

Unleashing the Potential of Metadata-Driven ELT Framework



Erwin de Kreuk

Principal Consultant - Lead Data & Analytics
InSpark



Erwin de Kreuk

Principal Consultant – Lead Data & Analytics
InSpark



Erwin de Kreuk

Principal Consultant - Lead Data & Analytics
InSpark

 @erwinedekreuk

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

 erwinedekreuk.com

 github.com/edkreuk

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



Let's
connect



We Are InSpark

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

THANK YOU



Platinum



smart
casual
datadesign



redgate TIME^XTENDER

Gold



b.telligent
smart data. smart decisions.

Lucient



Measure Killer



Silver



Bronze

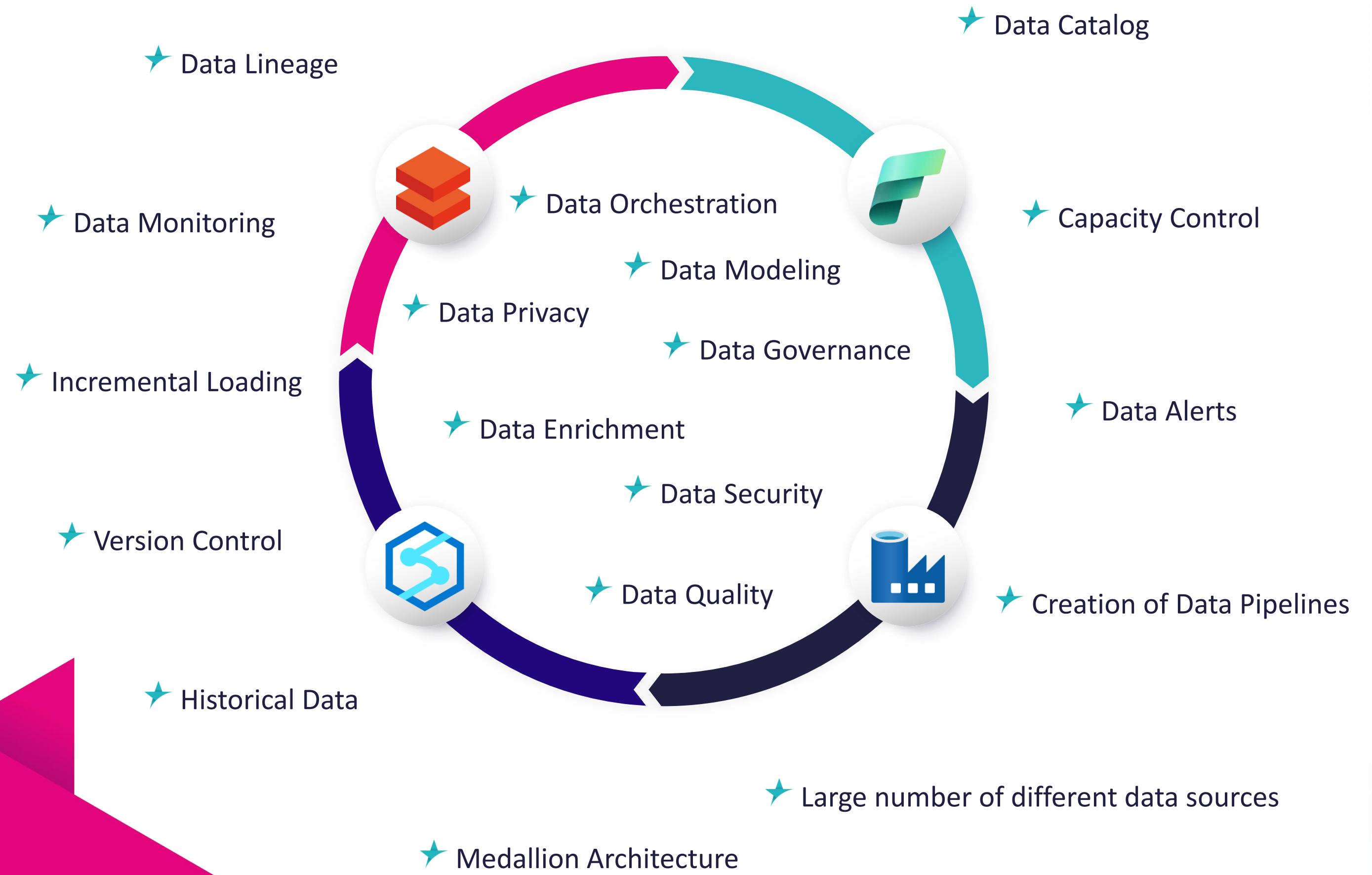


CUBIDO
Digital Solutions



Data platform Challenges

'From data source to data model' to report



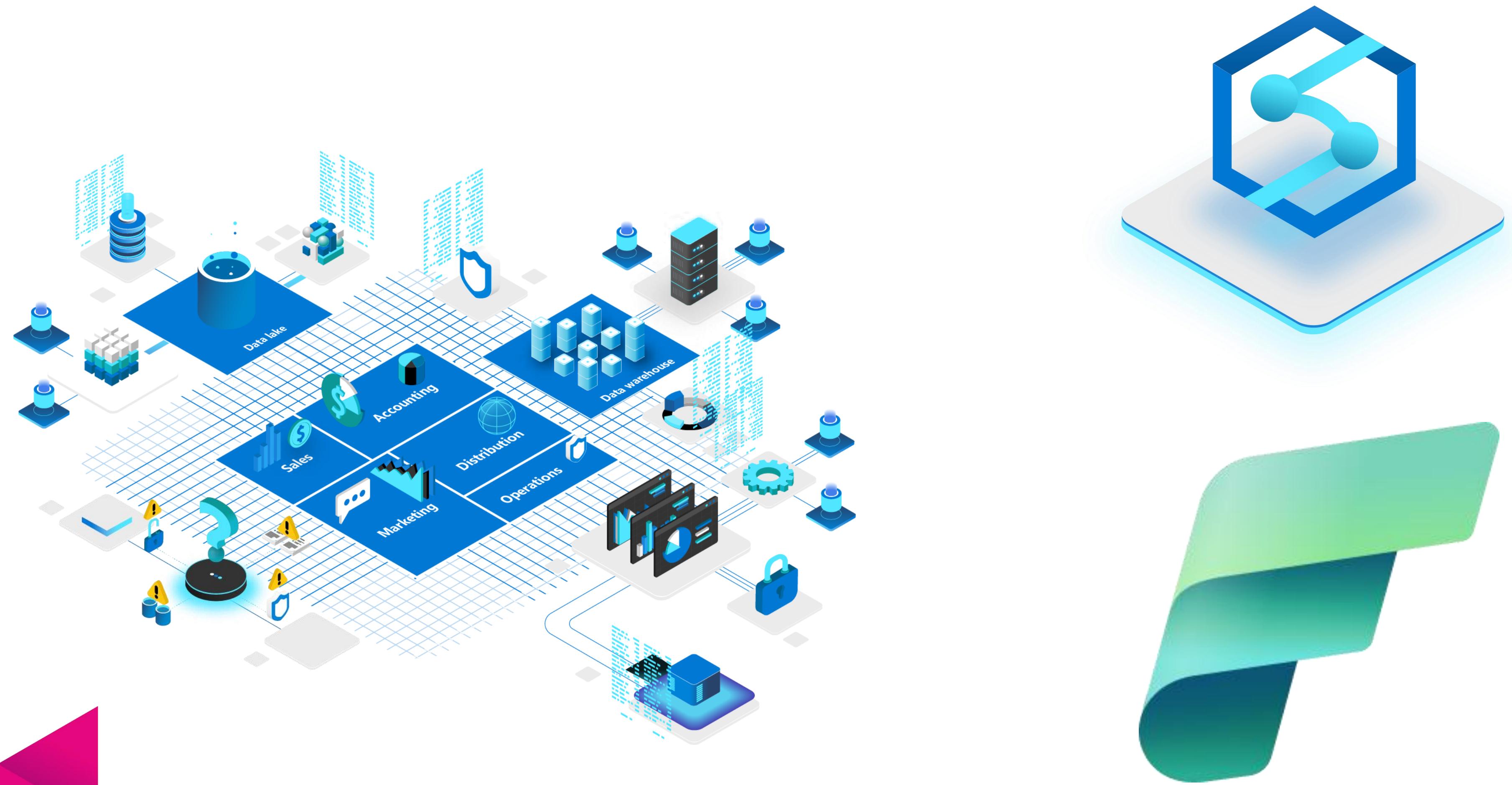
What would help

- ★ 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



Metadata driven ELT Framework

MSPAR

Objectives of today

- Why
- Out-of-the-box Framework in Synapse
- Custom-made Framework in Fabric
- Recap



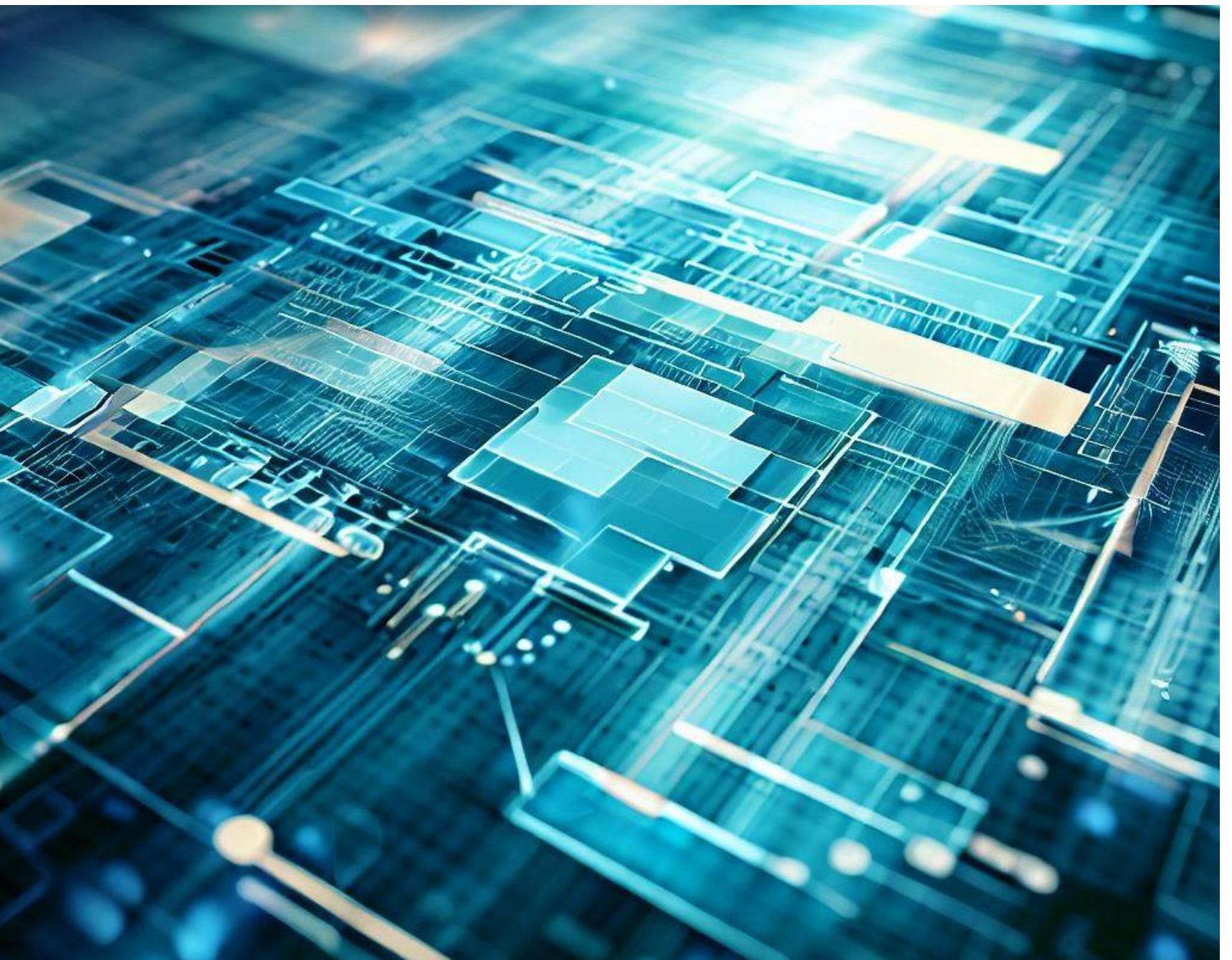
Why

- **Automation:** It allows for the automation of data pipeline creation and management, reducing manual effort. For example, data engineers can use metadata to generate code, schedule tasks, monitor performance, and handle errors.
- **Flexibility:** Enables easy modifications and adaptations to changes in data sources, formats, or business requirements. For example, data engineers can use metadata to update data pipelines without changing the underlying code, or to switch between different data platforms or tools.
- **Scalability:** Facilitates the scaling of data pipelines as data volumes or complexity increase. For example, data engineers can use metadata to distribute data processing across multiple nodes, or to optimize data storage and retrieval.
- **Traceability:** Enhances traceability and auditability of data as it moves through the pipeline. For example, data engineers can use metadata to track data provenance, lineage, and quality, or to provide documentation and reports.
- **Quality Assurance:** Improves data quality by enabling better monitoring and validation of data at each stage. For example, data engineers can use metadata to define data quality rules, metrics, and thresholds, or to perform data cleansing and transformation.

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

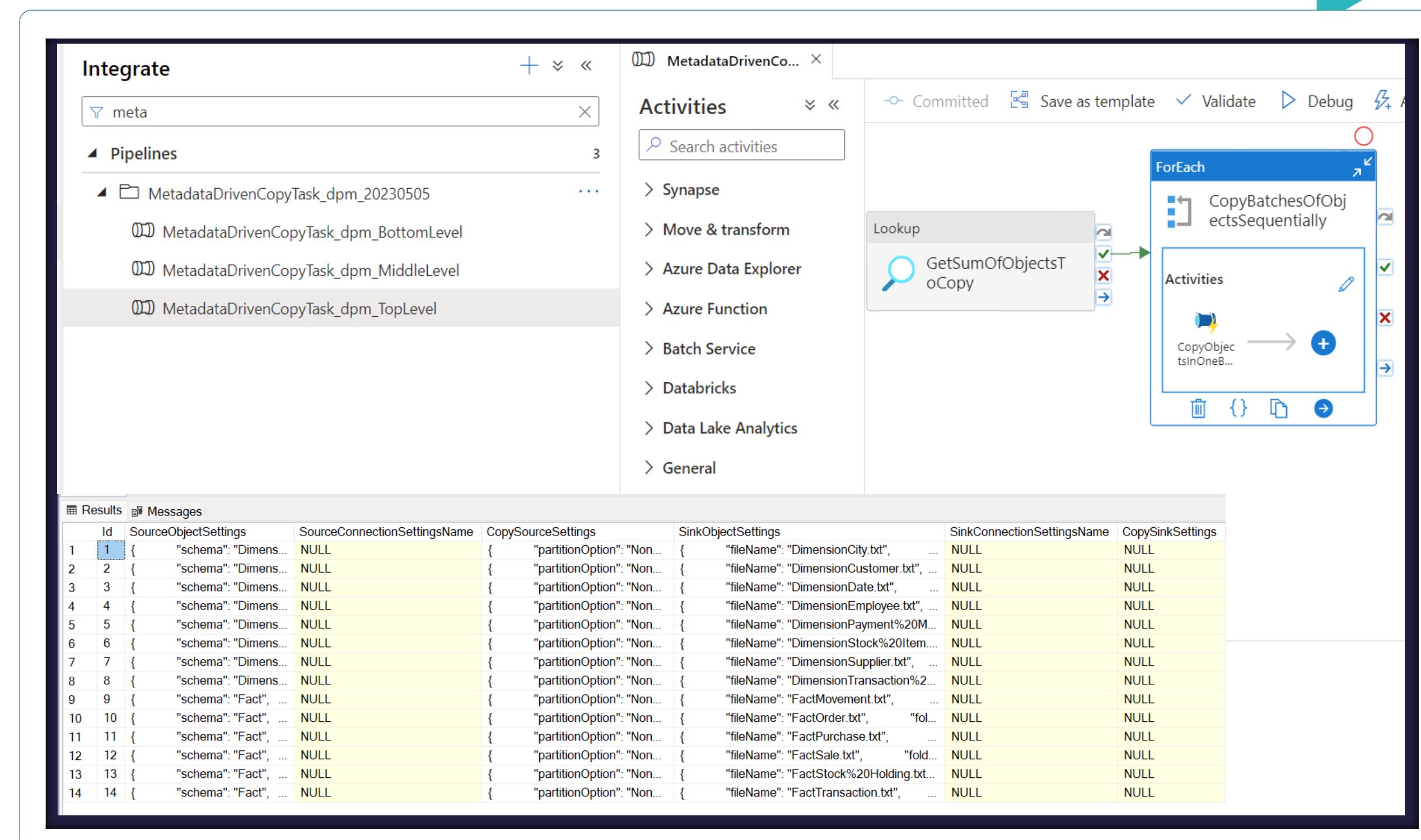
[Copy to clipboard](#)

```
***** 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 / Json /
- ★ Microsoft Fabric but also on Azure Synapse Analytics and Azure Data Factory
- ★ Based on the Medaillon Architecture



Microsoft Fabric



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

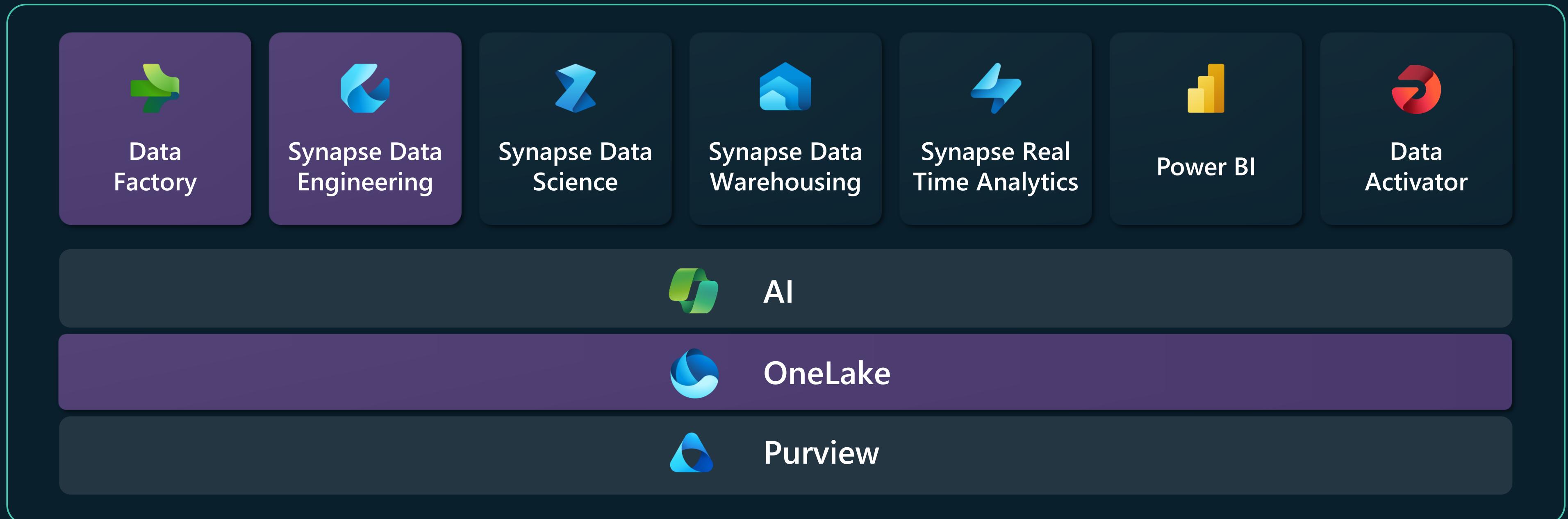
Unified
governance

Unified
business model



Microsoft Fabric

The unified data platform for the era of AI



Unified
architecture

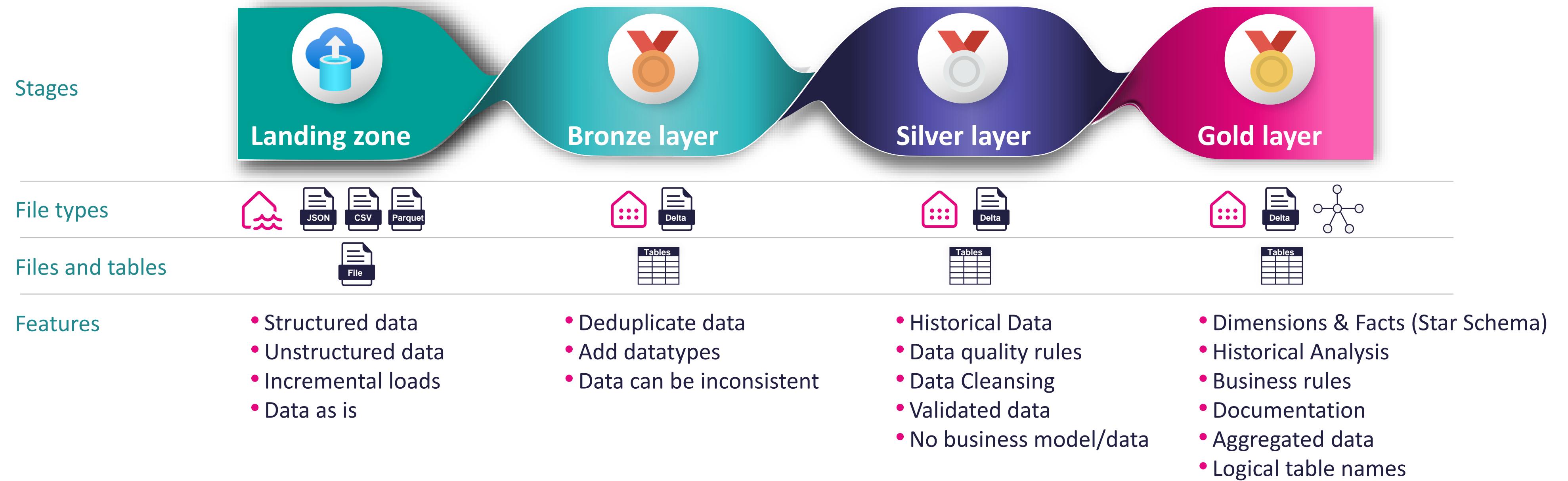
Unified
experience

Unified
governance

Unified
business model

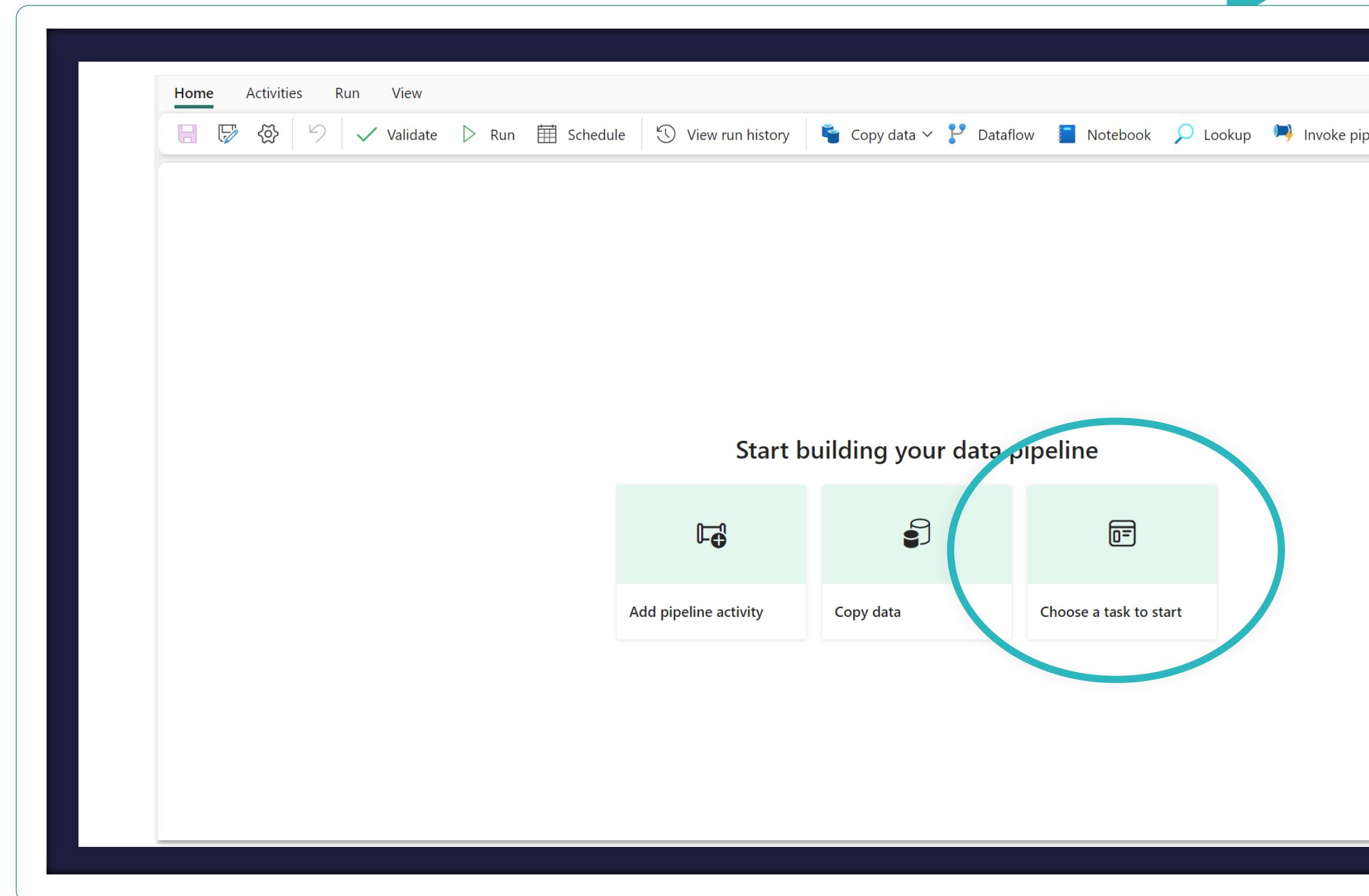
Medaillon Architecture

'Data processing in different stages'



Metadata-driven Framework

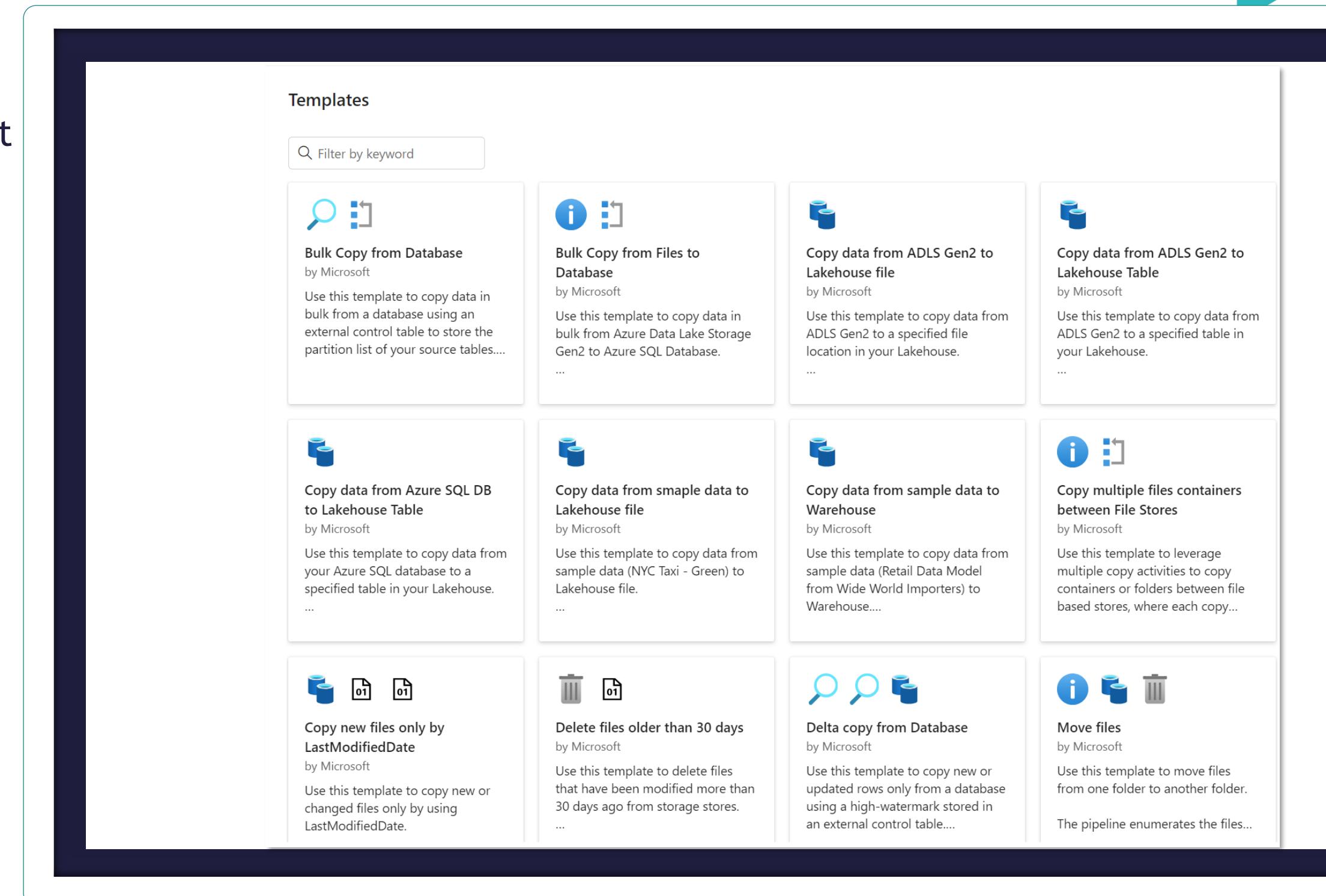
Understand Parameters



Metadata-driven Framework

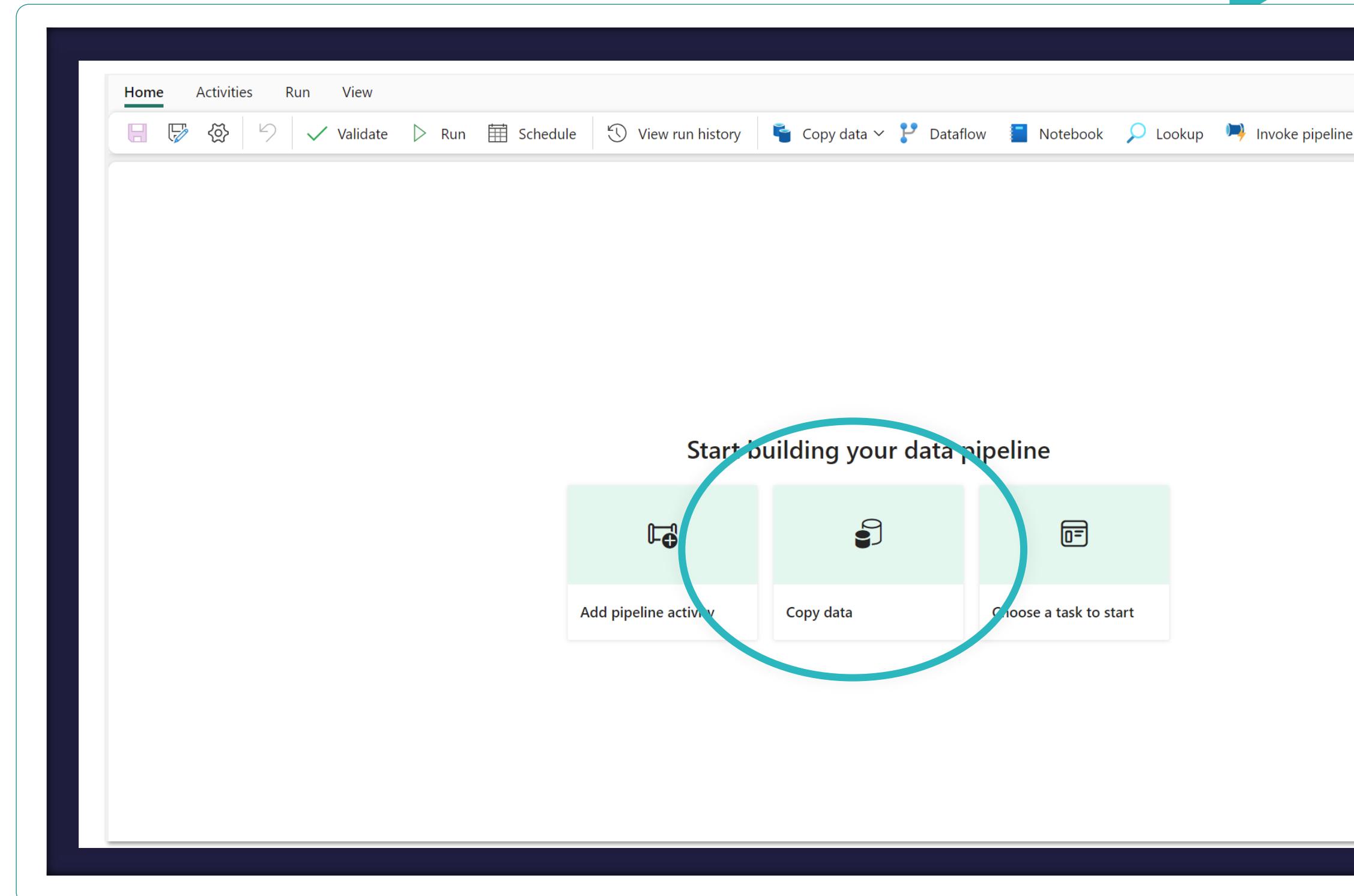
Templates

- ★ 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



Metadata-driven Framework

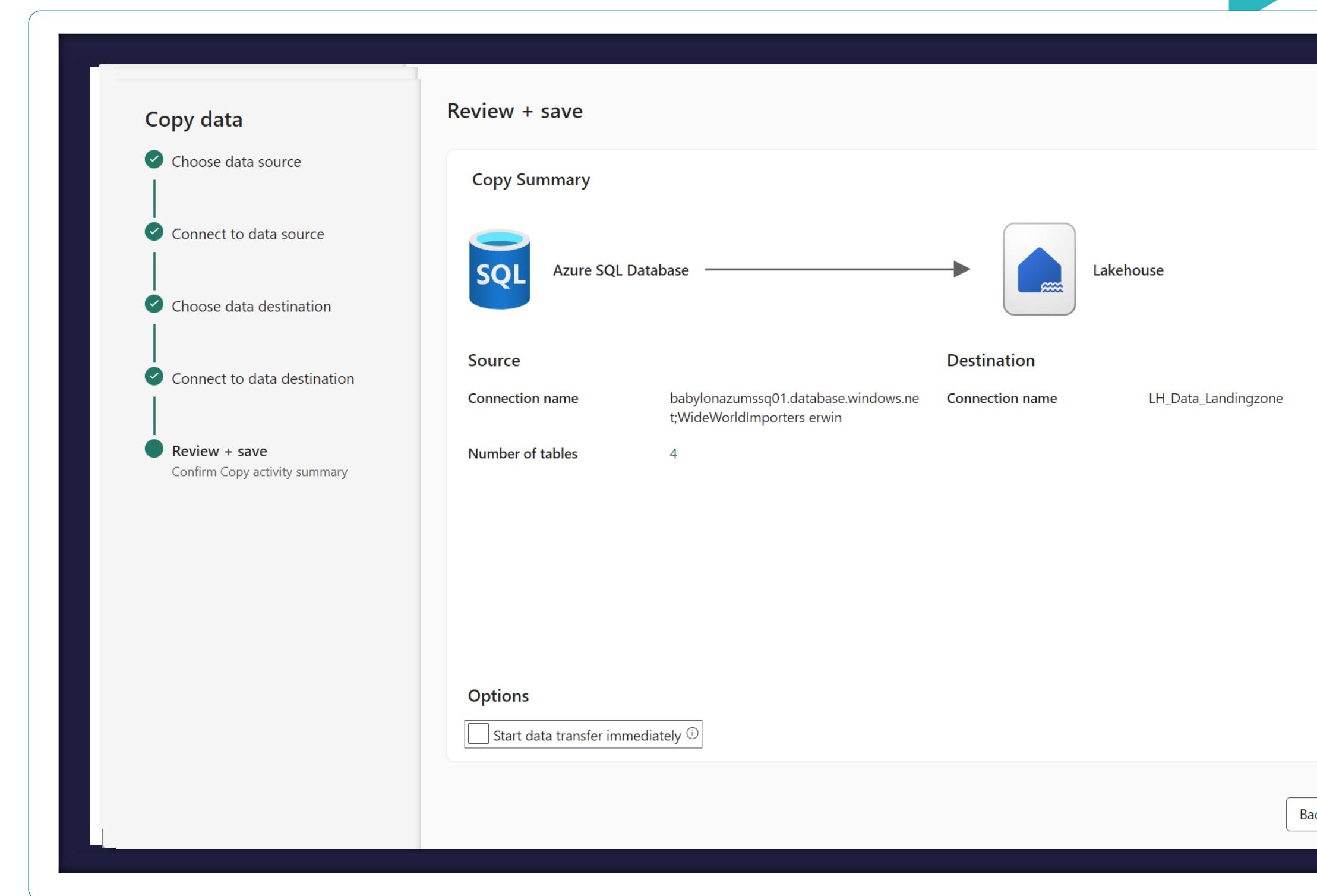
Understand Parameters



Metadata-driven Framework

Copy Assistant

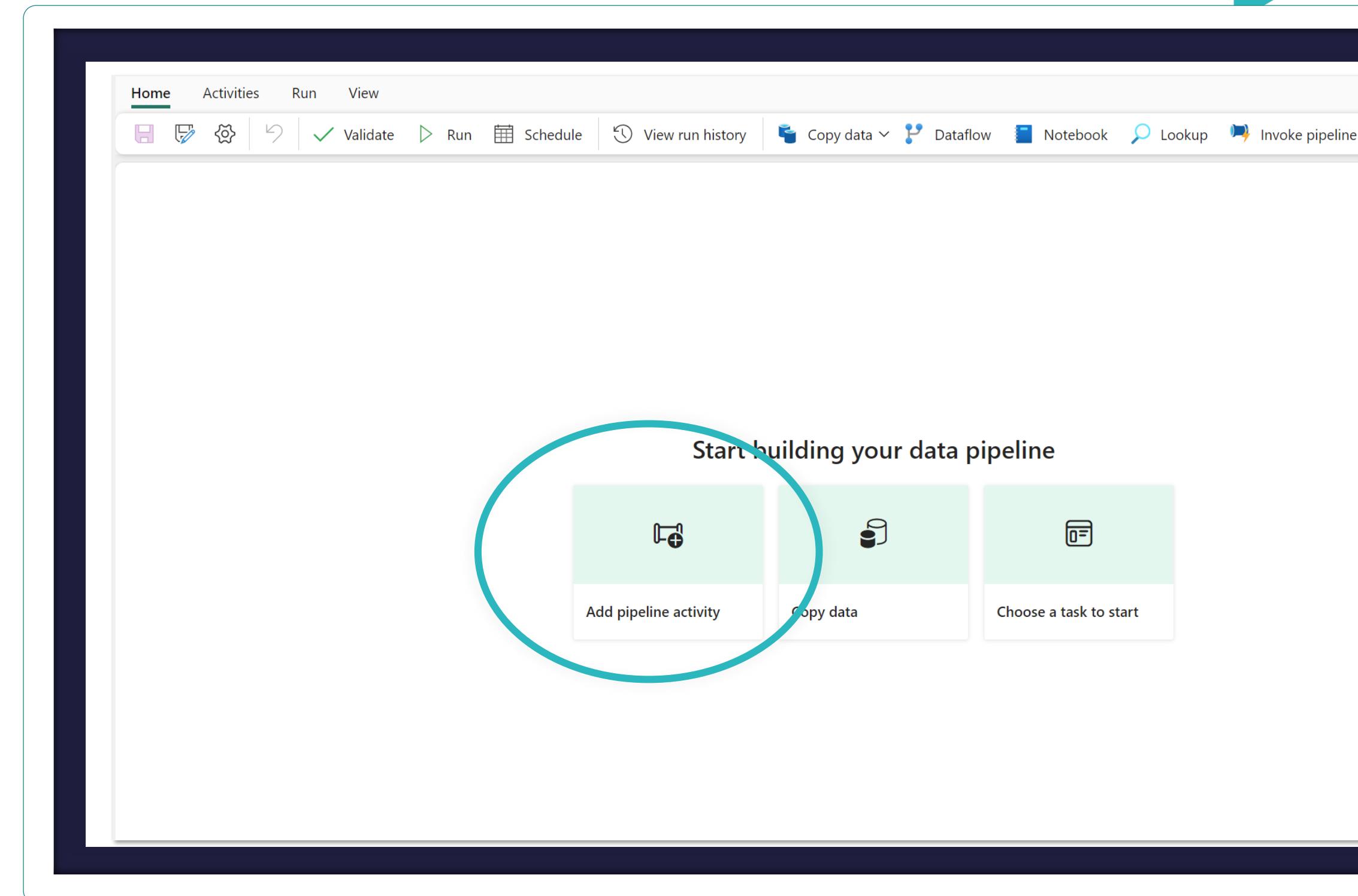
- ★ Select Data Source
- ★ Select Tables
- ★ Select Destination
- ★ Select Lakehouse
- ★ Select Filetypes
- ★ Data is stored as Delta Parquet



Metadata-driven Framework

Understand Parameters

★ Or just start from blank canvas



Innovate to accelerate

DEMO

Custom-Made Framework

- ★ Pipeline Parameters
- ★ Notebook Parameters

IMPLEMENTING
DEFAULT PARAMETERS
THAT DEPEND ON
OTHER PARAMETERS

Pipeline Parameters

- ★ Define Parameters
 - ★ Pass through from Pipeline to Pipeline
 - ★ Define Metadata
 - ★ Versions

Parameters			
	Name	Type	Default value
<input type="checkbox"/>	LoadDataLanding	Array	[{"source": {}}
<input type="checkbox"/>	LoadBronze	Array	[{"source": {"source_file_path": "bronze"}, "target": {"target_file_path": "bronze"}, "transform": {"script": "bronze"}, "type": "transformed"}]
<input type="checkbox"/>	LoadSilver	Array	[{"source": {"source_schema": "silver"}, "target": {"target_file_path": "silver"}, "transform": {"script": "silver"}, "type": "transformed"}]

Metadata-driven Framework

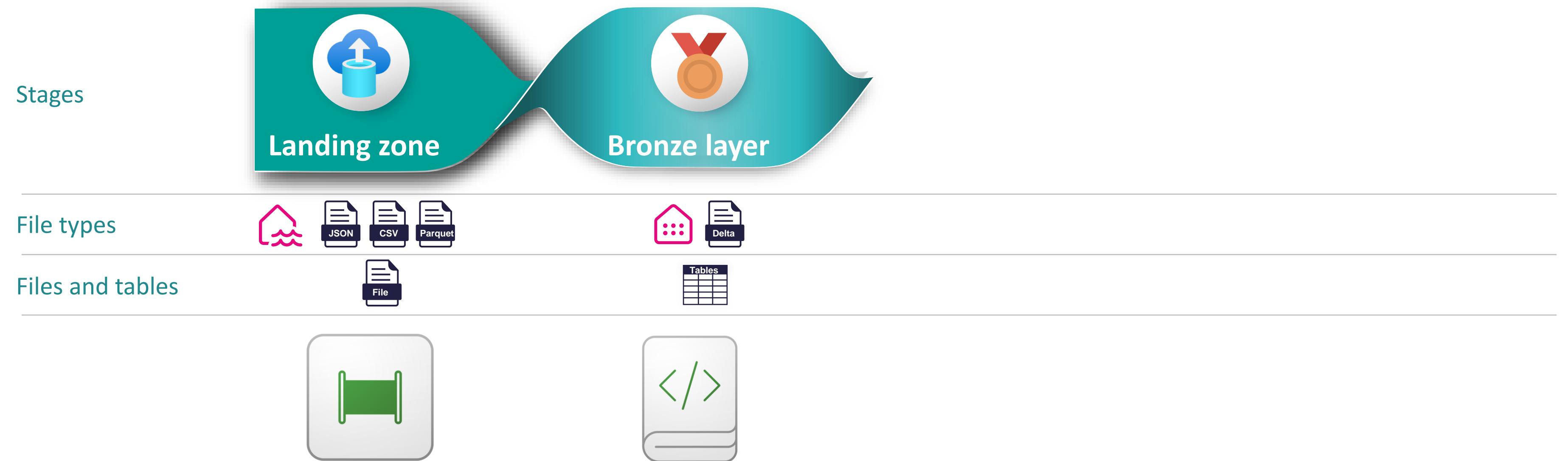
Lakehouse

- ★ 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.



Medaillon Architecture

'Data processing in different stages'



Metadata-driven Framework

Notebook Parameters

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



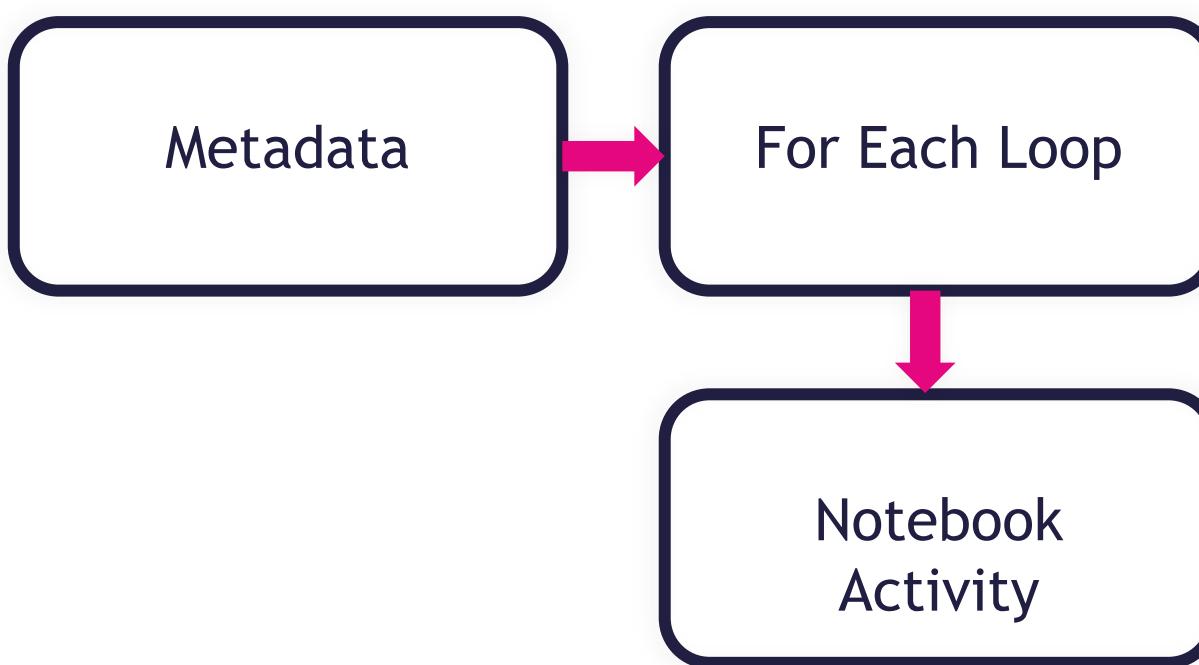
```
1 # Set arguments
2 PrimaryKeys = ""
3 IsIncremental = False
4
5 SourceWorkspace= ""
6 SourceLakehouse =""
7 SourceLakehouseName ='LH_Data_Landingzone'
8 SourceFilePath = ""
9 SourceFileName = ""
10
11 TargetWorkspace= ""
12 TargetLakehouse =""
13 TargetLakehouseName ='LH_Bronze_Layer'
14 TargetSchema = ""
15 TargetName = ""
16
17
```

The screenshot shows a Jupyter Notebook interface. A code cell containing Python code is selected. A context menu is open to the right of the cell, with the 'Toggle parameter cell' option highlighted. The status bar at the bottom indicates 'PySpark (Python) Parameters'. The code cell contains variables like PrimaryKeys, IsIncremental, SourceWorkspace, SourceLakehouse, SourceLakehouseName, SourceFilePath, SourceFileName, TargetWorkspace, TargetLakehouse, TargetLakehouseName, TargetSchema, and TargetName.

Metadata-driven Framework

Notebook Parameters

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



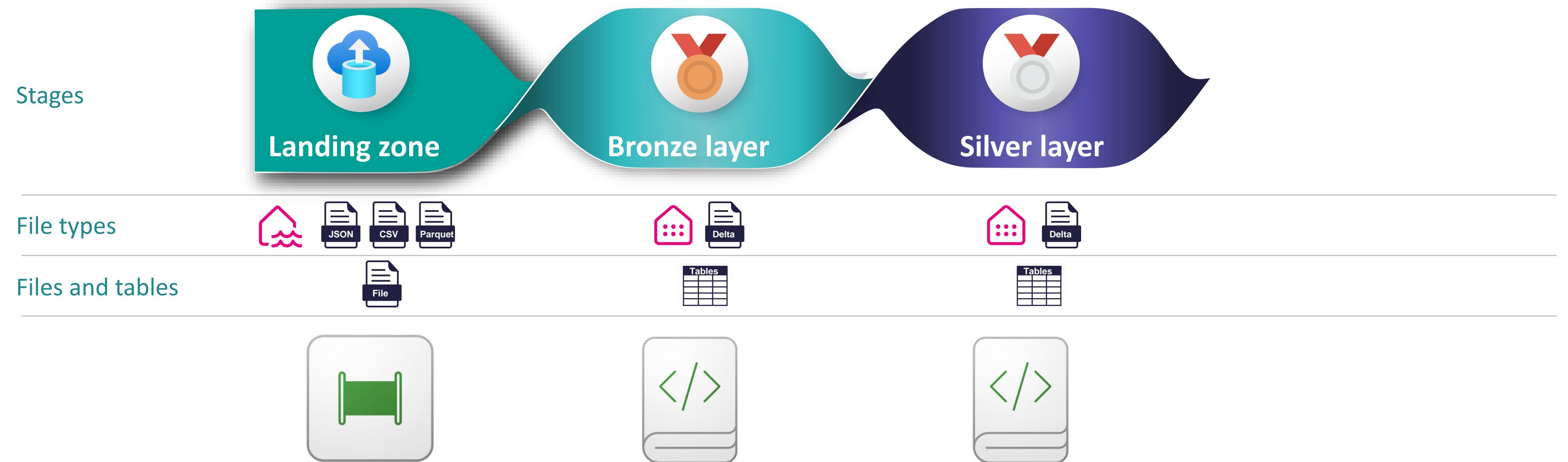
Base parameters

Name	Type	Value
SourceLakehouse	String	@item().source.SourceLakehouse
source_file_path	String	@item().source.source_file_path
source_file_name	String	@item().source.source_file_name
PrimaryKeys	String	@item().target.PrimaryKeys
TargetLakehouse	String	@item().target.TargetLakehouse
target_schema	String	@item().target.target_schema
target_name	String	@item().target.target_name

```
{
  "source": {
    "source_file_path": "WideWorldImporters",
    "source_file_name": "ApplicationPeople.parquet",
    "SourceLakehouse": "xxxxxxxxxxxxxxxxxxxx"
  },
  "target": {
    "target_schema": "Application",
    "target_name": "People",
    "PrimaryKeys": "PersonID",
    "targetLakehouse": "xxxxxxxxxxxxxxxxxxxx"
  }
},
```

Medaillon Architecture

'Data processing in different stages'



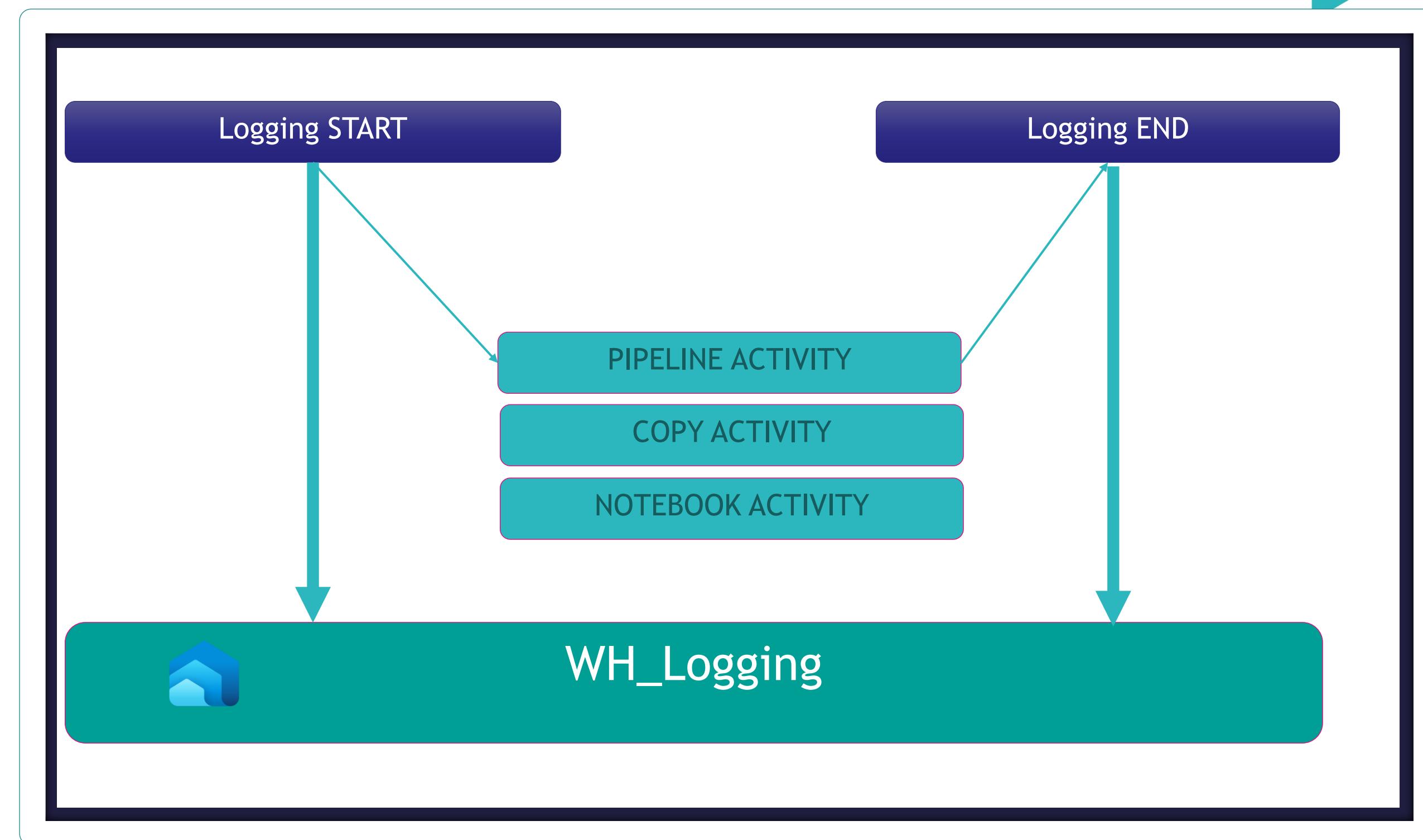
Innovate to accelerate

DEMO

Custom-Made Framework

Logging

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



Logging

- ★ Add Information about pipelines
- ★ Adding System Variables

The screenshot shows the 'General' tab of the Pipeline Settings page. Under 'Data store type', 'Workspace' is selected. The 'Warehouse' dropdown is set to 'WH_Logging'. The 'Stored procedure name' is '[logging].[sp_AuditPipeline]'. Below this, the 'Stored procedure parameters' section is expanded, showing a table with 11 rows:

Name	Type	Value
LogData	String	{ "LogType": "PipelineRun", "LogTime": "2023-10-01T12:00:00Z", "LogMessage": "Pipeline run started" }
LogType	String	Start
PipelineGuid	Guid	@{pipeline().Id}
PipelineName	String	@{pipeline().Name}
PipelineParameters	String	Value
PipelineParentRunGuid	Guid	@{pipeline().RunId}
PipelineRunGuid	Guid	@{pipeline().RunId}
TriggerGuid	Guid	@{pipeline().TriggerId}
TriggerTime	DateTime	@{pipeline().TriggerTime}
TriggerType	String	Manual
WorkspaceGuid	Guid	@{pipeline().WorkspaceId}

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  
}
```

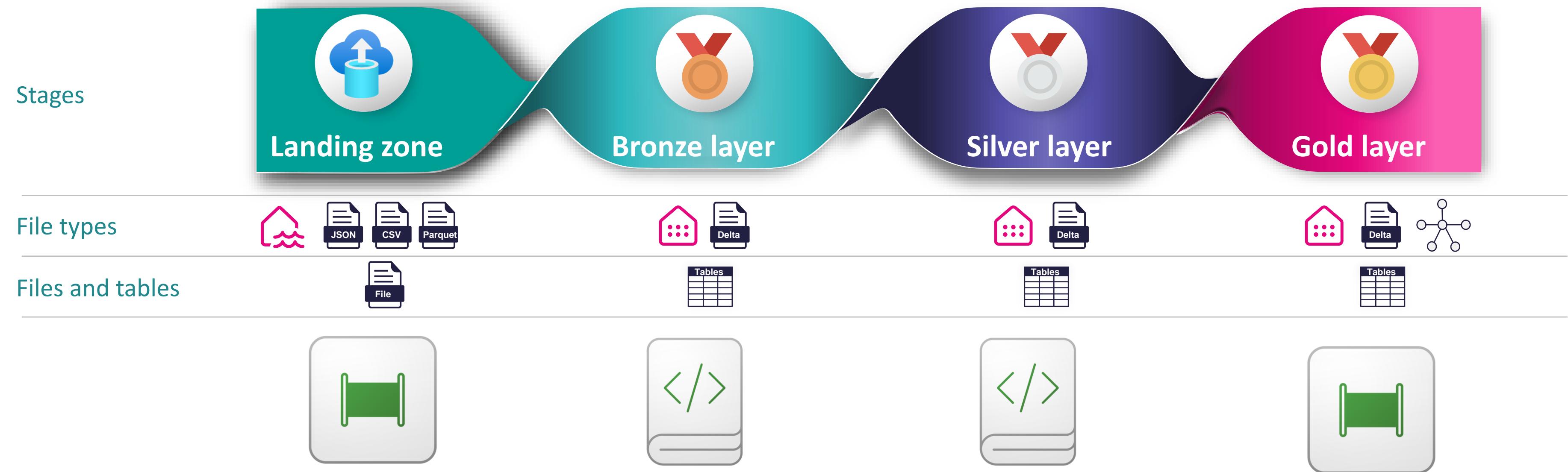
```
{  
    "status": "Succeeded",  
    "result": {  
        "runId": "f0c69e6c-e28e-4db4-a8af-ac2892780xx2",  
        "runStatus": "Succeeded",  
        "sessionId": "fc55xx38-6fd8-47b0-863b-65c8c8db9878",  
        "sparkPool": "3xx61f99-edc7-4d6c-a866-f3bf70bc7235",  
        "error": null,  
        "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=power-  
    bi&tab=data",  
    "executionDuration": 49  
}
```

Innovate to accelerate

DEMO

Medaillon Architecture

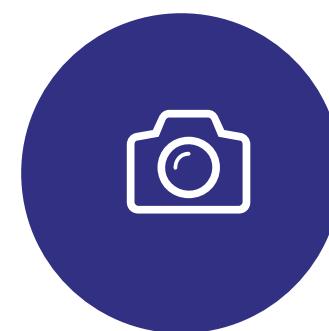
'Data processing in different stages'



Recap



Can we build 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?



Recap

- ★ 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
- ★ Views can't not be called from a Notebook (Silver to Gold)
- ★ Schedule can't be Parameterized like in ADF/Synapse
- ★ Build in retry to Notebook Activity

Questions?



THANK YOU



Platinum



smart
casual
datadesign



redgate TIME^XTENDER

Gold



b.telligent
smart data. smart decisions.

Lucient



Measure Killer



Silver



Bronze



CUBIDO
Digital Solutions



Thank you

Evaluations, evaluations...



https://evals.datagrillen.com/evals_vienna.aspx



Erwin de Kreuk

Principal Consultant - Lead Data & Analytics
InSpark

@erwinedekreuk

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

erwinedekreuk.com

github.com/edkreuk

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



Let's
connect

