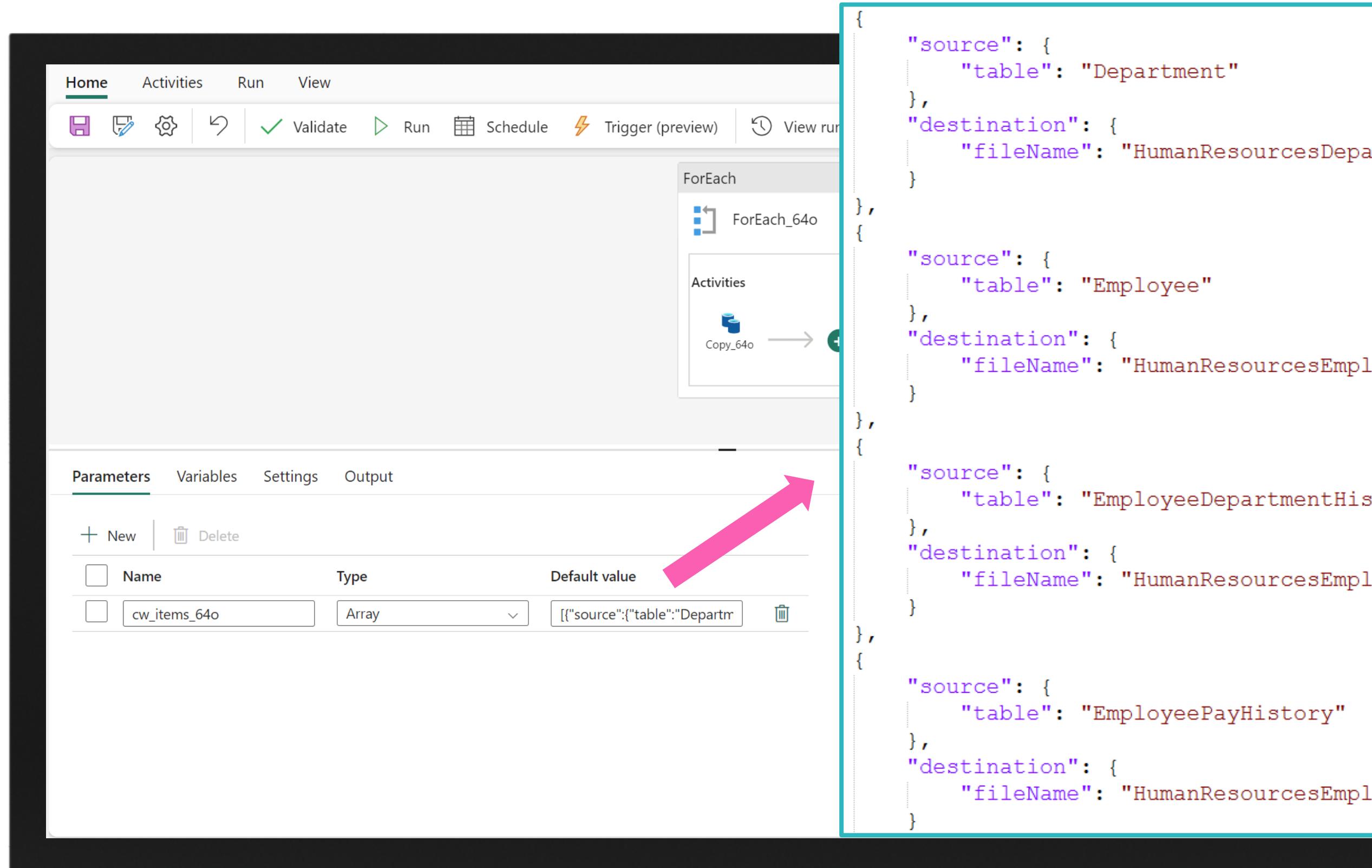
Copy Assistant

- ★ Reads data from a source data store.
- ★ Performs serialization/deserialization, compression/decompression, column mapping, and so on.
- ★ Writes data to the destination data store.



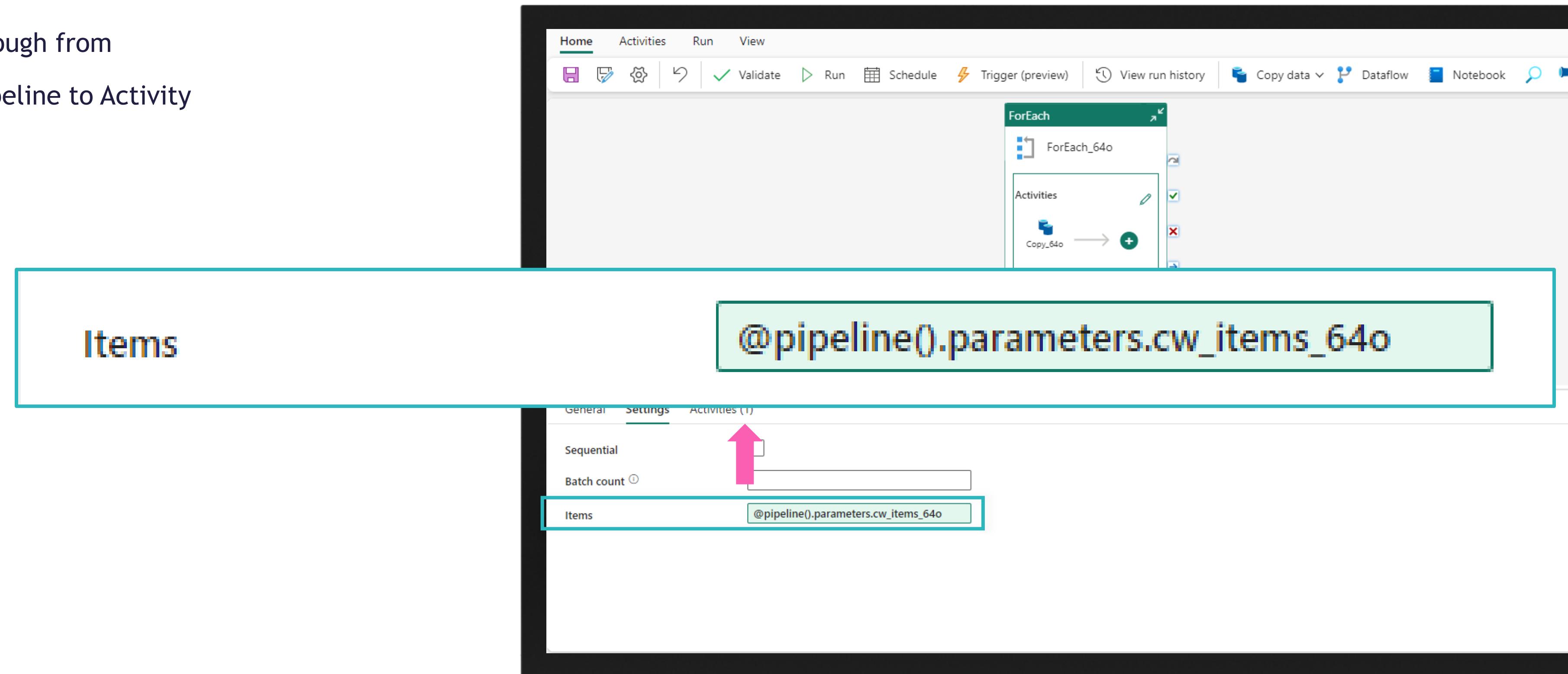
Understand Parameters

- Copy data assistant
 - Select Data Source
 - Select Tables
 - Select Destination
 - Select Filetypes
 - Review actions



Pipeline Parameters

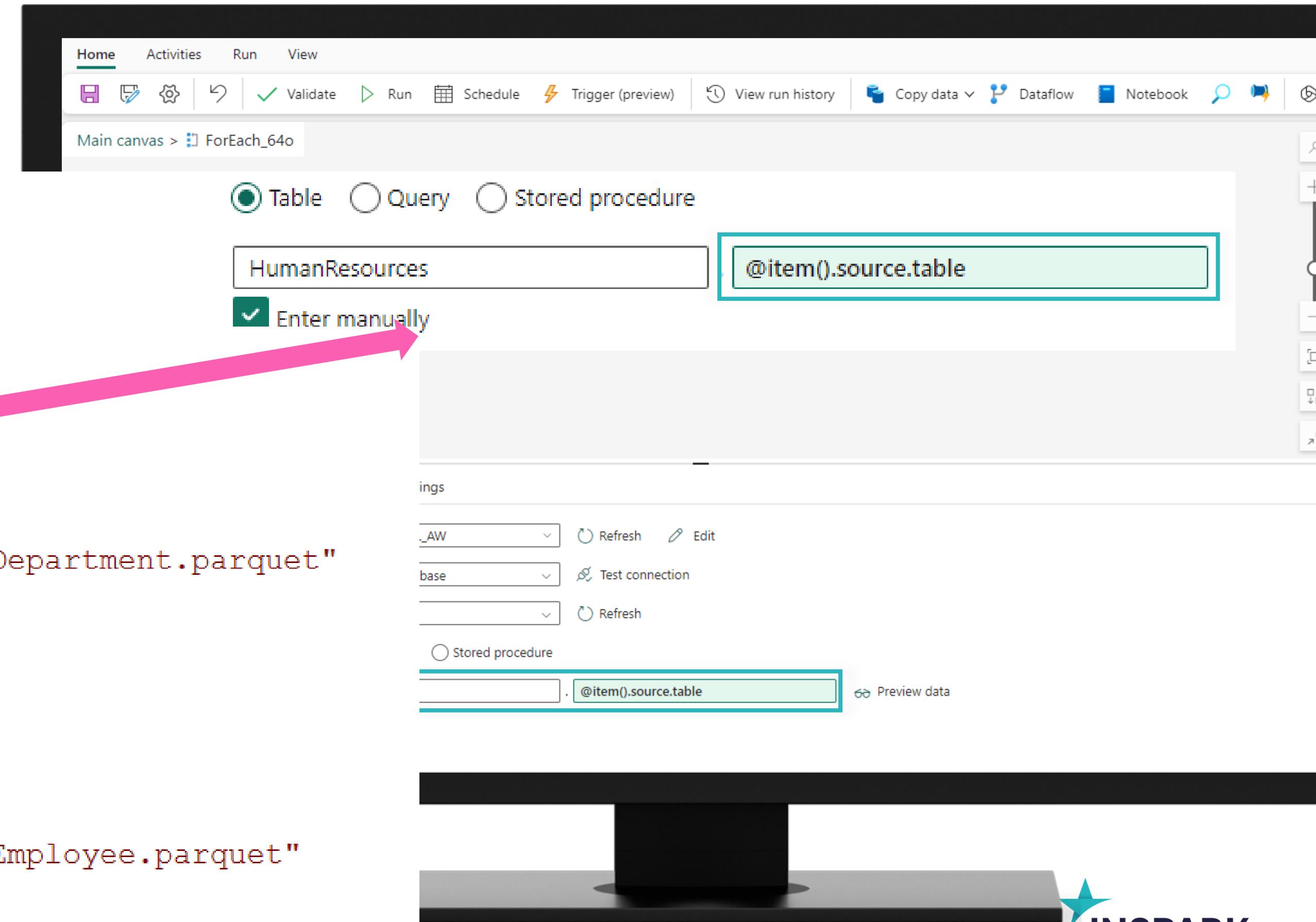
- Pass through from
 - Pipeline to Activity



Copy Activity Parameters

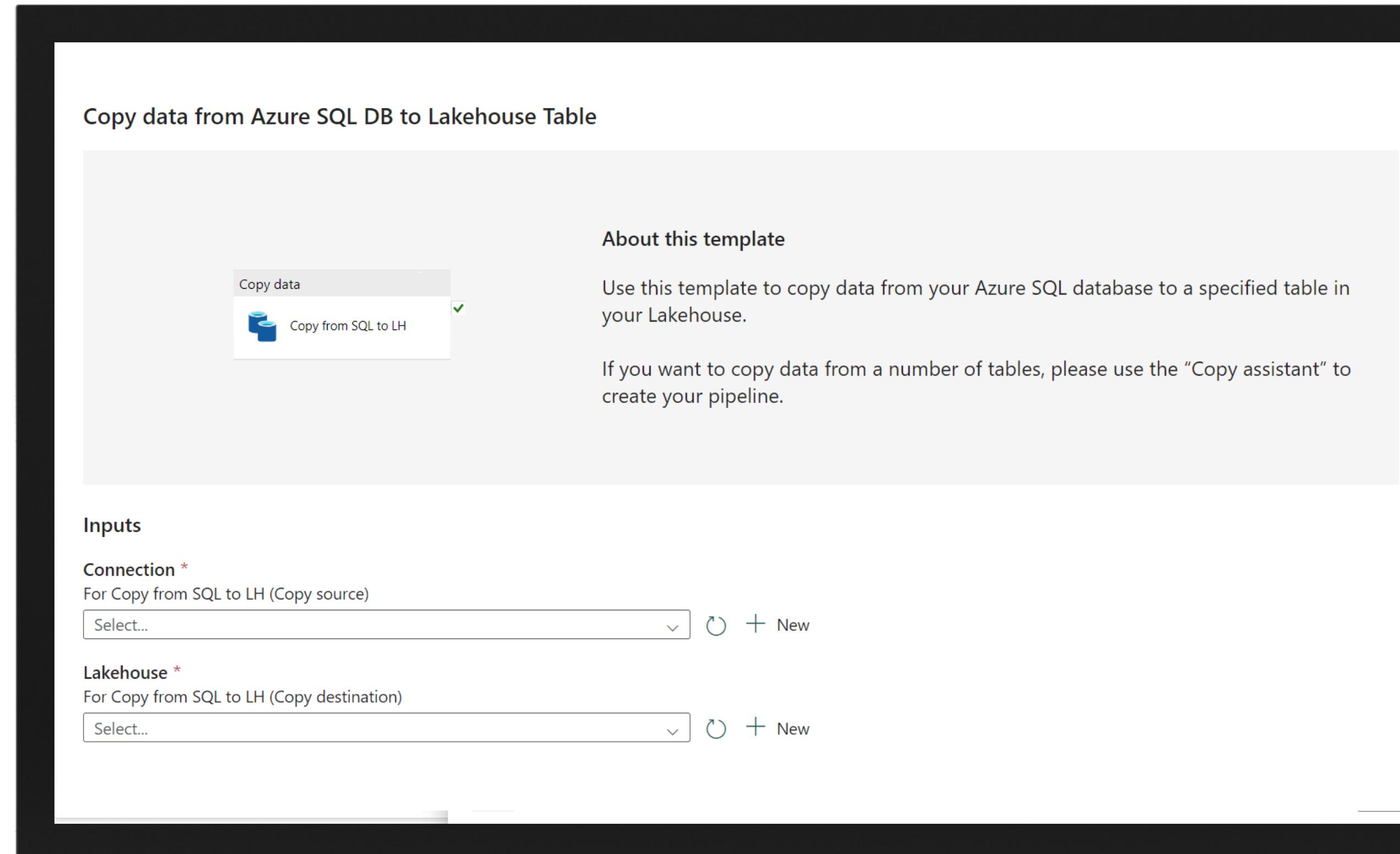
- Pass Parameters from Pipeline to Copy activity
- Use Parameters from For Each Activity

```
[  
  {  
    "source": {  
      "table": "Department"  
    },  
    "destination": {  
      "fileName": "HumanResourcesDepartment.parquet"  
    }  
  },  
  {  
    "source": {  
      "table": "Employee"  
    },  
    "destination": {  
      "fileName": "HumanResourcesEmployee.parquet"  
    }  
  },  
  {  
    "source": {  
      "table": "Employee"  
    },  
    "destination": {  
      "fileName": "HumanResourcesEmployee.parquet"  
    }  
  }]
```



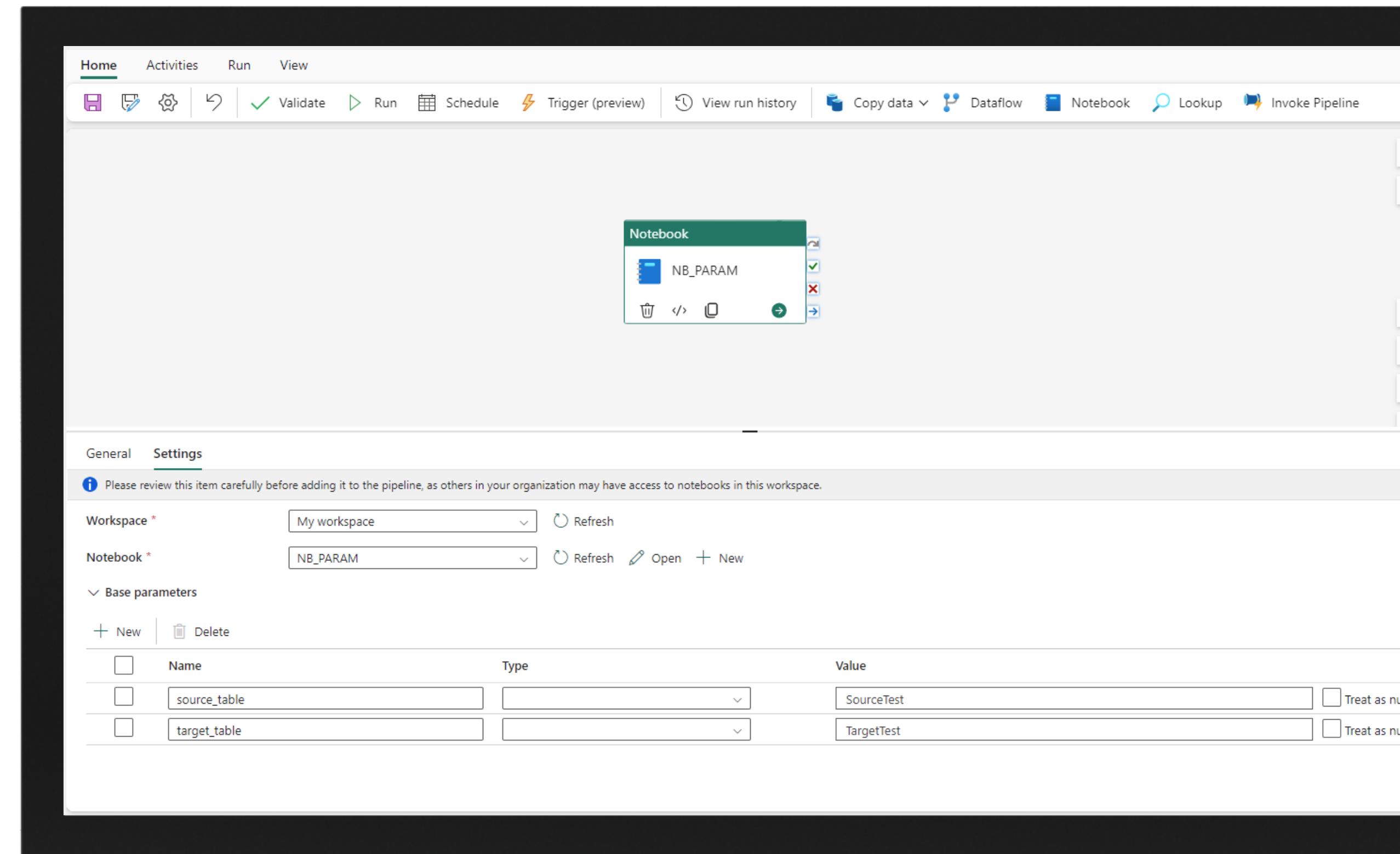
Templates

- ★ Templates are pre-defined pipelines that allow you to get started quickly with Data Factory.
- ★ These templates help to reduce development time by providing an easy way to create pipelines.
- ★ Templates are available for common data integration scenarios.
- ★ Templates can be customized to meet specific requirements



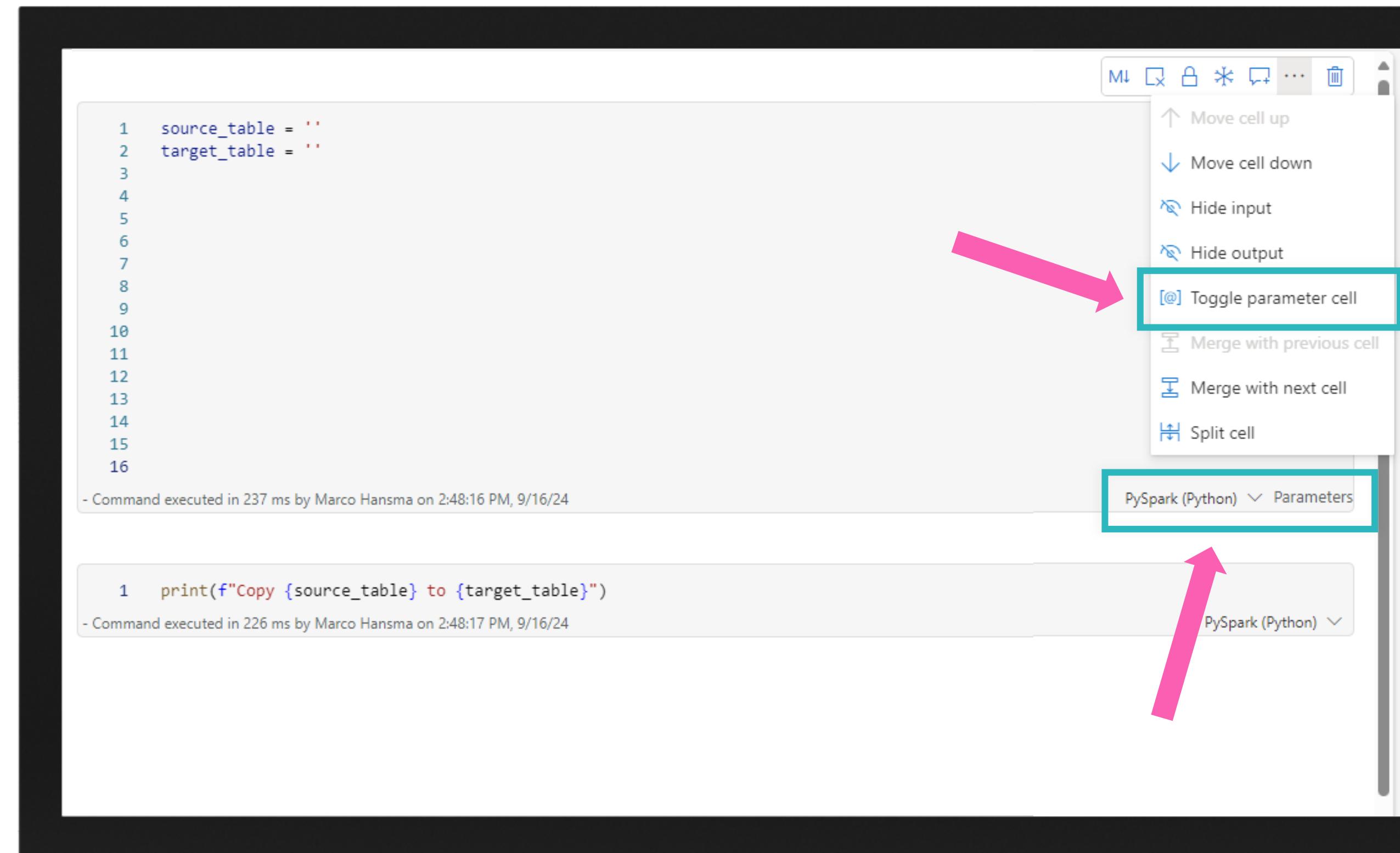
Blank canvas

★ Everything build from scratch with no predefined options



Notebook Parameters

- Pass Parameters from Data Pipeline to Notebook
 - Toggle parameter cell



```
1 source_table = ''  
2 target_table = ''  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16
```

- Command executed in 237 ms by Marco Hansma on 2:48:16 PM, 9/16/24

```
1 print(f"Copy {source_table} to {target_table}")
```

- Command executed in 226 ms by Marco Hansma on 2:48:17 PM, 9/16/24

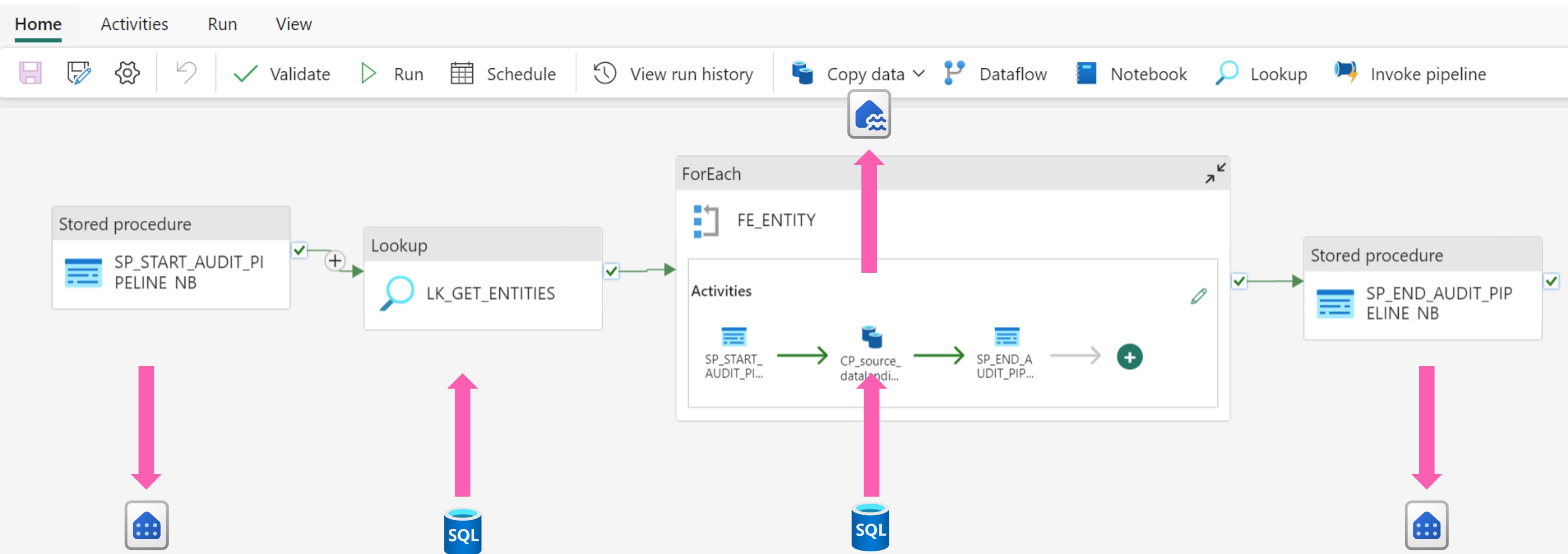
PySpark (Python) ▾ Parameters

PySpark (Python) ▾

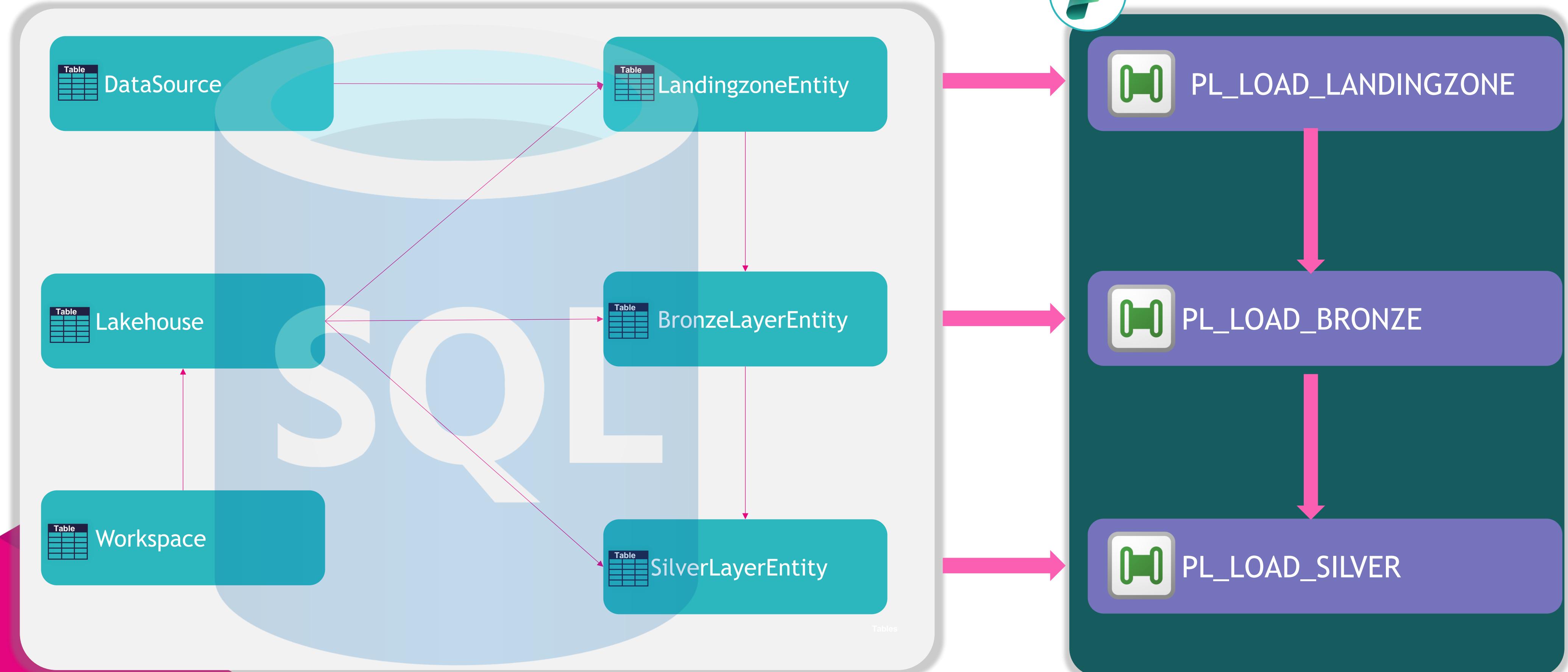
Innovate to accelerate

DEMO

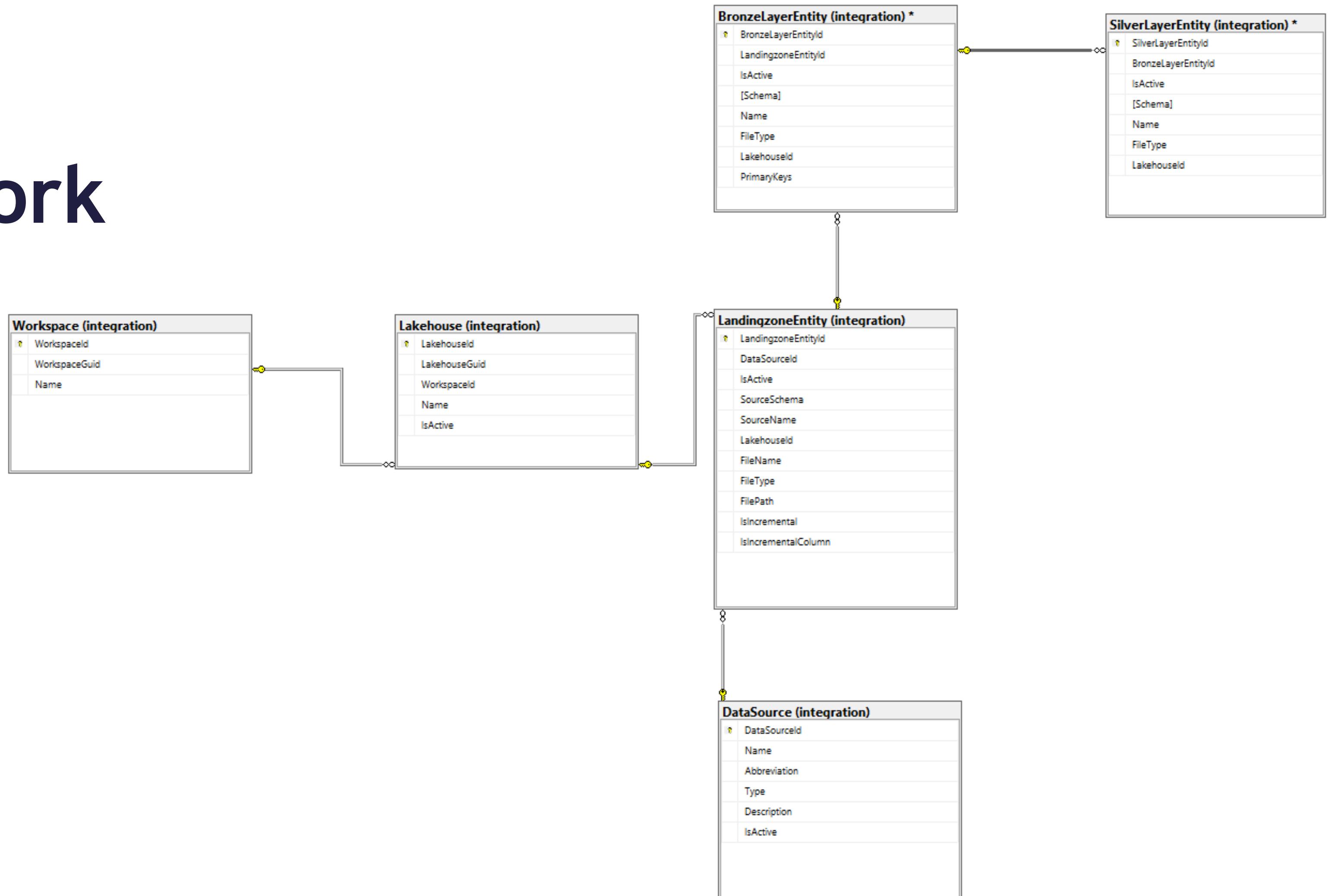
Framework in Azure SQL



Framework



Framework



Who is using a Medallion architecture?



‘Uniform data architecture’
From data “Spaghetti to Lasagna”



'Data processing in different stages'

Medallion Architecture



'Data processing in different stages'

Stage:



Gold layer



Silver layer



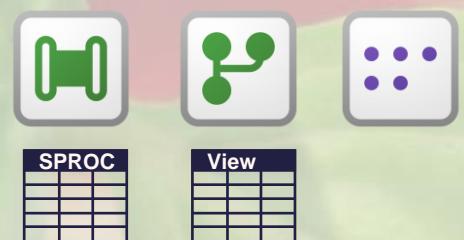
Bronze layer



Landing zone

Medallion Architecture

- Dimensions & Facts (Star Schema)
- Historical Analysis
- Business rules
- Documentation
- Aggregated data
- Logical table names



- Historical Data (Type 1 or 2)
- Data quality rules
- Data Cleansing
- Validated data
- No business model/data



- Deduplicate data
- Add datatypes
- Data can be inconsistent
- Mostly a copy of the source
- Schema



- Structured data
- Unstructured data
- Incremental loads
- Data as is
- Stored in Datetime folder structure
- No Schema

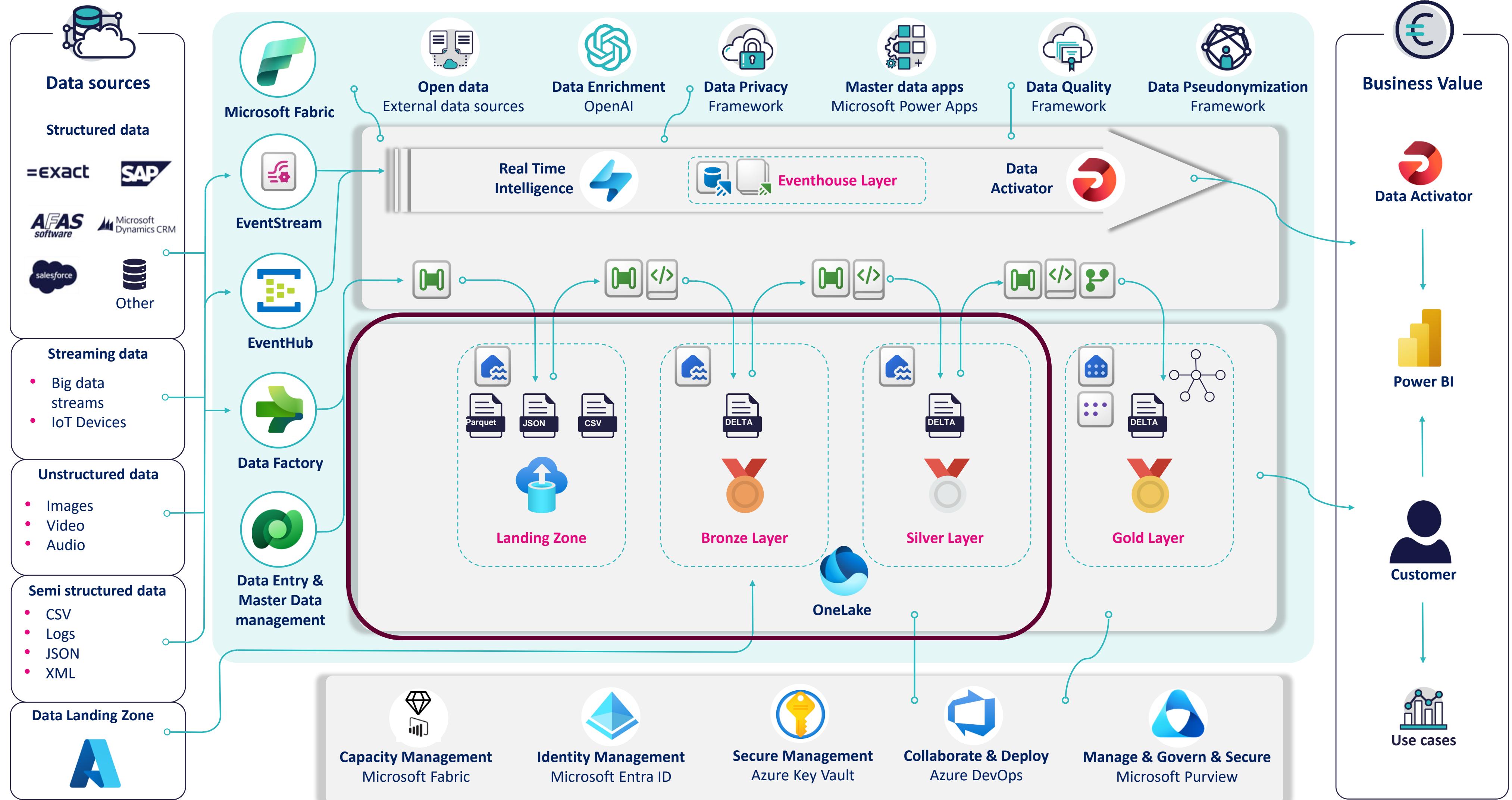


Definition:

Filertype:

Files/Tables:

Fabric:



Lakehouse

	Name	Type
📁	LH_Bronze_Layer	Lakehouse
└─	📁 LH_Bronze_Layer	Semantic model (default)
└─	🏠 LH_Bronze_Layer	SQL analytics endpoint
📁	LH_Data_Landingzone	Lakehouse
└─	📁 LH_Data_Landingzone	Semantic model (default)
└─	🏠 LH_Data_Landingzone	SQL analytics endpoint
📁	LH_Silver_Layer	Lakehouse
└─	📁 LH_Silver_Layer	Semantic model (default)
└─	🏠 LH_Silver_Layer	SQL analytics endpoint

Copy Activity Parameters

- ★ Pass Parameters from Pipeline to Copy activity
- ★ Use Parameters from For Each Activity

Preview data

ceName	DataSourceAbbreviation	DataSourceType	IsActive	SourceSchema	SourceName	TargetFilePath	TargetFileName	TargetFileType	TargetFileFormat
JwdvImdl01	ADLS	ADLS	true		customers.csv	demo/2024/09/10	customers_2024091015.csv	csv	Binary
JwdvImdl01	ADLS	ADLS	true		organizations.csv	demo/2024/09/10	organizations_2024091015.csv	csv	Binary
JwdvImdl01	ADLS	ADLS	true		people.csv	demo/2024/09/10	people_2024091015.csv	csv	Binary

General Source **Destination** Mapping Settings

Connection * `@item().TargetLakehouseGuid`

Connection type * `Lakehouse`

Workspace ID `@item().WorkspaceGuid`

Root folder Tables Files

File path `@item().TargetFilePath` / `@item().TargetFileName`

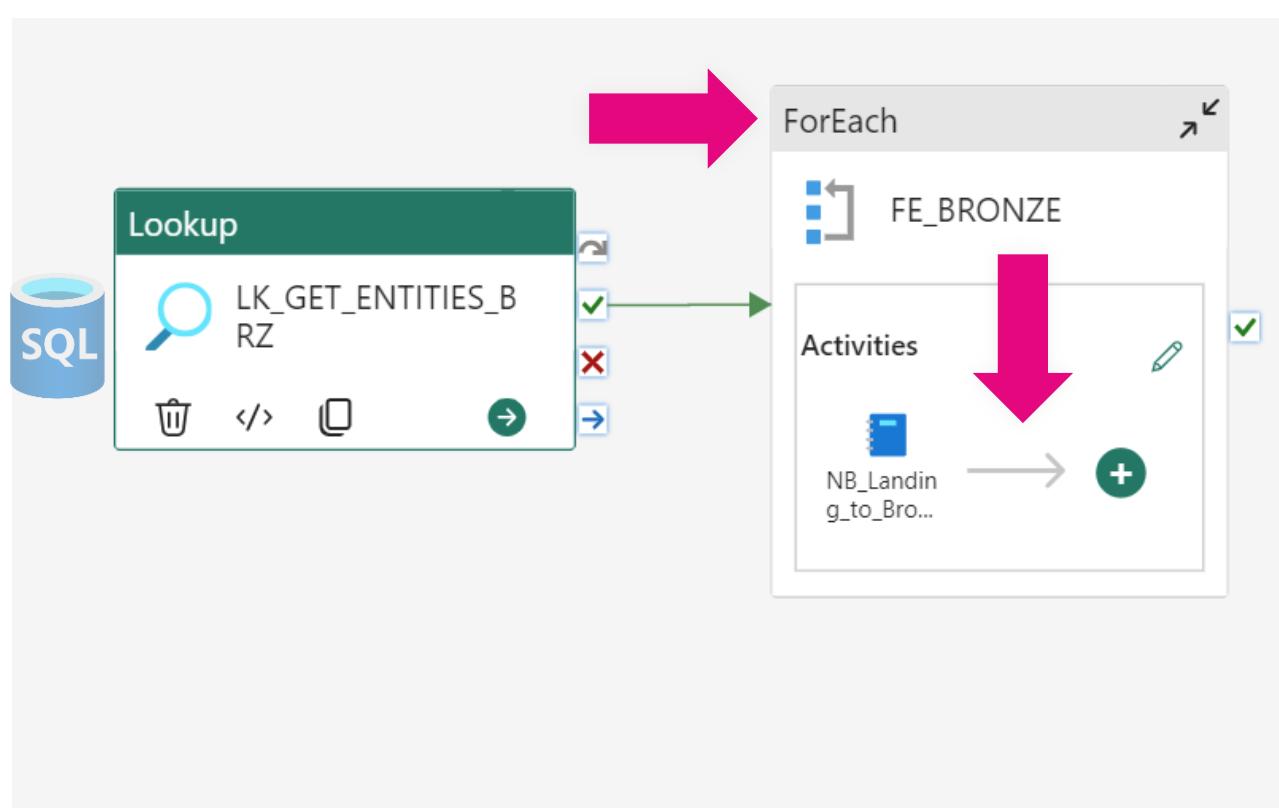
File format * `Binary`

> Advanced

A screenshot of the Azure Data Factory UI. On the left, a pipeline diagram shows a 'Lookup' activity (SQL source) connected to a 'ForEach' activity. Inside the 'ForEach' loop, a 'Copy CSV' activity is listed. A pink arrow points from the 'ForEach' activity to the 'Destination' tab of the 'Copy CSV' activity settings. On the right, the 'Destination' tab is selected, showing configuration for a Lakehouse connection, workspace ID, file path, and file format. Two pink arrows point upwards from the 'File path' and 'File format' fields to the corresponding parameters in the pipeline diagram, illustrating how pipeline parameters are passed to the copy activity.

Notebook Parameters

- ★ Pass Parameters from Data Pipeline to Notebook
- ★ Toggle parameter cell



Preview data

#	EntityId	SourceFilePath	SourceFileName	SourceFileType	TargetSchema	TargetName	TargetWorkspaceId	SourceWorkspaceId	TargetLakehouseId	SourceLakehouseId	TargetLakehouseName	SourceLakehouseName
1	1	demo/2024/09/10	customers_2024091017.csv	csv	demo	Customers	586fc19d-fa6a-4cb1-9ca3-e518a524f5da	586fc19d-fa6a-4cb1-9ca3-e518a524f5da	e130dba7-c8c3-438a-85ad-2cd4c9d59a09	009058ac-b71c-4774-a8ee-7c7d945c3972	LH_Bronze_Layer	LH_Data_Landingzor
2	3	demo/2024/09/10	organizations_2024091017.csv	csv	demo	Organizations	586fc19d-fa6a-4cb1-9ca3-e518a524f5da	586fc19d-fa6a-4cb1-9ca3-e518a524f5da	e130dba7-c8c3-438a-85ad-2cd4c9d59a09	009058ac-b71c-4774-a8ee-7c7d945c3972	LH_Bronze_Layer	LH_Data_Landingzor
3	4	demo/2024/09/10	people_2024091017.csv	csv	demo	People	586fc19d-fa6a-4cb1-9ca3-e518a524f5da	586fc19d-fa6a-4cb1-9ca3-e518a524f5da	e130dba7-c8c3-438a-85ad-2cd4c9d59a09	009058ac-b71c-4774-a8ee-7c7d945c3972	LH_Bronze_Layer	LH_Data_Landingzor

Notebook * NB_LANDING_BRONZE

Base parameters

Name	Type	Value
SourceLakehouse	String	@item().SourceLakehouseId
source_file_path	String	@item().SourceFilePath
source_file_name	String	@item().SourceFileName
PrimaryKeys	String	@item().PrimaryKeys
TargetLakehouse	String	@item().TargetLakehouseId
target_schema	String	@item().TargetSchema
target_name	String	@item().TargetName
SourceWorkspace	String	@item().SourceWorkspaceId
TargetWorkspace	String	@item().TargetWorkspaceId
source_file_type	String	@item().SourceFileType

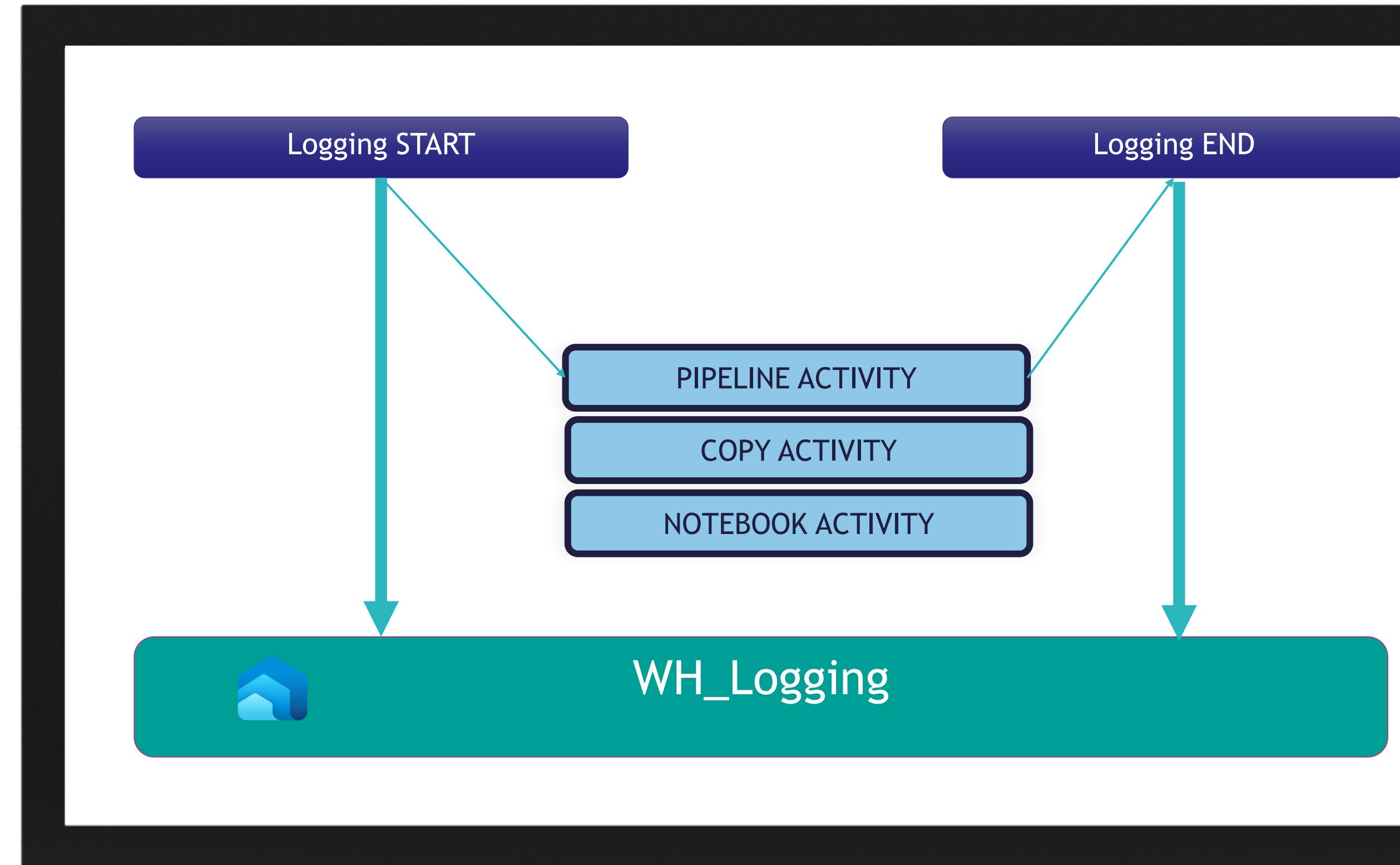
Innovate to accelerate

DEMO

Logging

Logging

- ★ Log Start and End Time of records
- ★ Log Extracted Records
- ★ Log Execution Failure



Logging

- ★ Log Start and End Time of records
- ★ Log Extracted Records
- ★ Log Execution Failure

The screenshot shows the Pipeline expression builder interface. On the left, there's a large black rectangular redaction box covering most of the left side of the screen. To its right is a smaller redaction box at the bottom. The main area contains a form for configuring a pipeline:

General Settings

Data store type: Workspace (selected) External
Warehouse: WH_Logging
Stored procedure name: [logging].[sp_AuditPipeline]

Stored procedure parameters

Name	Type
LogData	String
LogType	String
PipelineGuid	Guid
PipelineName	String
PipelineParameters	String
PipelineParentRunGuid	Guid
PipelineRunGuid	Guid
TriggerGuid	Guid
TriggerTime	DateTime
TriggerType	String
WorkspaceGuid	Guid

Pipeline expression builder

Add dynamic content below using any combination of expressions, functions and system variables.

@pipeline().Pipeline

Clear contents

Parameters System variables Functions Variables

Search

Pipeline ID

ID of the pipeline

Pipeline Name

Name of the pipeline

Pipeline group ID

ID of the group to which the pipeline run belongs

Pipeline run ID

ID of the specific pipeline run

Pipeline trigger ID

ID of the trigger that invokes the pipeline

Pipeline trigger time

Time when the trigger that invoked the pipeline. The trigger time is the actual fired time, not the sched...

Pipeline trigger type

Type of the trigger that invoked the pipeline (Manual, Scheduler)

Pipeline triggered by pipeline ID

ID of the pipeline that triggered this pipeline. Applicable when a pipeline run is triggered by an Execut...

Pipeline triggered by pipeline name

Name of the pipeline that triggered this pipeline. Applicable when a pipeline run is triggered by an Ex...

Pipeline triggered by pipeline run ID

Run ID of the pipeline that triggered this pipeline. Applicable when a pipeline run is triggered by an Ex...

Workspace ID

ID of the workspace the pipeline run is running within

Logging

- ★ Add Information about pipelines
- ★ Adding System Variables
- ★ Add Information about Notebooks

Pipeline expression builder

Add dynamic content below using any combination of [expressions](#), [functions](#) and [system variables](#).

```
{  
    "Action" : "End",  
    @{{activity('NB_Landing_to_Bronze').output.result.exitValue}}  
}
```



Innovate to accelerate

DEMO

Challenges



Out of the box fast and easy, less flexible



Parameterize of connections with Azure Key Vault
(Like ADF/Synapse with Linked Services)



Schedule can't be Parameterized like in
ADF/Synapse



Build in retry to Notebook Activity

Key Learnings and Best Practices



Metadata-driven approach is a best practice for managing data pipelines



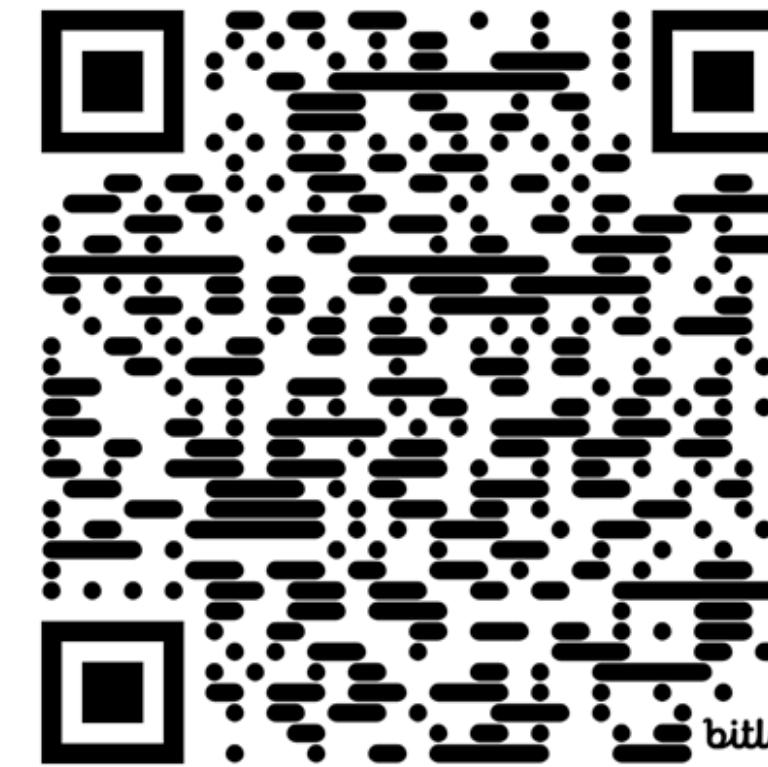
Medallion Lakehouse architecture is a proven framework for implementing metadata-driven approach



Start small and gradually expand the metadata-driven approach across the organization



Session Feedback



https://bit.ly/dMC2024_SessionFeedback

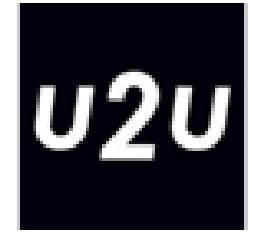
Thank you, partners



PLAINSIGHT



LACO/



delaware

EpicData.



KOHERA



TriFinance
BEYOND ADVISORY



@erwindekreuk



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



erwindekreuk.com



github.com/edkreuk



<https://sessionize.com/erwin-de-kreuk/>



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

