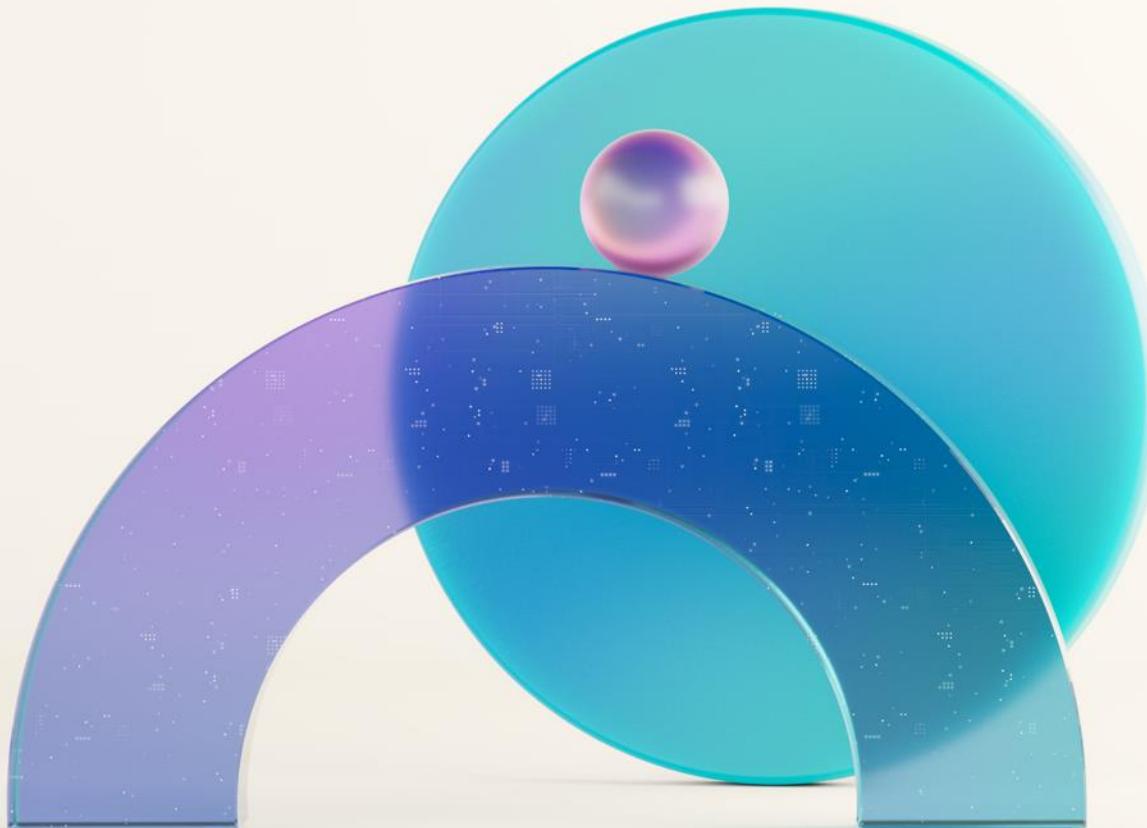
A minimalist, abstract background featuring a large, translucent pink triangle on the left and a series of overlapping, semi-transparent rectangular planes in shades of cyan, green, and blue. A single, smooth sphere sits atop one of the blue-green planes.

Microsoft Fabric

COMMUNITY CONFERENCE



Drive Efficiency and Reliability Using Metadata-Driven Frameworks

March 26-28, 2024

Erwin de Kreuk

Erwin de Kreuk

- Principal Consultant
- Lead Data & Analytics InSpark



Let's connect



 @erwindekreuk

 linkedin.com/in/erwindekreuk

 erwindekreuk.com

 github.com/edkreuk

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



Objectives



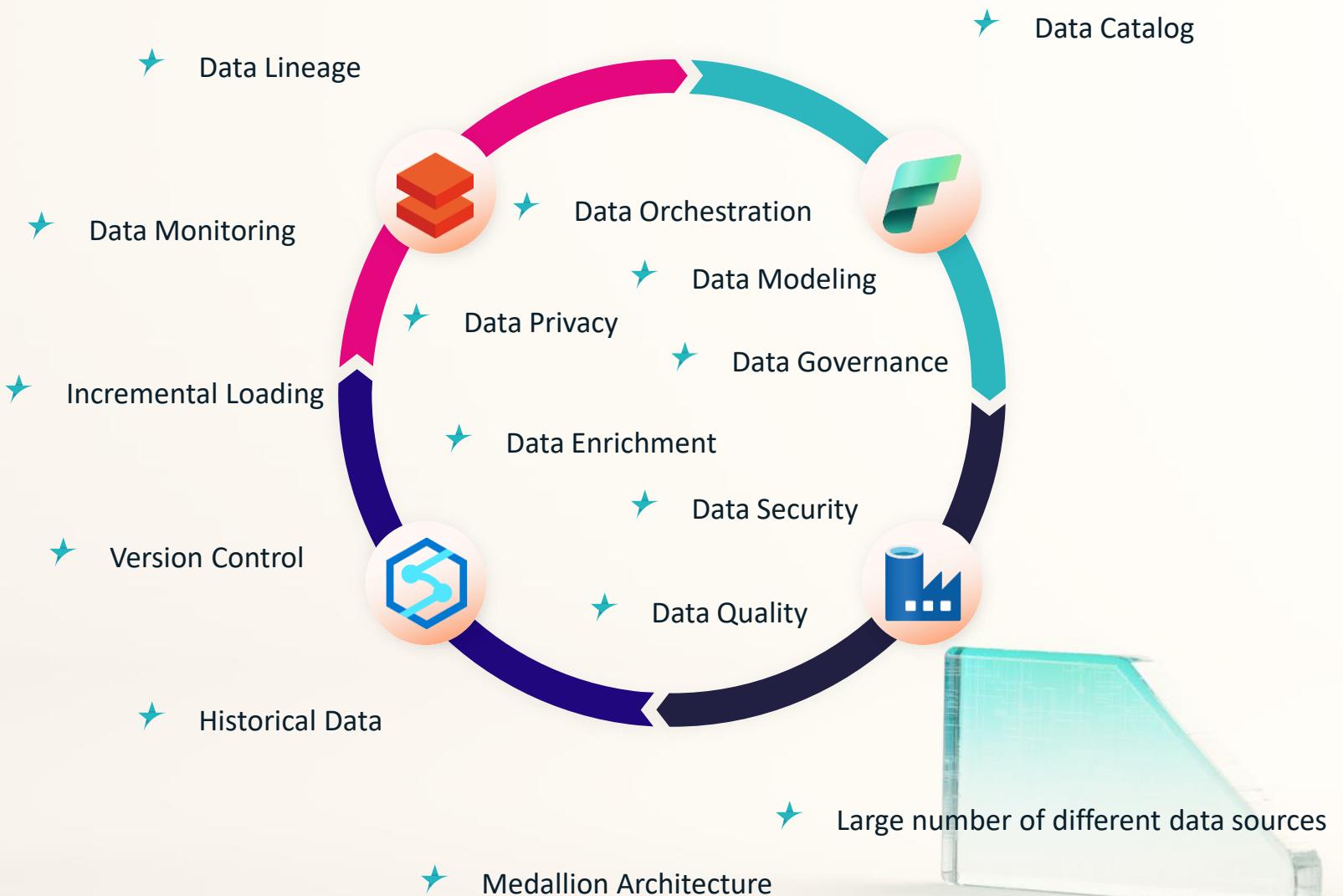
- Data platform challenges
- Why
- Medallion Architecture
- Parameters
- Framework
- Logging
- Recap





Data platform Challenges

'From data source to data model to report'





What would help

Simplicity in connecting data sources

Fast result in 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

Integration of solution like:

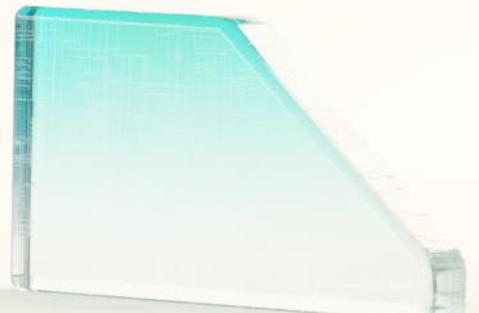
- Data Privacy
- Data Quality
- Data Governance

From data “Spaghetti to Lasagna”

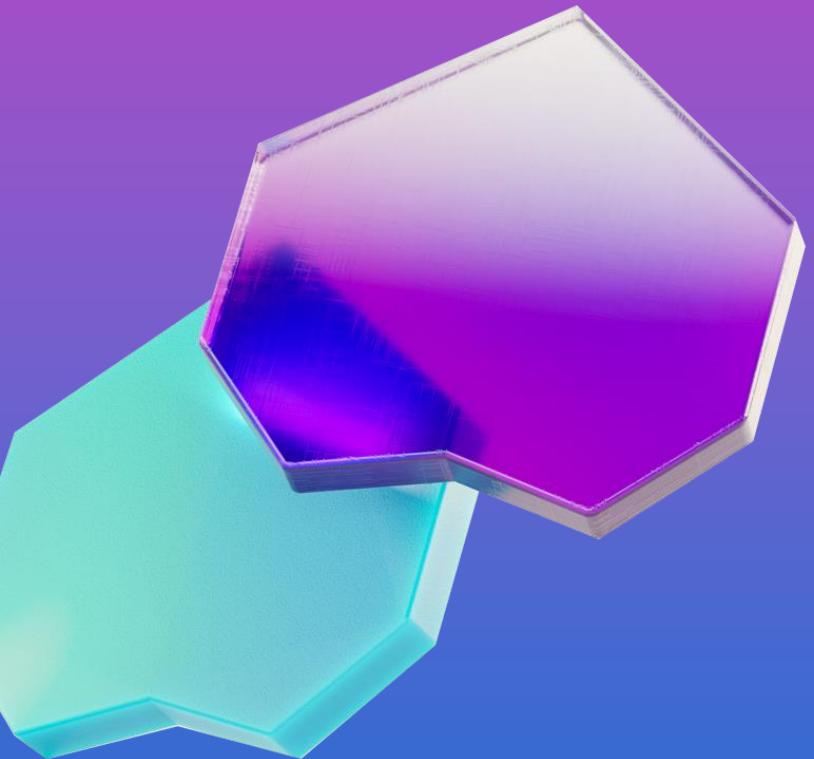
- Uniform data architecture

*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



Automation



Flexibility

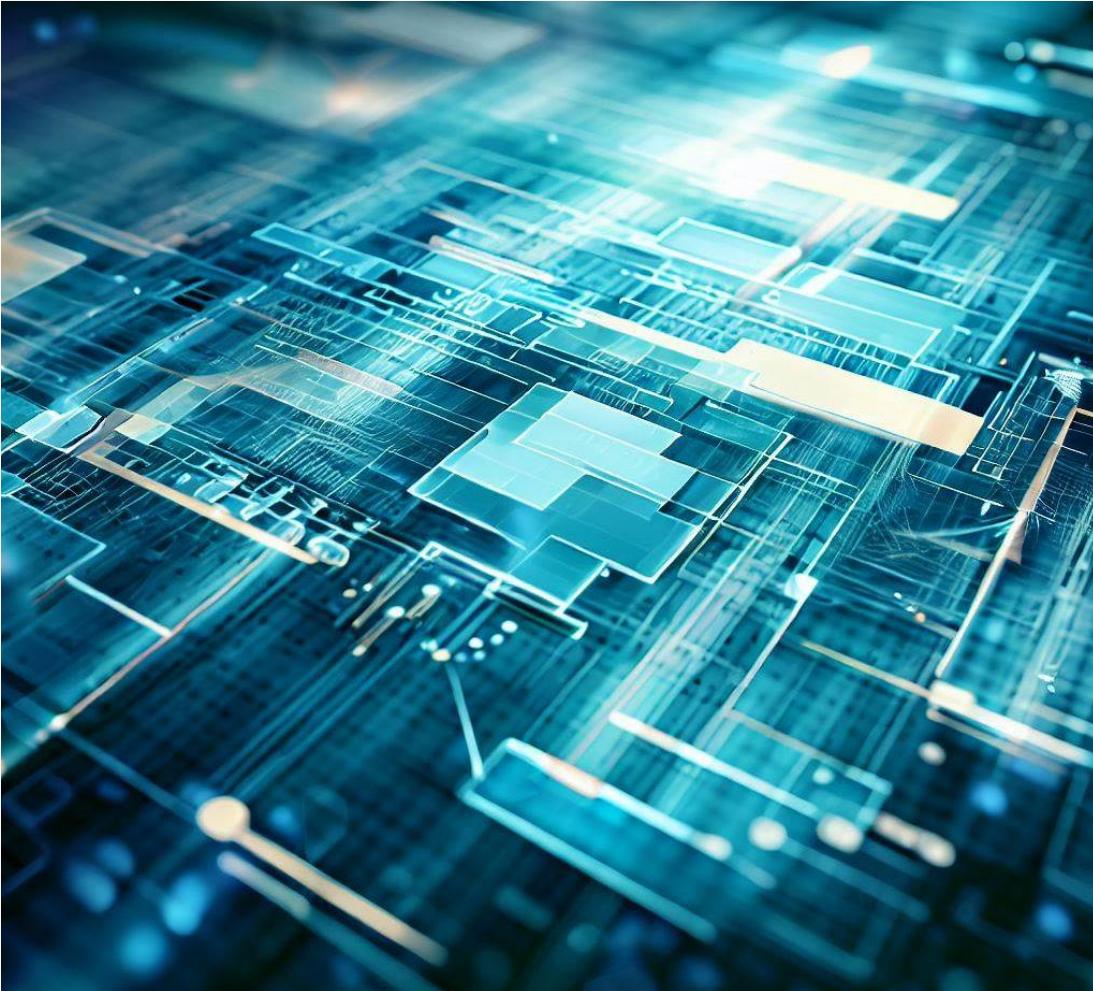


Scalability



Traceability

Out-of-the-Box Framework



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



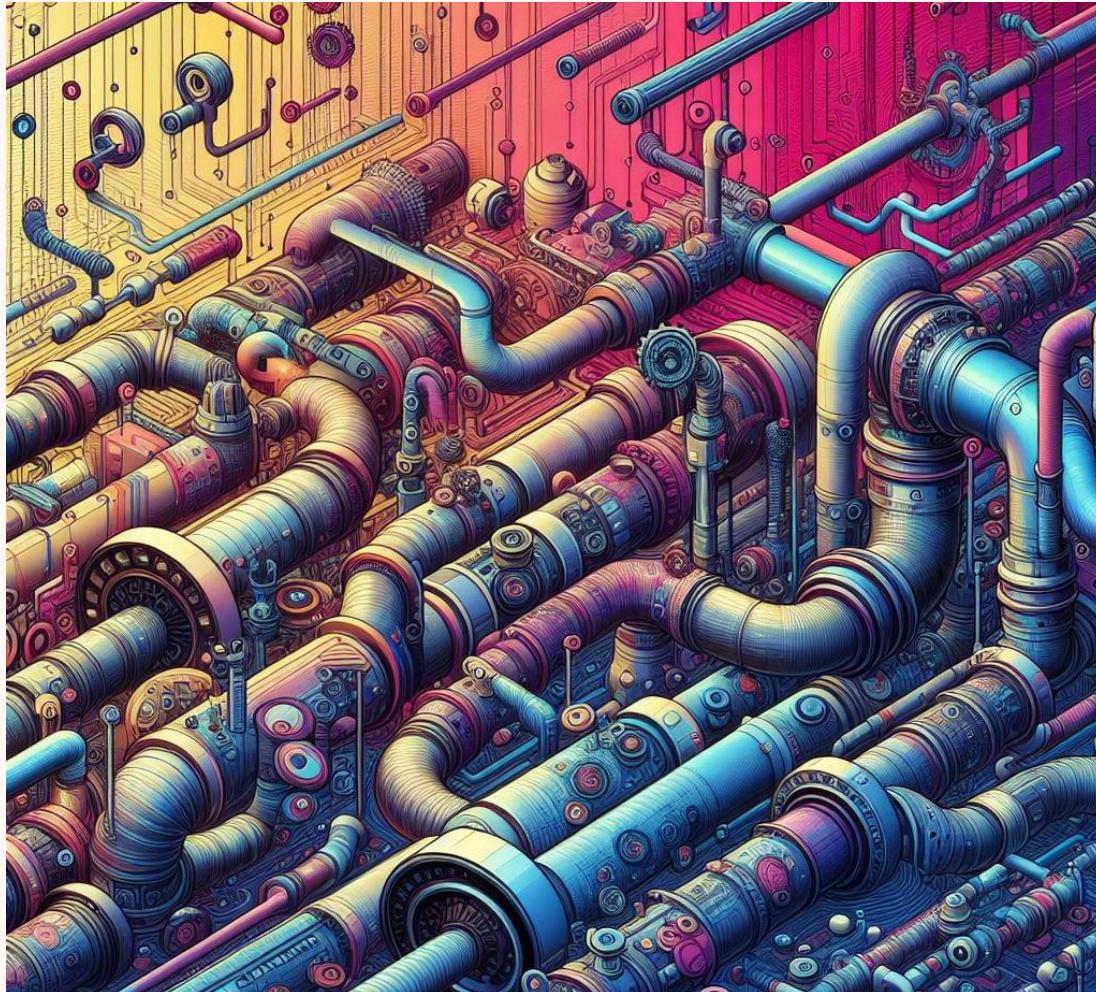
Custom-Made Framework



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



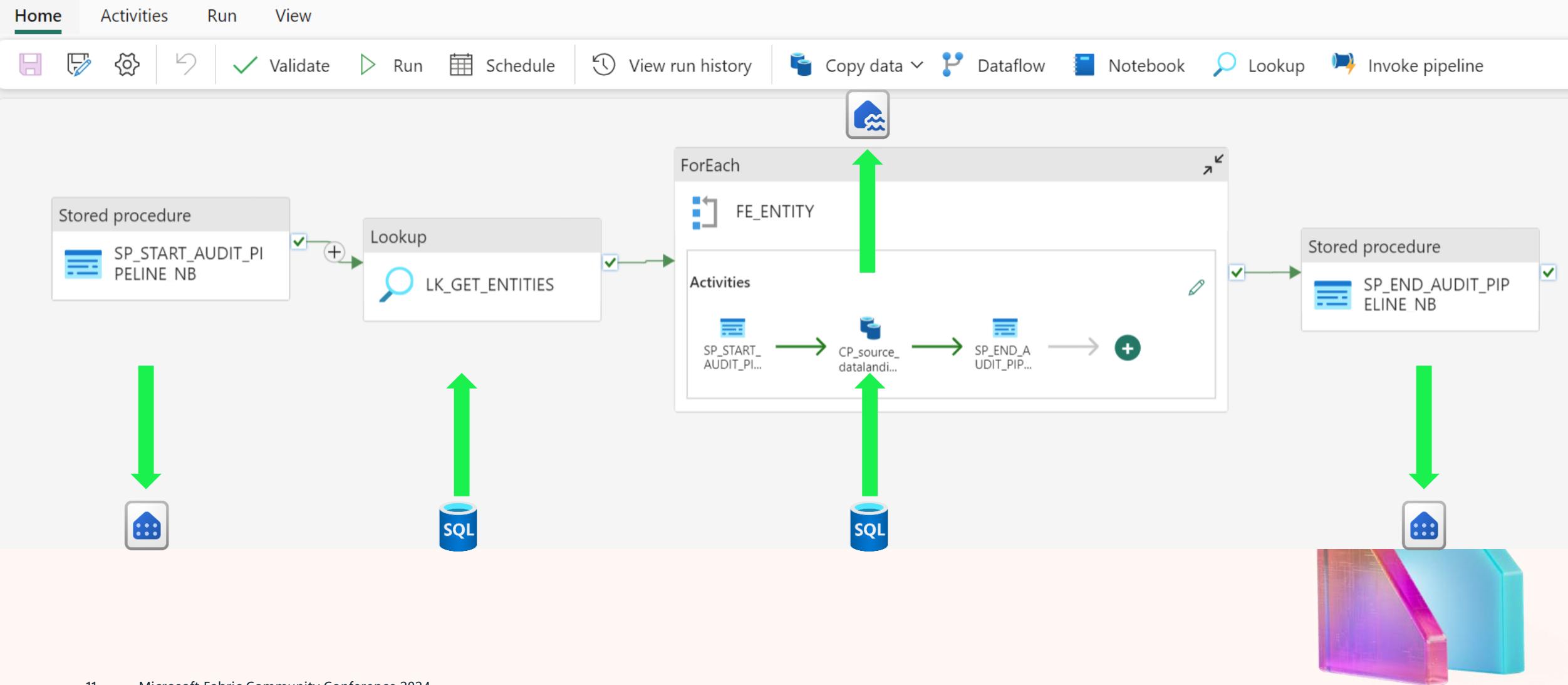
Custom-Made Framework



- 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



FINAL GOAL for today's session



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



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

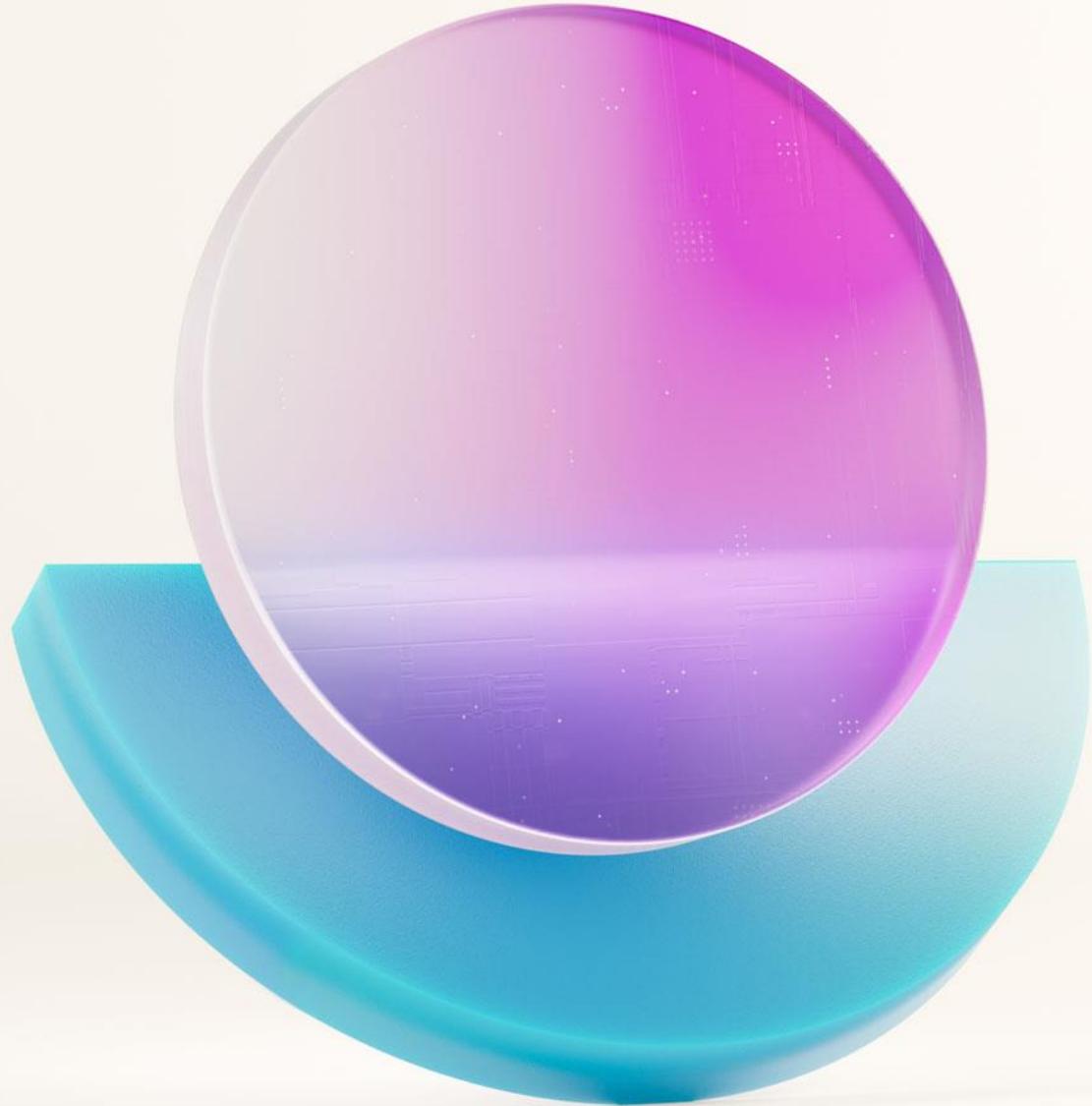
Unified
governance

Unified
business model

Lakehouse

- The foundation of Microsoft Fabric is a **Lakehouse**, which is built on top of the **OneLake** scalable storage layer and uses **Apache Spark and SQL** compute engines for big data processing.
- **Lakehouses use Spark and SQL engines** to process large-scale data and support machine learning or predictive modeling analytics.
- **Lakehouse data is organized in a *schema-on-read format***, which means you define the schema as needed rather than having a predefined schema.
- **Lakehouses support ACID** (Atomicity, Consistency, Isolation, Durability) transactions through Delta Lake formatted tables for data consistency and integrity.
- **Lakehouses are a single location** for data engineers, data scientists, and data analysts to access and use data.

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



Medallion architecture

Who is using a Medallion architecture?



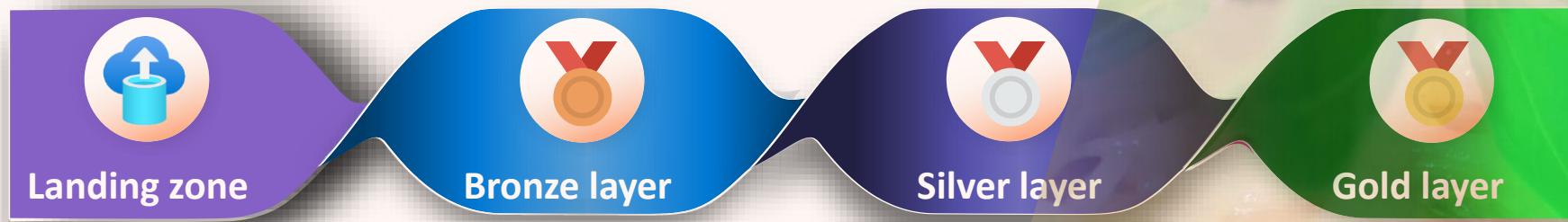
*‘Uniform data architecture’
From data “Spaghetti to Lasagna”*



Medallion Architecture

'Data processing in different stages'

Stages



Medallion Architecture

'Data processing in different stages'

Stage:



Gold layer



Silver layer



Bronze layer



Landing zone

Definition:

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

Filetype:



Files/Tables:



Fabric:



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

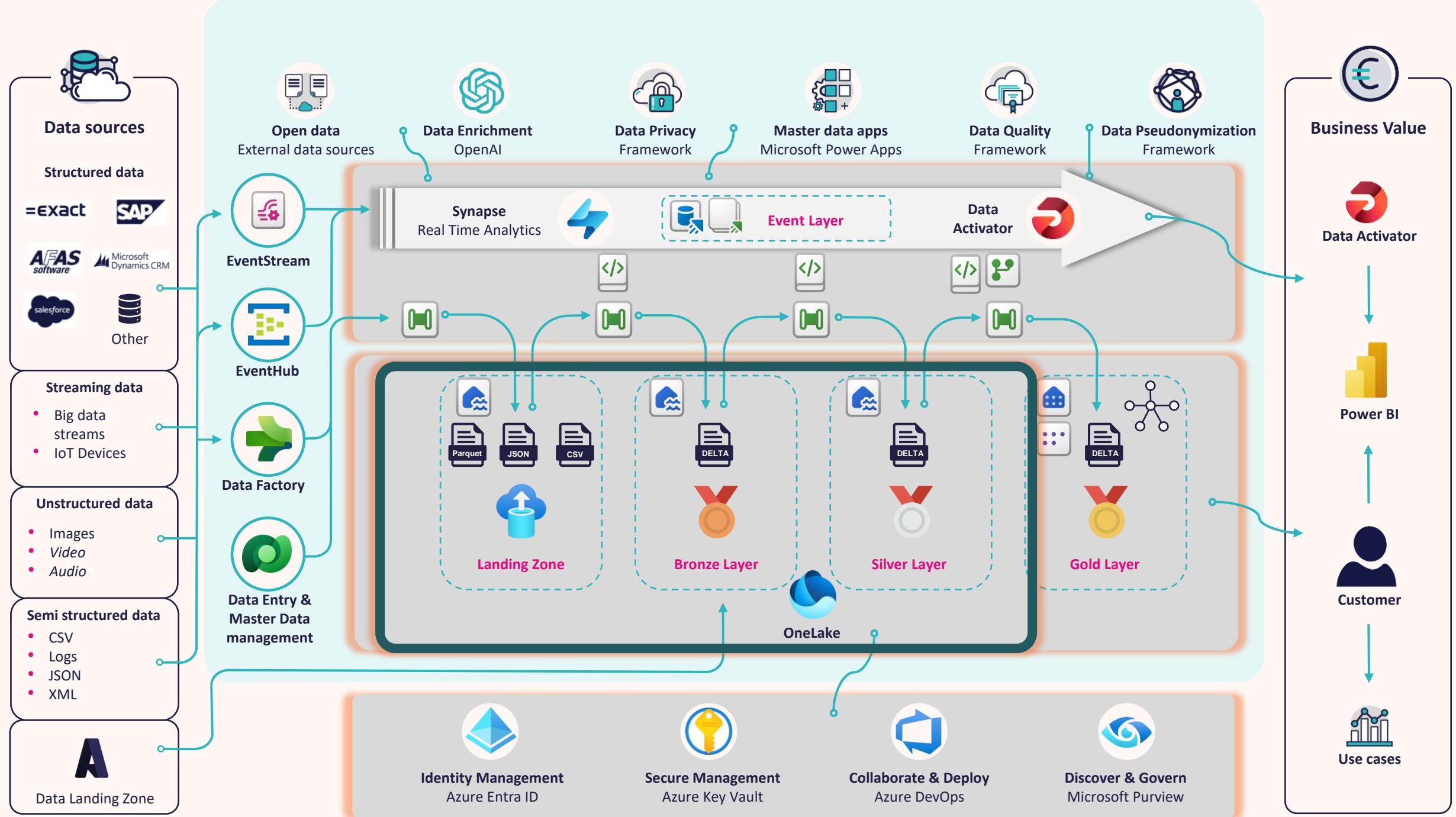


- Deduplicate data
- Add datatypes
- Data can be inconsistent
- Mostly a copy of the source
- Schema



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





Understand Parameters

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

Copy data from Azure SQL DB to Lakehouse Table

About this template

Use this template to copy data from your Azure SQL database to a specified table in your Lakehouse.

If you want to copy data from a number of tables, please use the "Copy assistant" to create your pipeline.

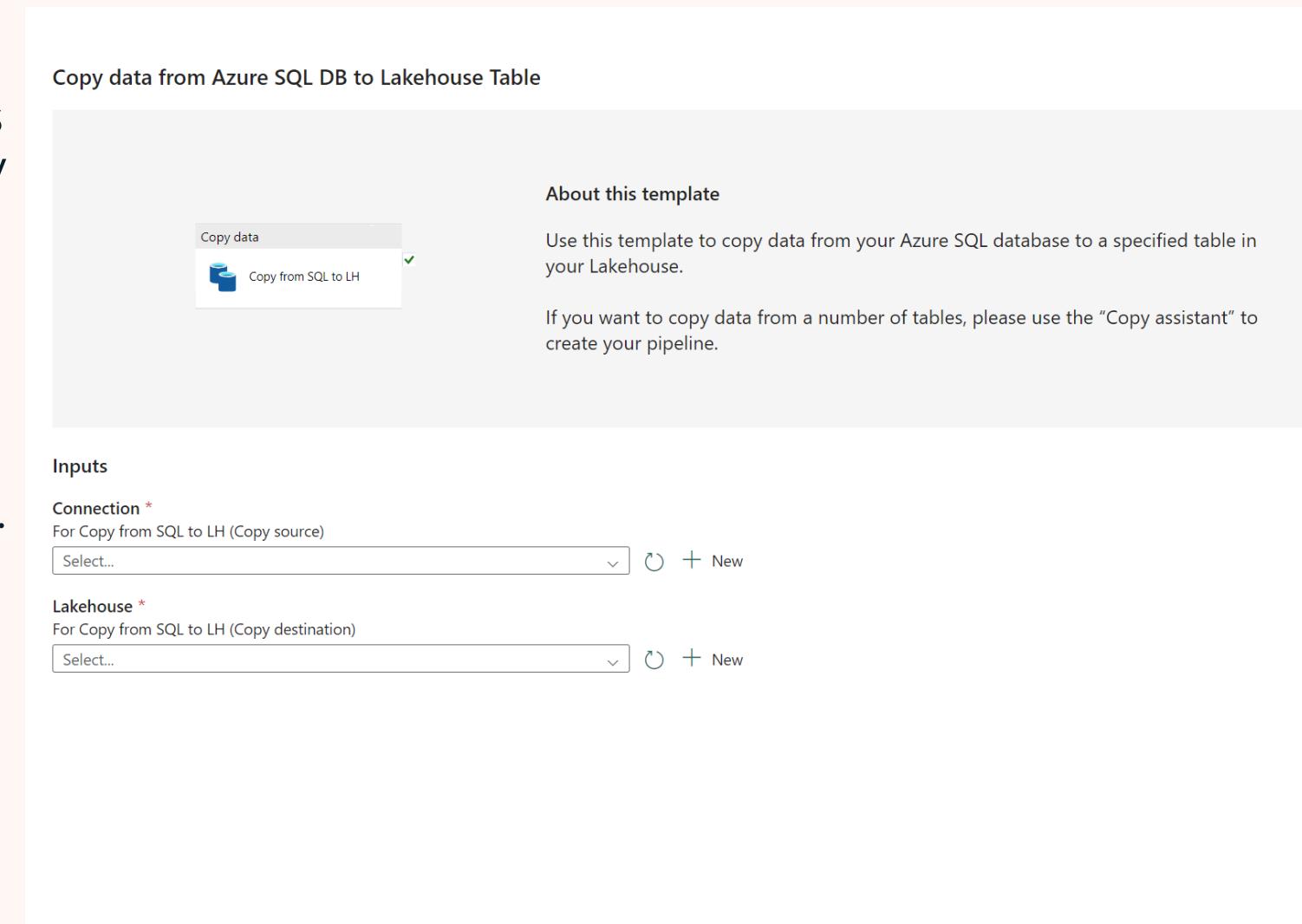
Inputs

Connection *
For Copy from SQL to LH (Copy source)

Select...  

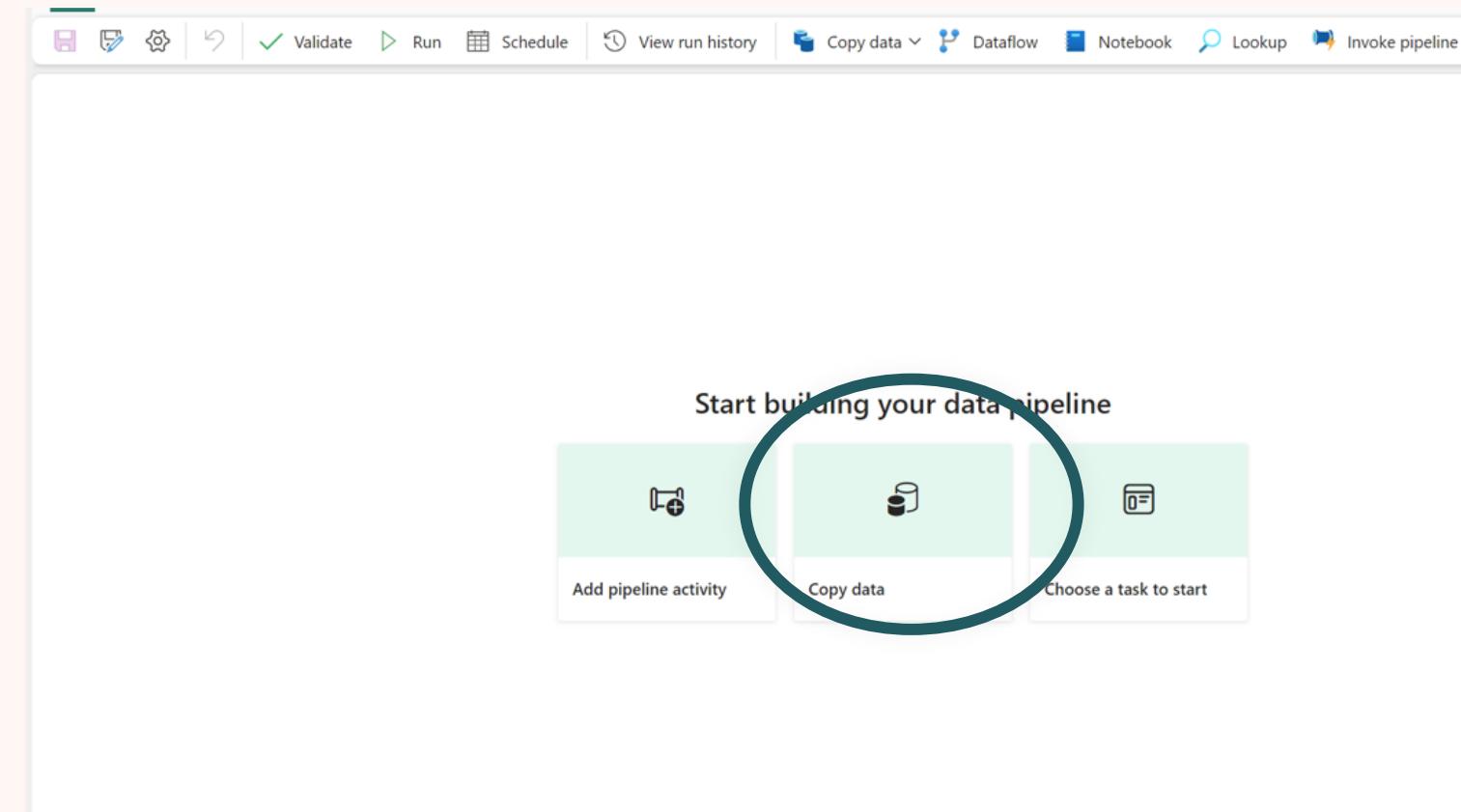
Lakehouse *
For Copy from SQL to LH (Copy destination)

Select...  



Understand Parameters

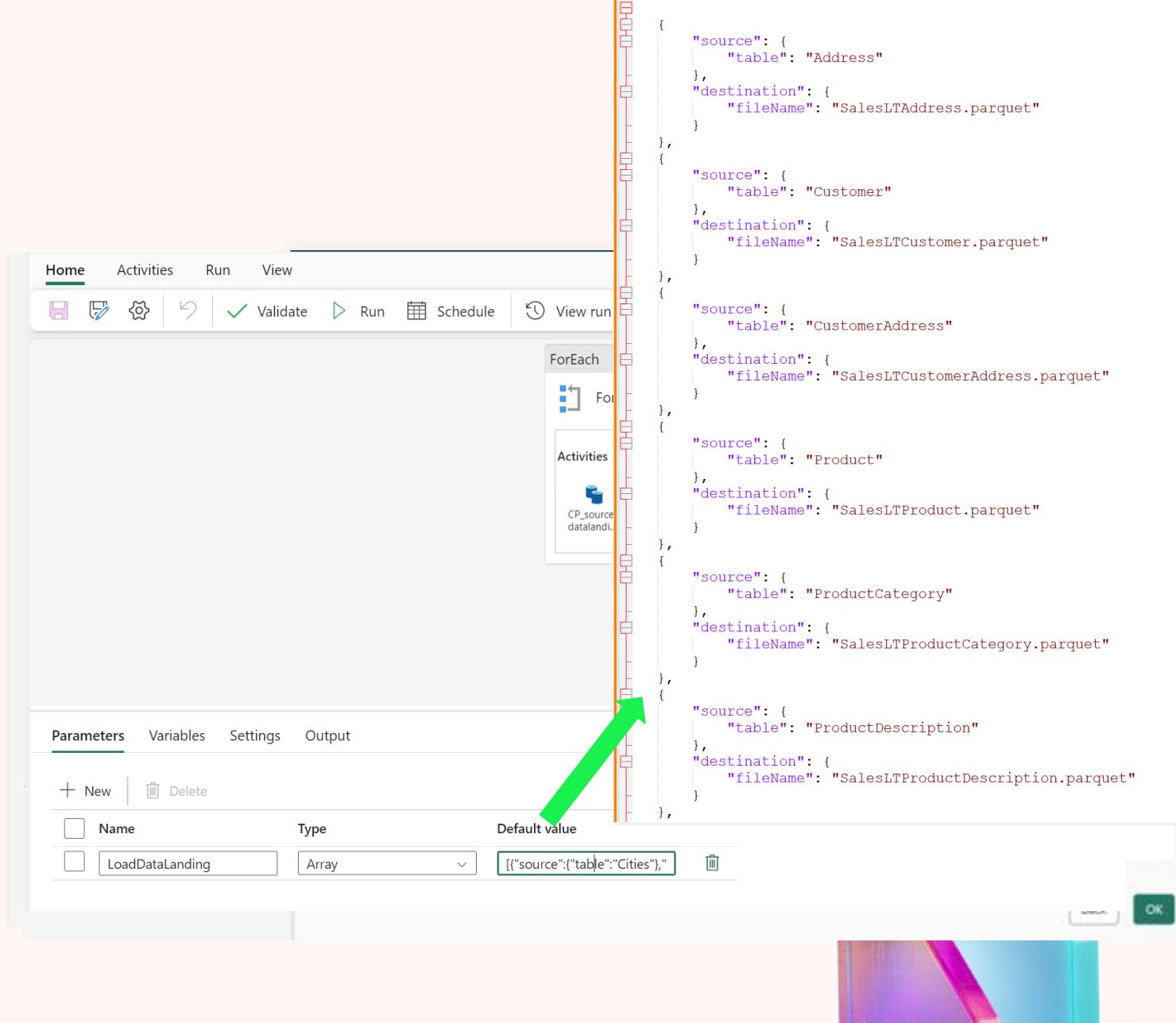
Copy Assistant



Understand Parameters

Copy Assistant

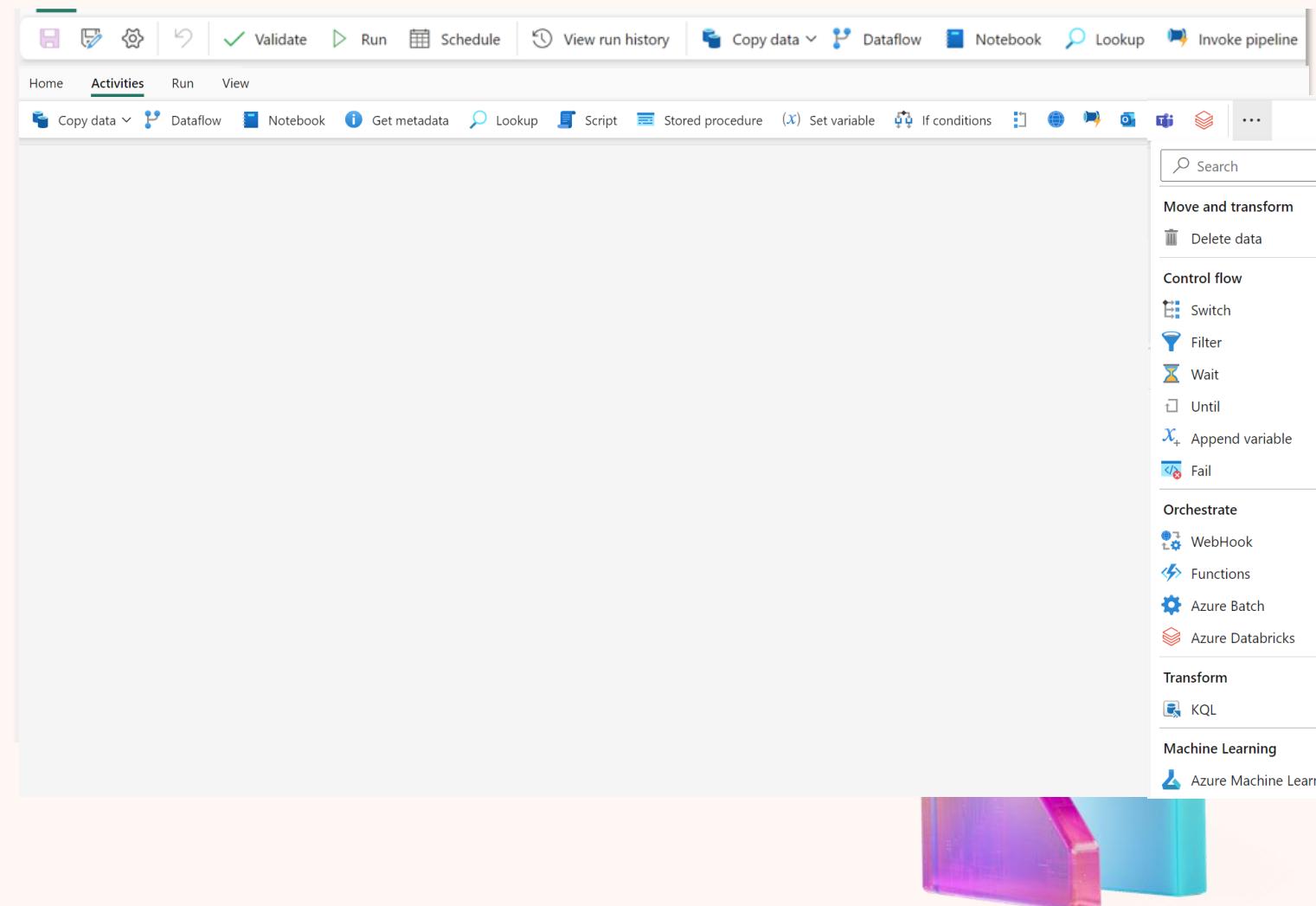
- Select Data Source
 - Select Tables
 - Select Destination
 - Select Lakehouse
 - Select Filetypes
 - Data is stored as Delta Parquet

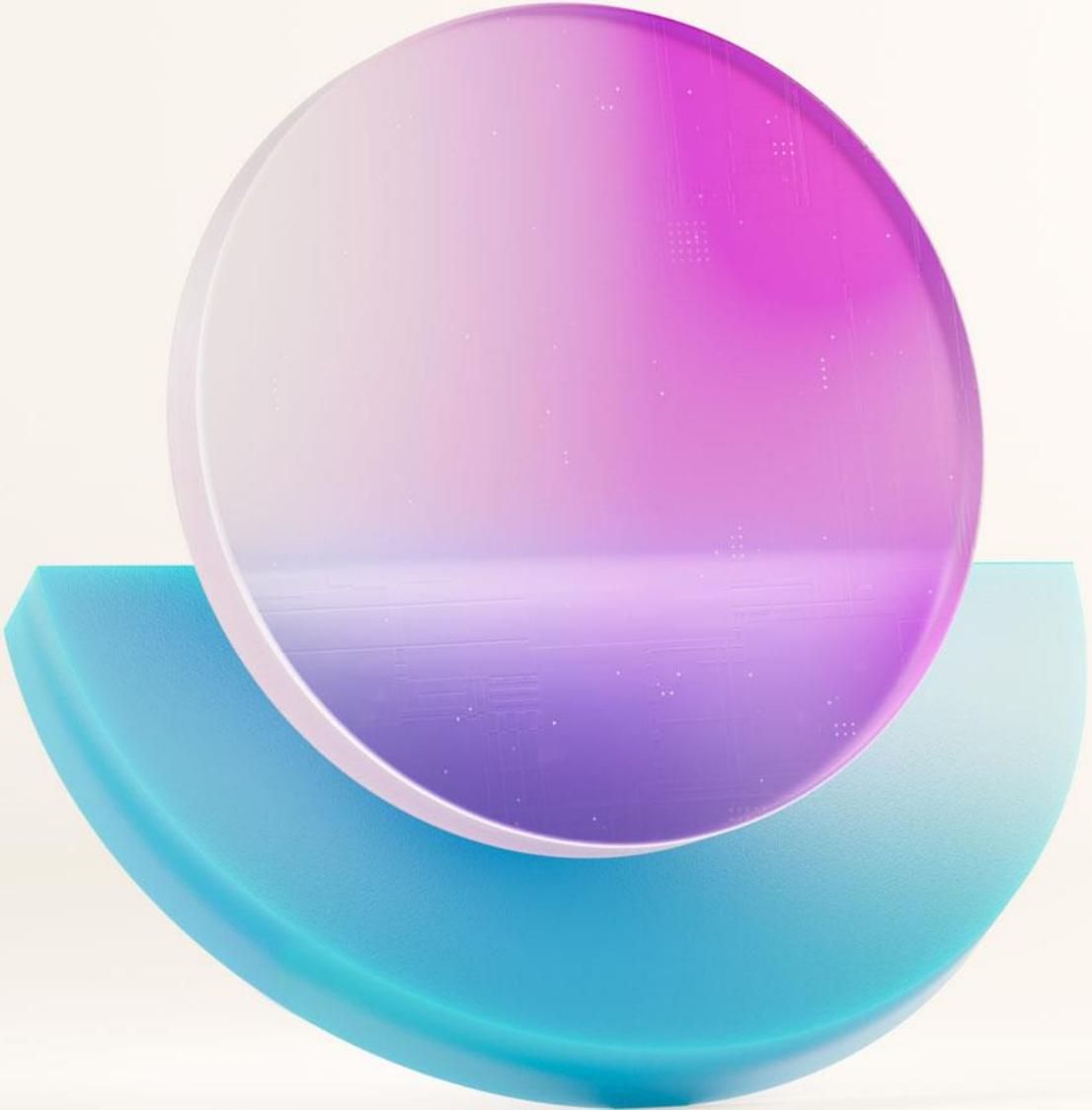


Understand Parameters

Blank Canvas

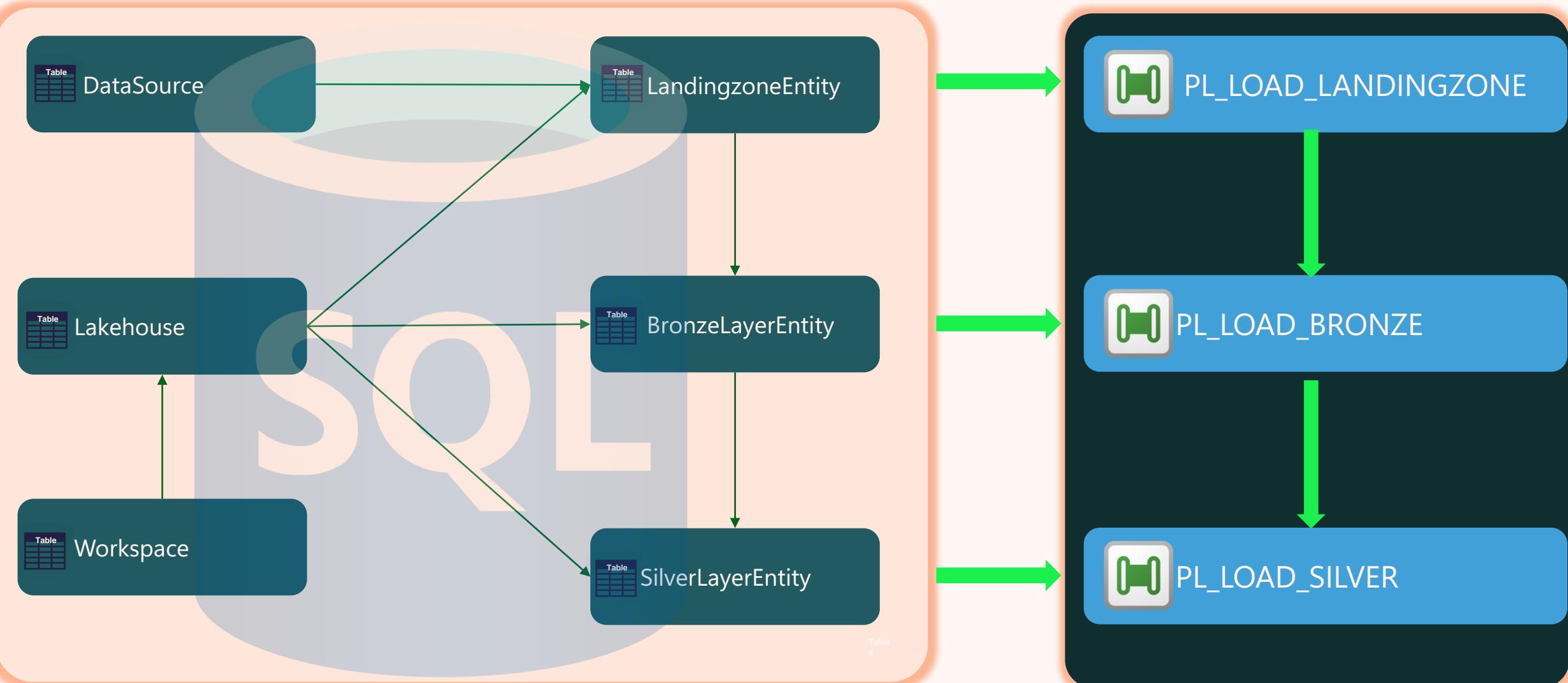
- Or just start from a blank Canvas



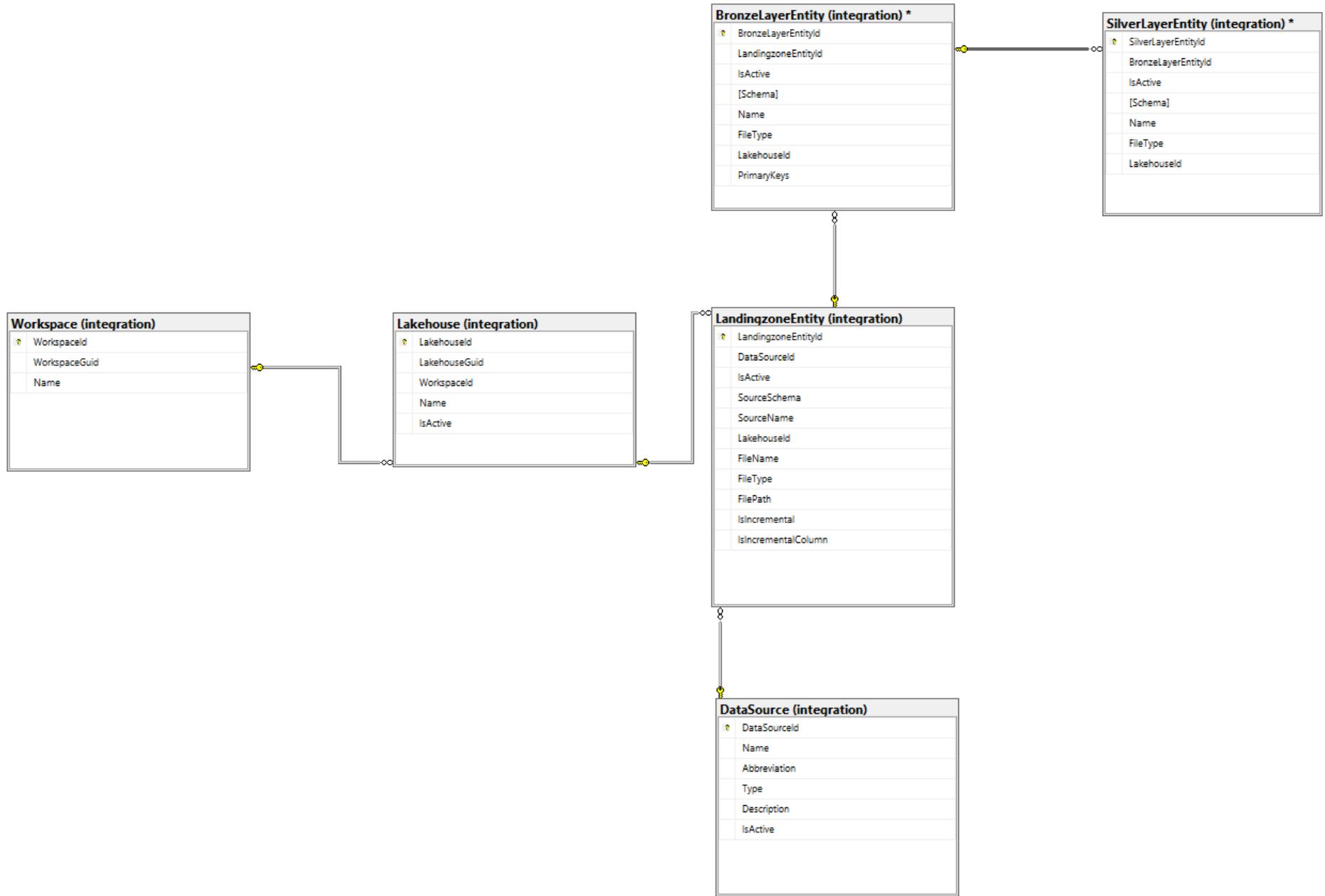


Framework in Azure SQL

Framework



Framework



Parameters

- Pipeline Parameters
- Notebook Parameters
- Copy Activity Parameters

IMPLEMENTING
DEFAULT PARAMETERS
THAT DEPEND ON
OTHER PARAMETERS

Define Pipeline Parameters

- Define Parameters
 - Pass through from Pipeline to Pipeline
 - Define Metadata
 - Versions
 - Key Vault
 - Secrets

Name	Type	Default value
location	String	West-Europe
key_vault_name	String	demokeyeuwdvilmvaultoxgnl
tenant_id	String	c183ff6a-8ca1-4185-9378-e
Environment	String	Development

Define Notebook Parameters

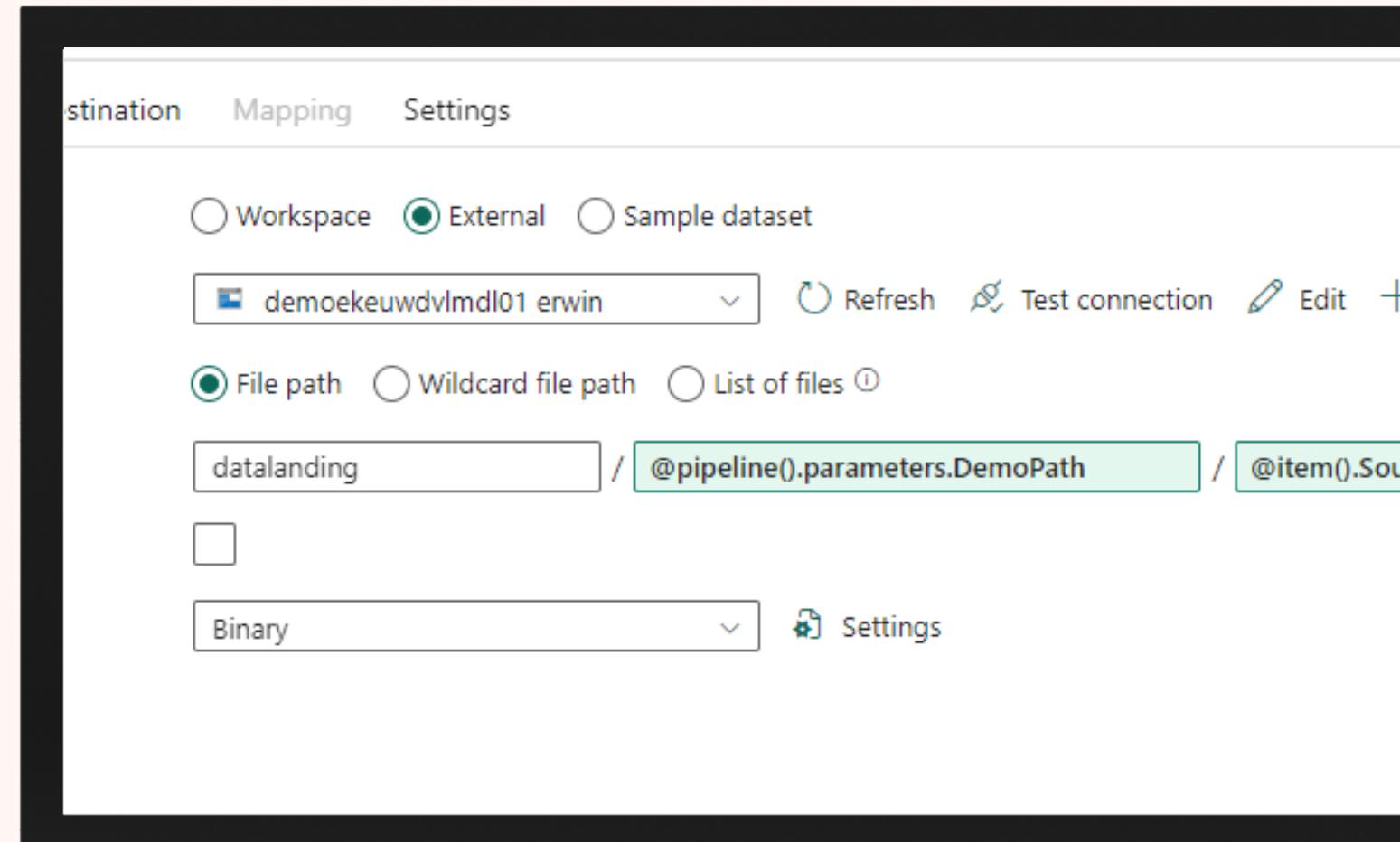
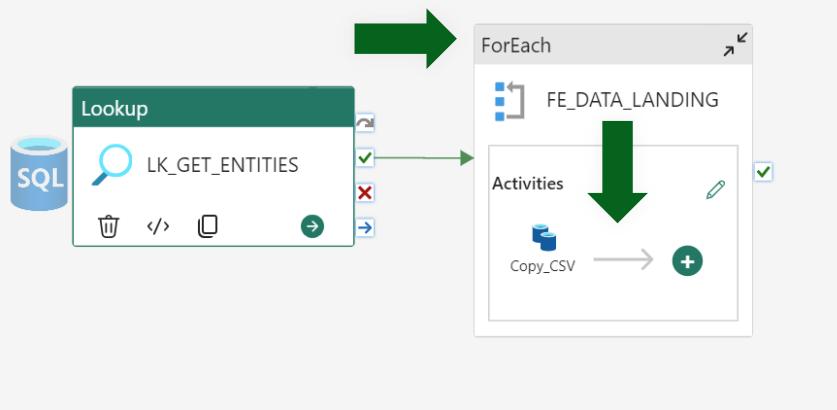
- Pass Parameters from Data Pipeline to Notebook
 - Toggle parameter cell

```
1 # Set arguments
2 PrimaryKeys = ""
3 IsIncremental = False
4
5 SourceWorkspace= ""
6 SourceLakehouse = ""
7 SourceLakehouseName = ""
8 source_file_path = ""
9 source_file_name = ""
10 source_file_type = ""
11
12 TargetWorkspace = ""
13 TargetLakehouse = ""
14 TargetLakehouseName = ""
15 target_schema = ""
16 target_name = ""
17
18
19 # # CSV
20 CompressionType = 'infer'
21 ColumnDelimiter = ','
22 RowDelimiter = '\n'
23 EscapeCharacter = ''
24 Encoding = 'UTF-8'
25 first_row_is_header = True
26 infer schema = True
```

PySpark (Python) ▾ Parameters

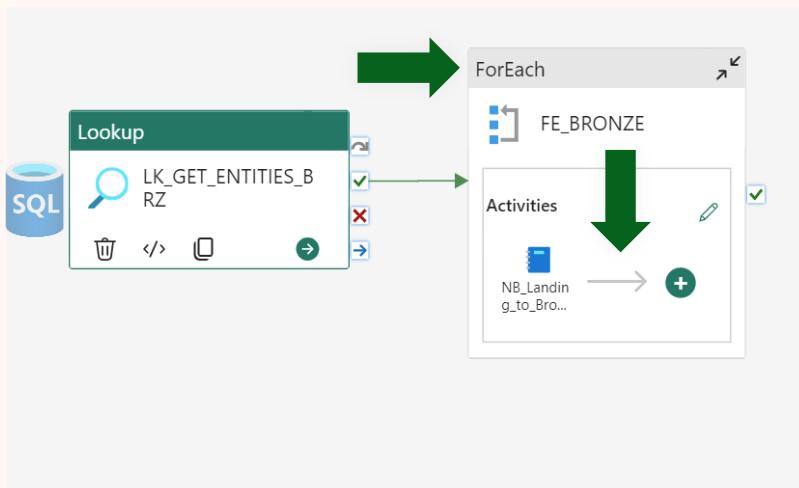
Copy Activity Parameters

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



Notebook Parameters

- Pass Parameters from Data Pipeline to Notebook
 - Set Base Parameters
 - Define Values



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

Microsoft Fabric
COMMUNITY CONFERENCE



DEMO

Fasten your seatbelts



'Optimize the notebook processing'

Fasten your seatbelts

- Notebook Activity
 - TooManyRequestForCapacity
 - Missing High Concurrency in Notebook Activity
 - Default 8 Vcores

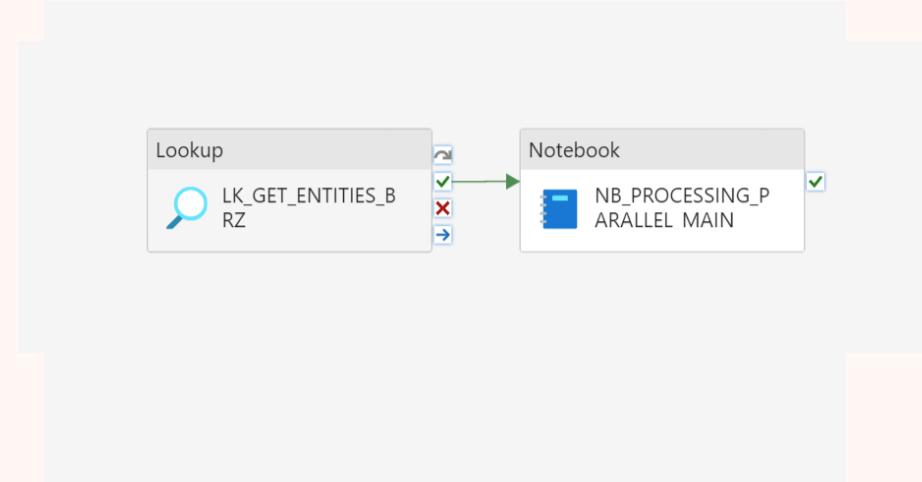
spark component or API have limits. To run this spark job, cancel an active Spark job through the Monitoring hub, choose a larger capacity SKU, or try again later. HTTP status code: 430 {Learn more} HTTP status code: 430.,"traceback":["Exception: Failed to create Livy session for executing notebook. Error:

[TooManyRequestsForCapacity] This spark job can't be run because you have hit a

Fasten your seatbelts

- Notebook Activity
 - TooManyRequestForCapacity
 - Missing High Concurrency in Notebook Activity
 - Default 8 Vcores
- `mssparkutils.notebook.runMultiple`
 - DAG (directed acyclic graph)
 - Notebook Name/Path
 - Arguments
 - Dependencies

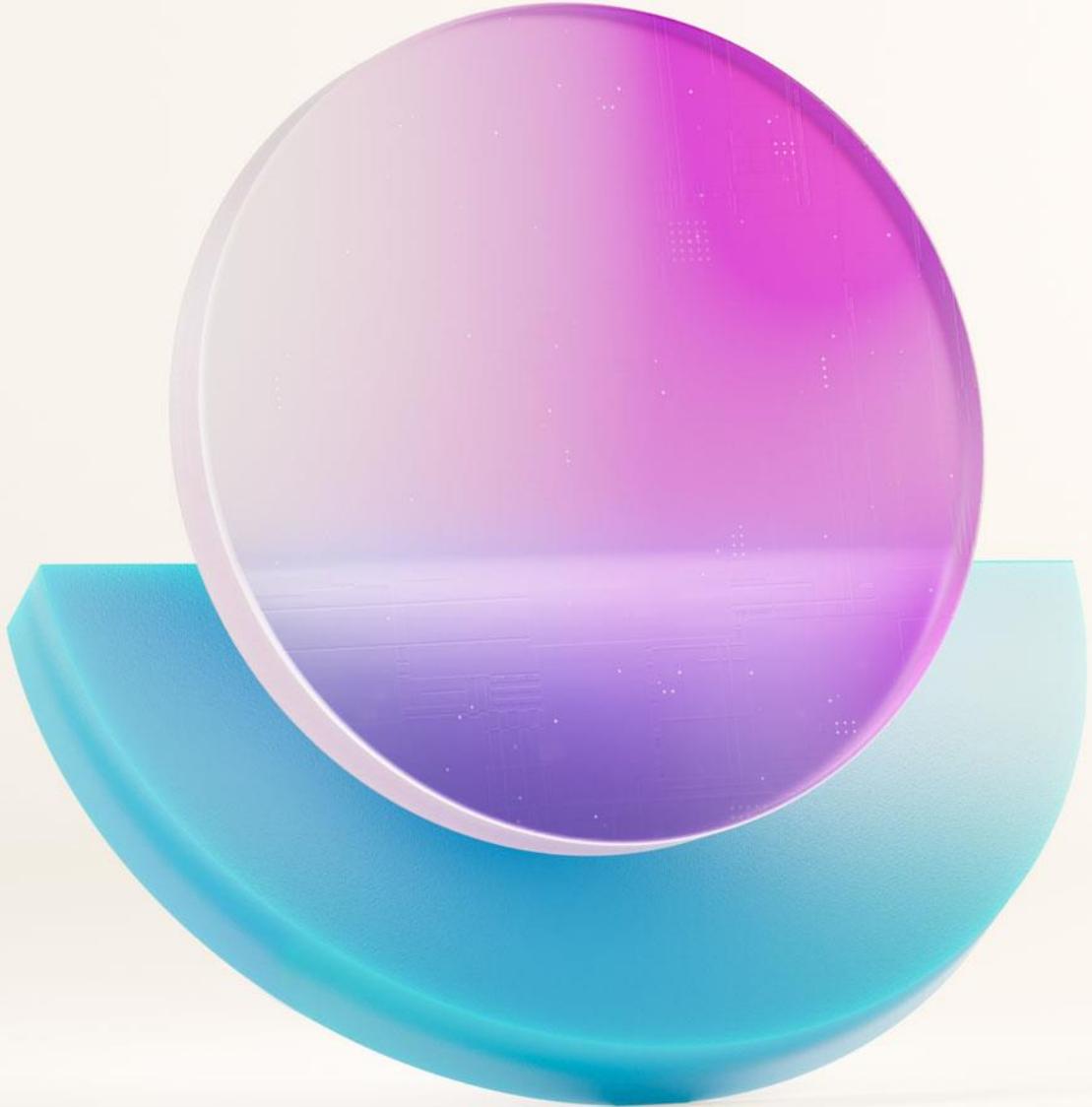
```
1  DAG = {  
2      "activities": [  
3          {  
4              "name": "NotebookSample", # activity name, must be unique  
5              "path": "NotebookSample", # notebook path  
6              "timeoutPerCellInSeconds": 90, # max timeout for each cell, default to 90 seconds  
7              "args": {"p1": "changed value", "p2": 666}, # notebook parameters  
8          },  
9          {  
10             "name": "NotebookSimple2",  
11             "path": "NotebookSimple2",  
12             "timeoutPerCellInSeconds": 120,  
13             "args": {"p1": "changed value 2", "p2": 777}  
14         },  
15         {  
16             "name": "NotebookSample2.2",  
17             "path": "NotebookSample2",  
18             "timeoutPerCellInSeconds": 120,  
19             "args": {"p1": "changed value 3", "p2": 888},  
20             "retry": 1,  
21             "retryIntervalInSeconds": 10,  
22             "dependencies": ["NotebookSample"] # list of activity names that this activity depends on  
23         }  
24     ]  
25 }  
26  
27 mssparkutils.notebook.runMultiple(DAG)
```



Microsoft Fabric
COMMUNITY CONFERENCE



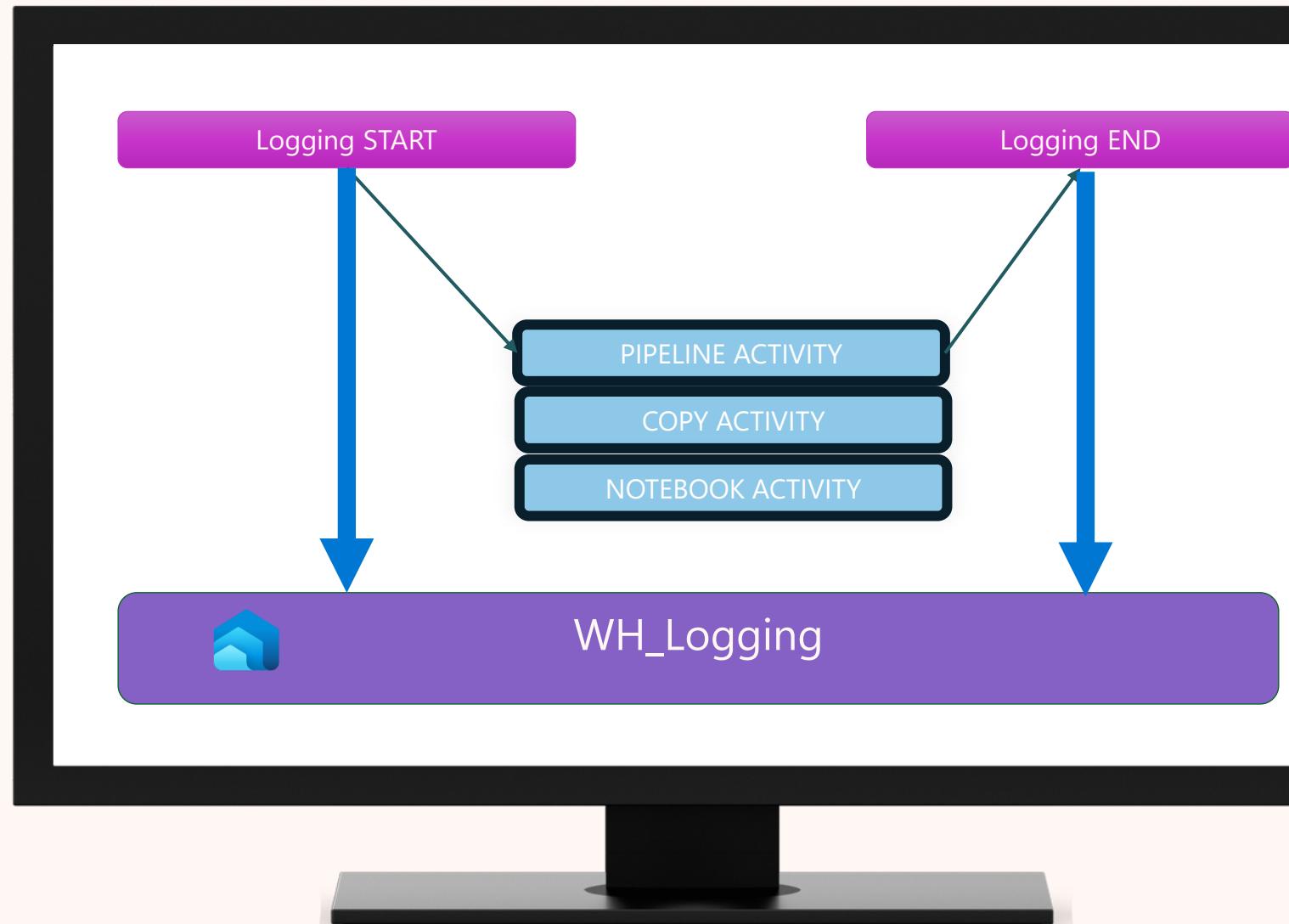
DEMO



Logging

Logging

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



Logging

- Add Information about pipelines
- Adding System Variables

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](#)

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 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

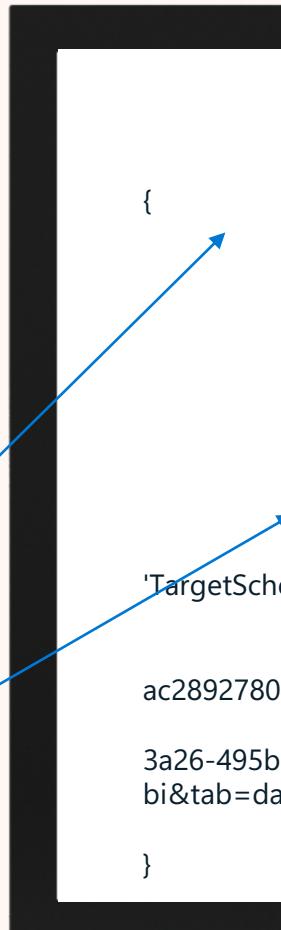
Logging

- 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}}  
}
```



```
{  
  "status": "Succeeded",  
  "result": {  
    "runId": "f0c69e6c-e28e-4db4-a8af-ac2892780xx2",  
    "runStatus": "Succeeded",  
    "sessionId": "fc55xx38-6fd8-47b0-863b-65c8c8db9878",  
    "sparkPool": "3xx61f99-edc7-4d6c-a866-f3bf70bc7235",  
    "error": null,  
    "lastCheckedOn": "2024-01-23T15:44:25.2733333Z",  
    "metadata": null,  
    "exitValue": "{CopyOutput}: {'Total Runtime': '0:00:30.493626',  
      'TargetSchema': 'Application', 'TargetName': 'People'}}"  
  },  
  "message": "Notebook execution is in Succeeded state, runId: f0c69e6c-e28e-4db4-a8af-ac2892780442",  
  "SparkMonitoringURL": "https://app.powerbi.com/workloads/de-ds/sparkmonitor/fecxxff4-3a26-495b-9ac9-475bbc59fae7/fc55be38-6fd8-47b0-863b-65c8c8db9878?trident=1&experience=powerbi&tab=data",  
  "executionDuration": 49  
}
```

Microsoft Fabric
COMMUNITY CONFERENCE



DEMO

Recap



Can we build Pipelines dynamically?



Can we extract data from my sources based on MetaData?



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



Can we build history from extracted data based on MetaData?

Challenges

- Out of the box fast and easy, less flexible
- Parameterize of connections with Azure Key Vault (Like ADF/Synapse with Linked Services)
- Need for High concurrency Notebook Activity in Pipeline (currently solved with the use of a Notebook executor or mssparkutils.notebook.runMultiple(DAG))
- Data gateway for Pipelines is coming
- Schedule can't be Parameterized like in ADF/Synapse
- Build in retry to Notebook Activity



Let's connect



Questions?





@erwindekreuk

linkedin.com/in/erwindekreuk

erwindekreuk.com

github.com/edkreuk

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

Microsoft Fabric
COMMUNITY CONFERENCE

Thank you

