ETL in Azure Made Easy

with Data Factory Data Flows





Paul Andrew | Principal Consultant & Solution Architect











Paul Andrew | Principal Consultant & Solution Architect









GitHub



https://github.com/mrpaulandrew

CommunityEvents

Demo code, content and slides from various community events.

C++

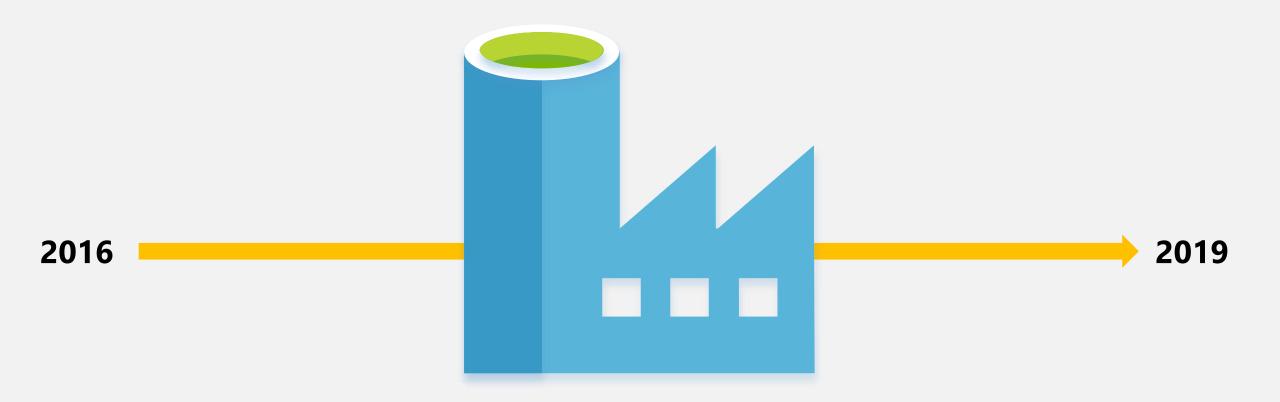
{Event/Location}-{Month}-{Year}



Azure Data Factory

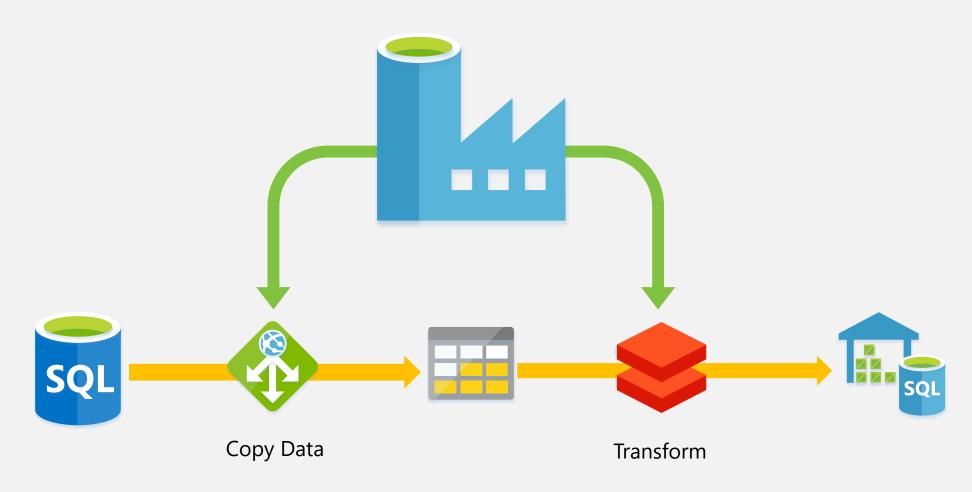


What is Azure Data Factory?



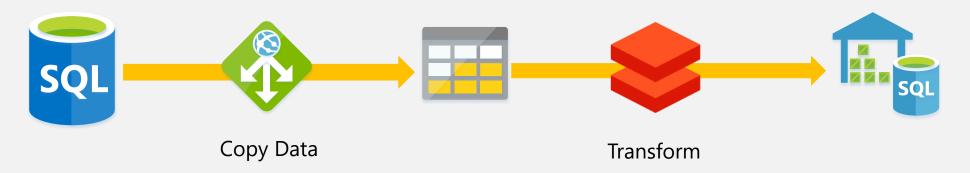


What is Azure Data Factory?

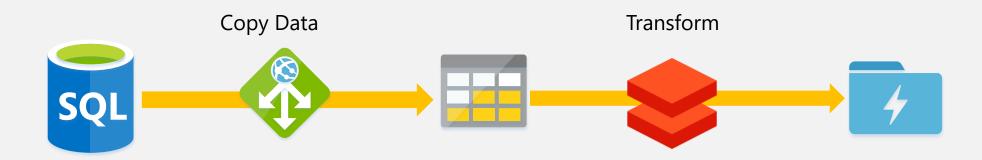




What is Azure Data Factory?



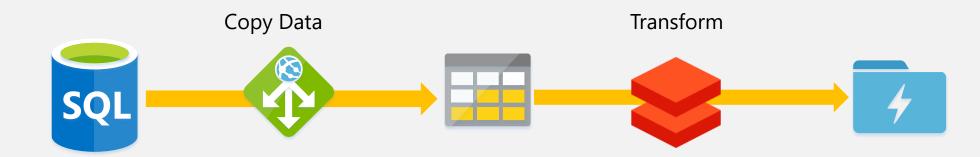




1 Linked Services – How do I connect? Like the SSIS connection manager.

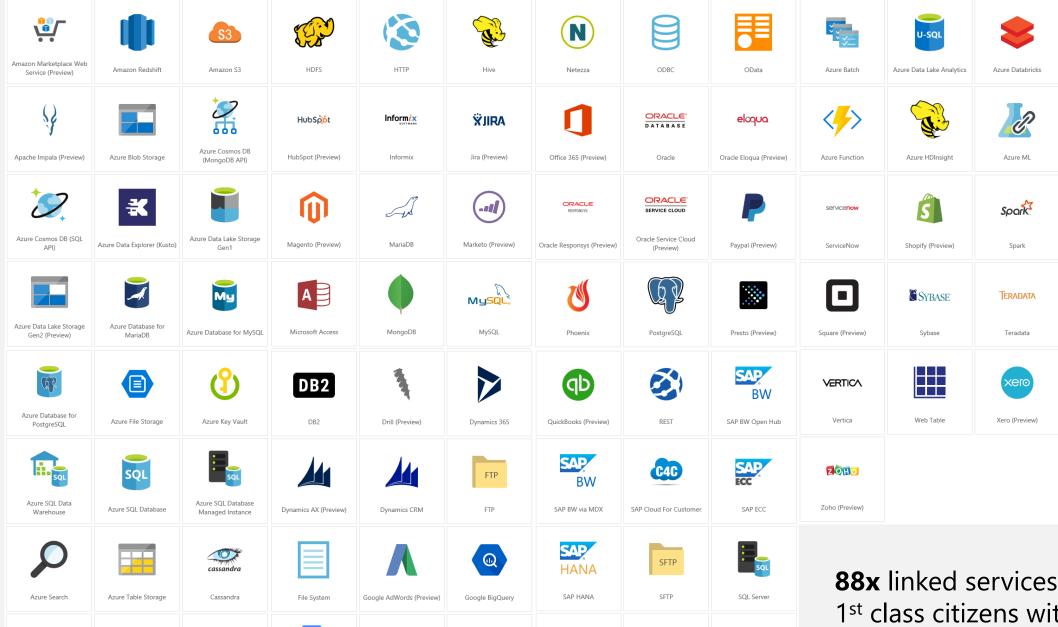












HBase

Greenplum

Common Data Service for

Concur (Preview)

Couchbase (Preview)

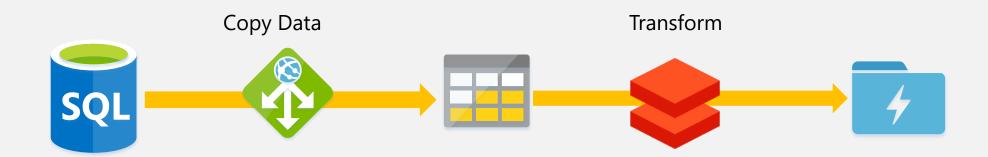
Google Cloud Storage (S3

Salesforce Marketing Cloud

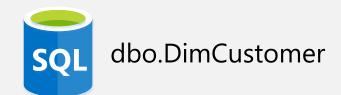
Salesforce Service Cloud

Salesforce

88x linked services supported as 1st class citizens within Azure Data Factory. As of 5th Feb 2019.

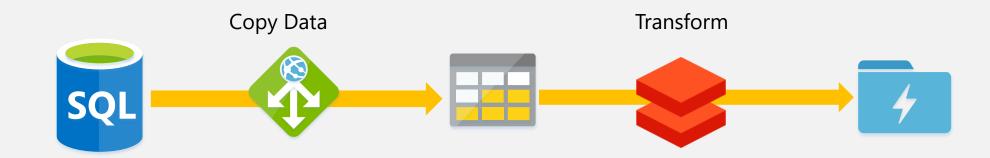


- Linked Services
- 2 Data Sets Where is my data? What format? What file path/table do I need?









- 1 Linked Services
- 2 Data Sets
- Activities What do we want to happen?
 With what conditions?



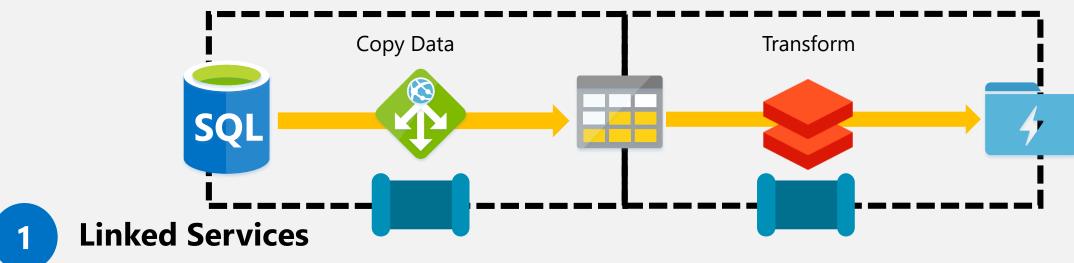
notebookPath: /Playground/Playing

baseParameters: Testing

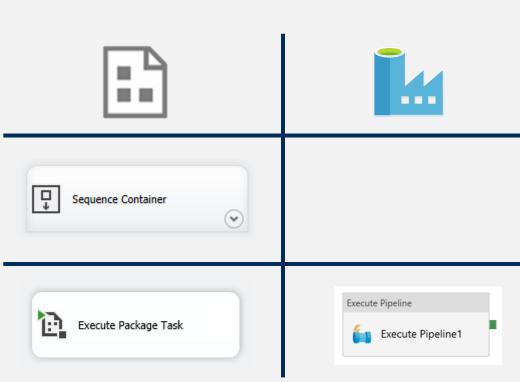
libraries[jar]: dbfs:/lib1.jar

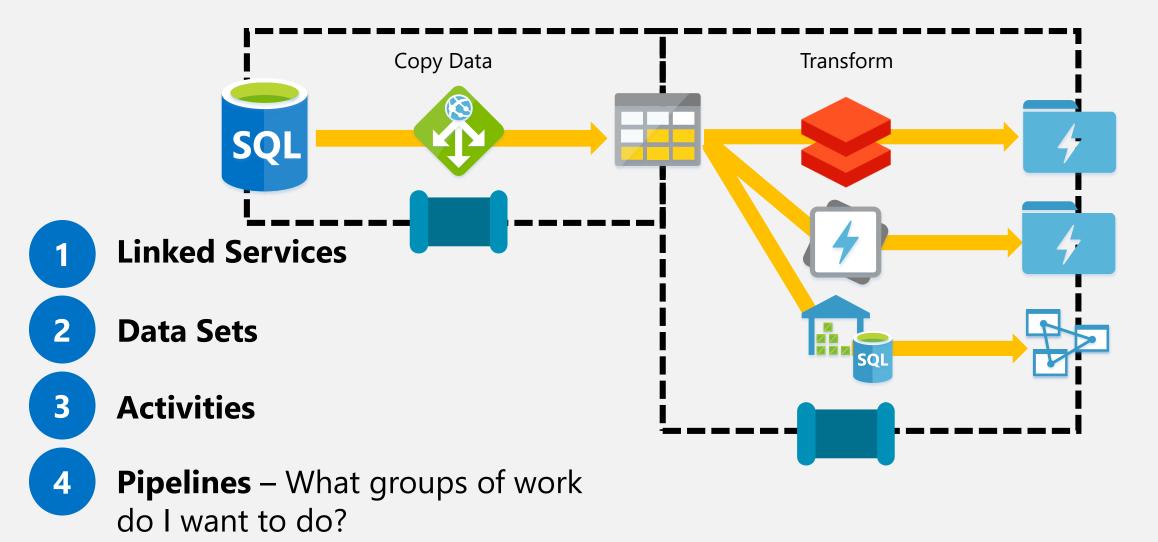
linkedServiceName: BricksOfData01



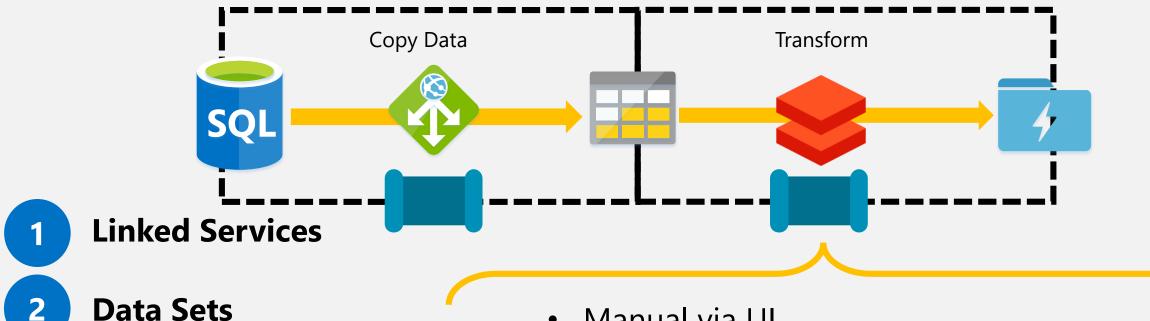


- 2 Data Sets
- 3 Activities
- 4 **Pipelines** What groups of work do I want to do?





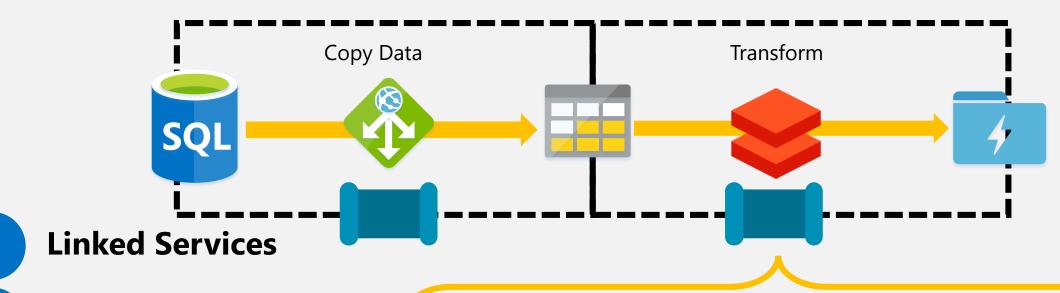




- **Activities**
- **Pipelines**

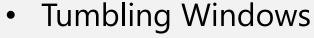
- Manual via UI
- **Tumbling Windows**
- Scheduled
- **Blob File Events**
- Logic App Calls
- **Triggers** How are we going to tell our pipeline(s) to execute?





- 2 Data Sets
- 3 Activities
- 4 Pipelines
- 5 Triggers





- Scheduled
- Blob File Events
- Logic App Calls

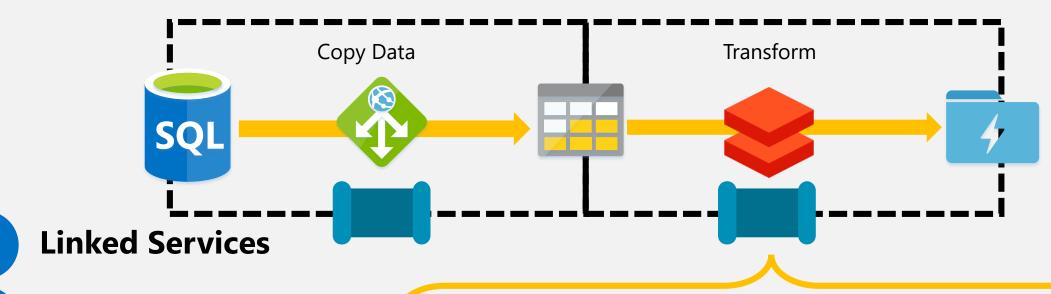




Invoke-AzureRmDataFactoryV2Pipeline

- -DataFactoryName \$dataFactoryName
- -ResourceGroupName \$resourceGroupName
- -PipelineName \$pipelineName



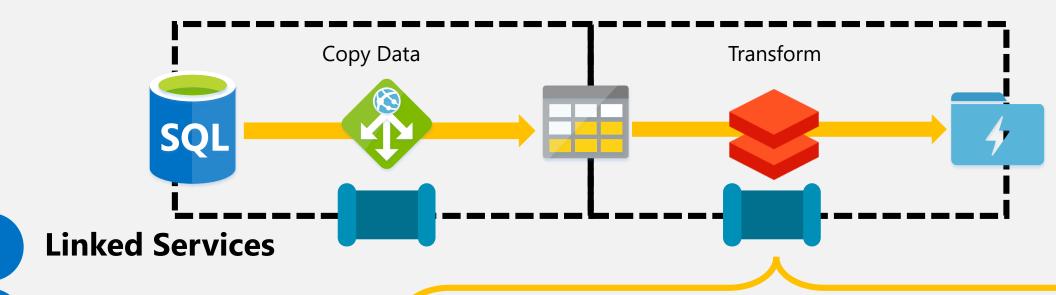


- 2 Data Sets
- 3 Activities
- 4 Pipelines
- 5 Triggers

- Manual via UI
- Tumbling Windows AKA Time Slices
- Scheduled
- Blob File Events
- Logic App Calls







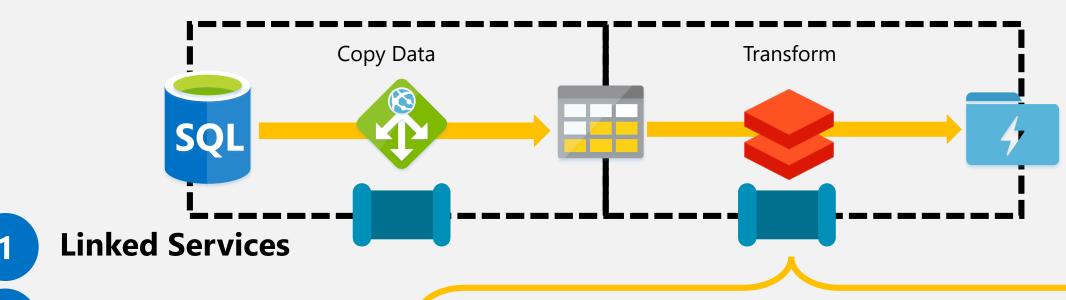
- 2 Data Sets
- 3 Activities
- 4 Pipelines
- 5 Triggers

- Manual via UI
- Tumbling Windows
- Scheduled
 - Blob File Events
 - Logic App Calls



- Every 1 minute.
- UTC





- 2 Data Sets
- 3 Activities
- 4 Pipelines
- 5 Triggers

- Manual via UI
- Tumbling Windows
- Scheduled
- Blob File Events
- Logic App Calls

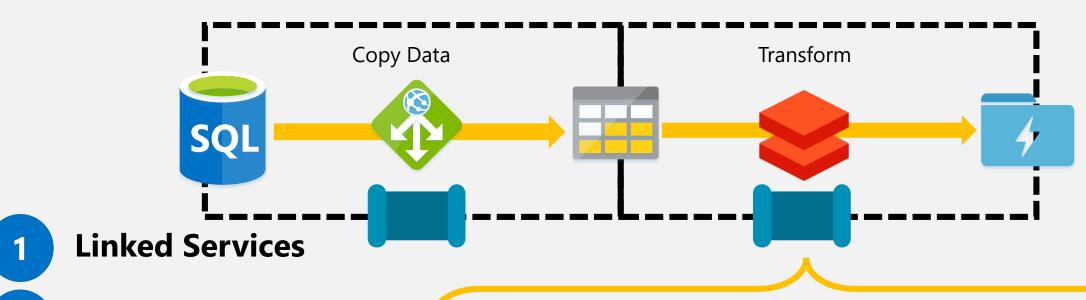




- {Path} Created

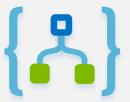
{Path} Deleted

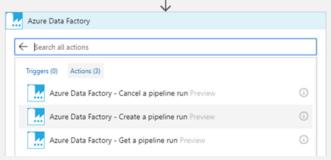


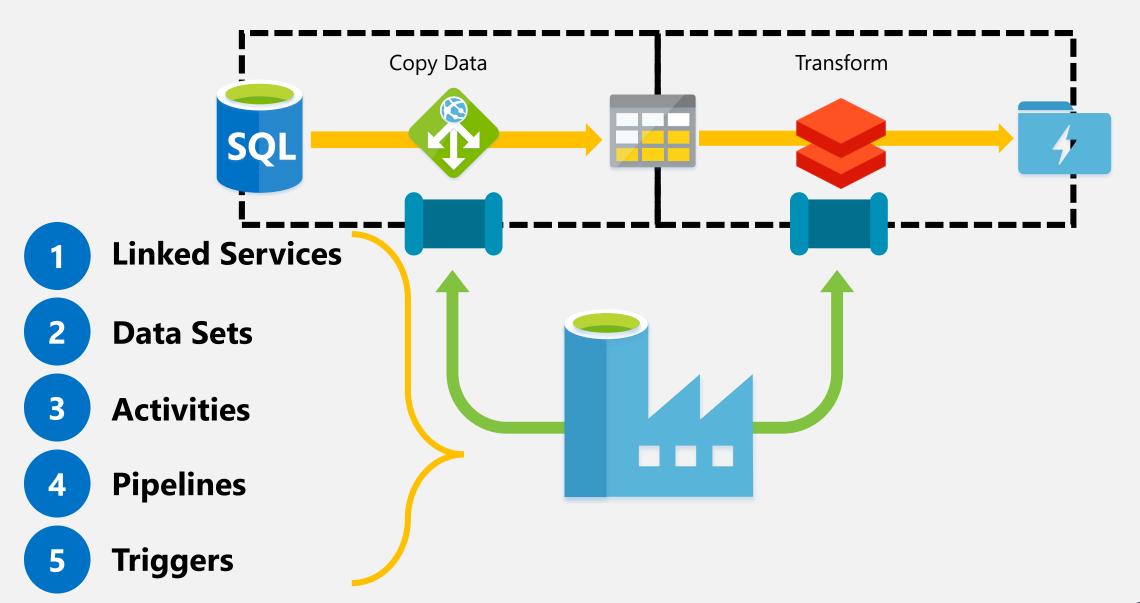


- 2 Data Sets
- 3 Activities
- 4 Pipelines
- 5 Triggers

- Manual via UI
- Tumbling Windows
- Scheduled
- Blob File Events
- Logic App Calls

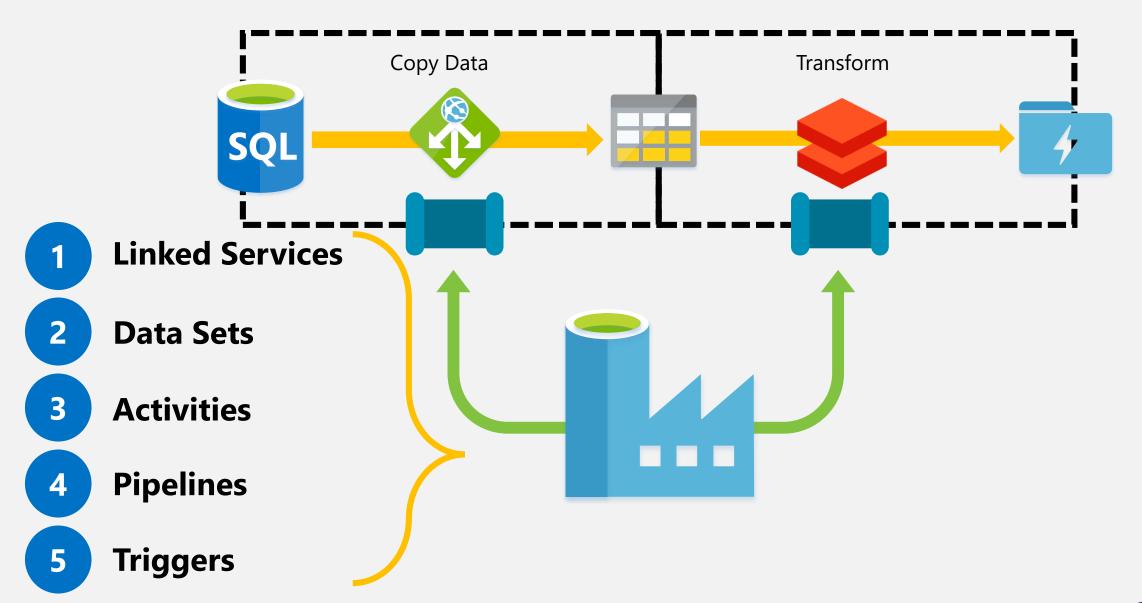






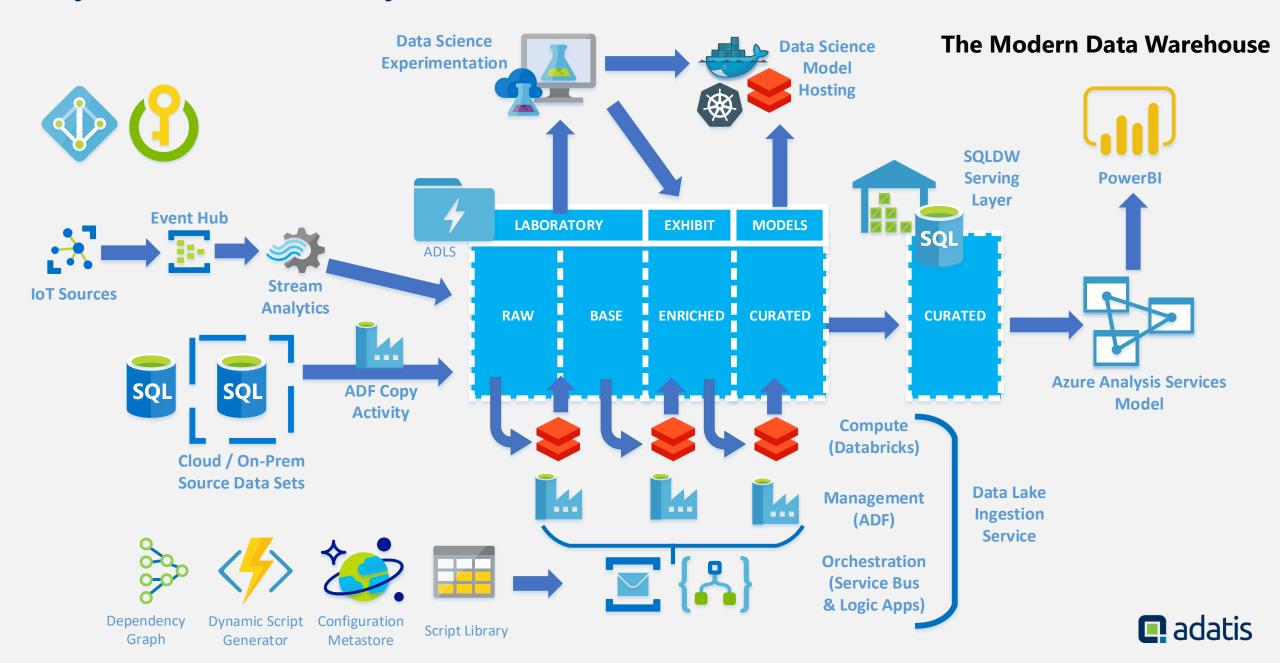


Data Factory Control Flow Components





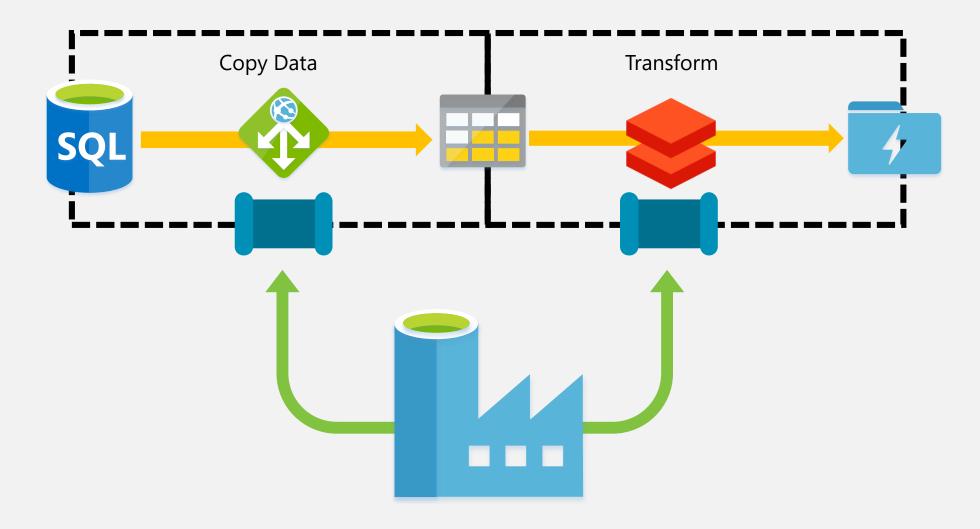
Why use Azure Data Factory?



Data Transformation in Azure

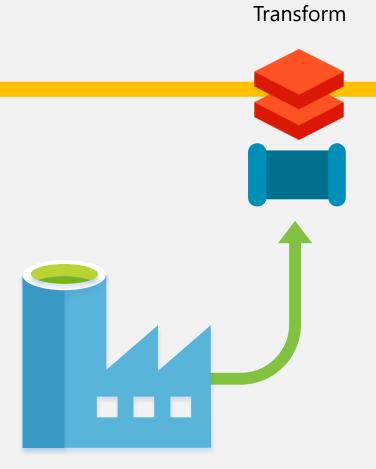


Data Factory Control Flow Components



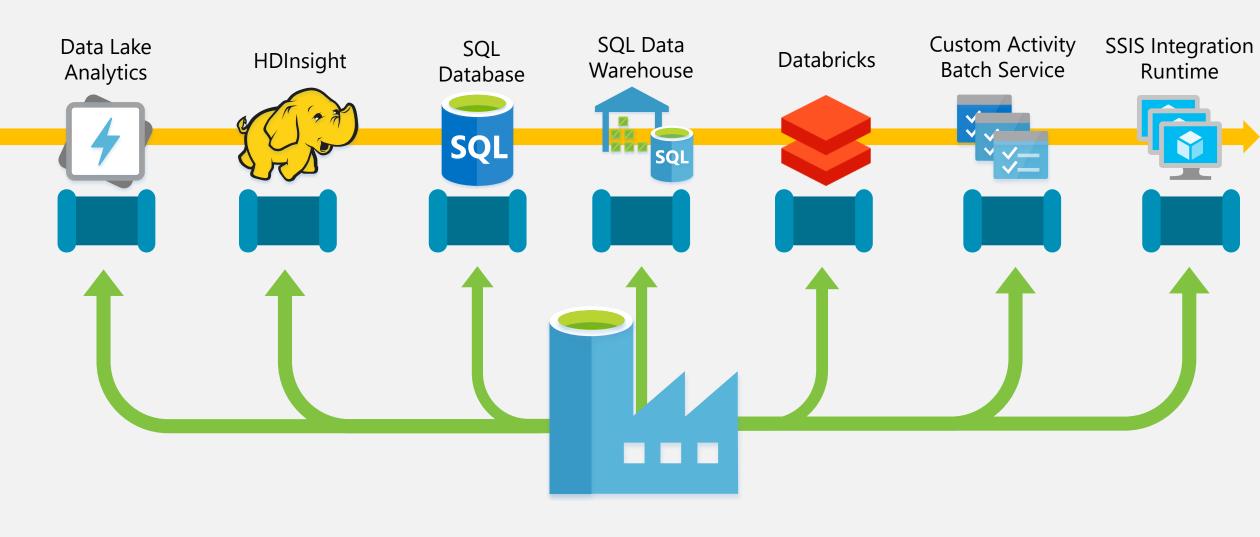


Data Transformation in **Azure**





Data Transformation in Azure



Future Uncertain

Expensive Clusters

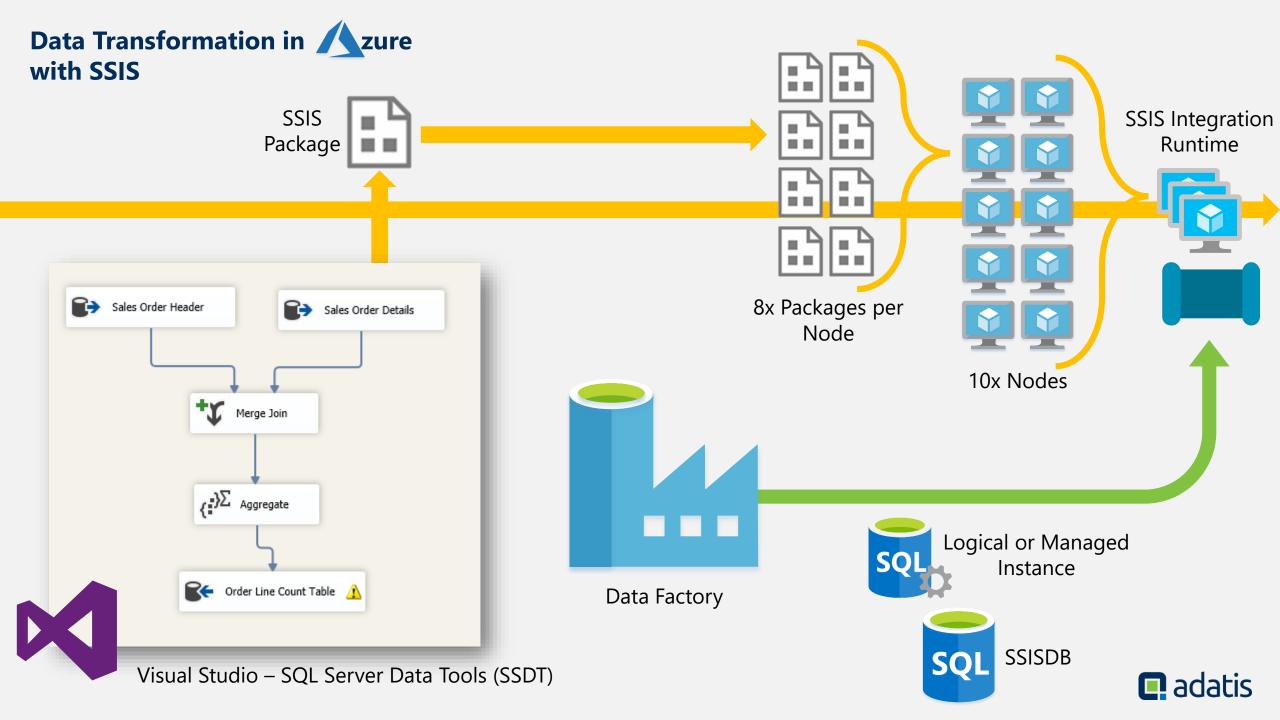
Only Scales Up

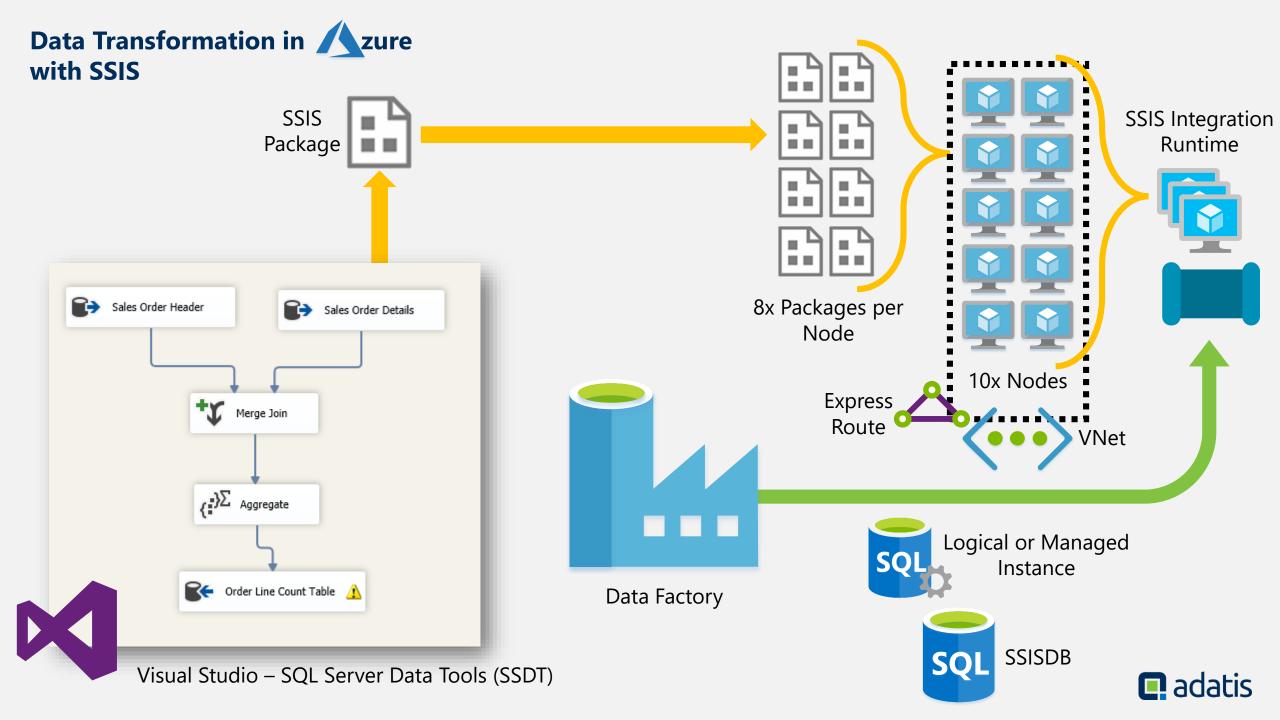
Requires at least 1TB

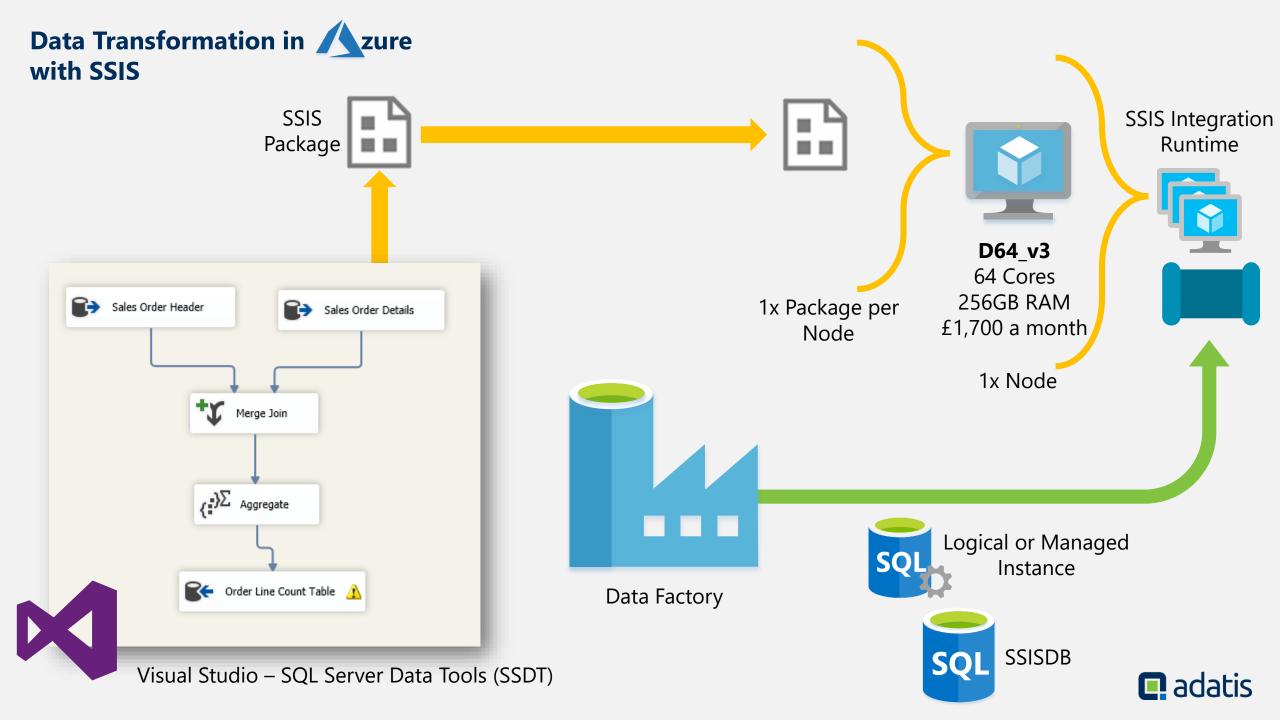
Learn Spark, Python/Scala

Custom Apps on laaS

IaaS VMs
20min start
adatis



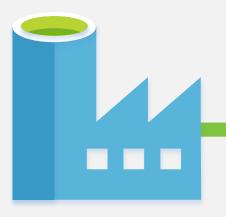




SSIS Integration Runtime

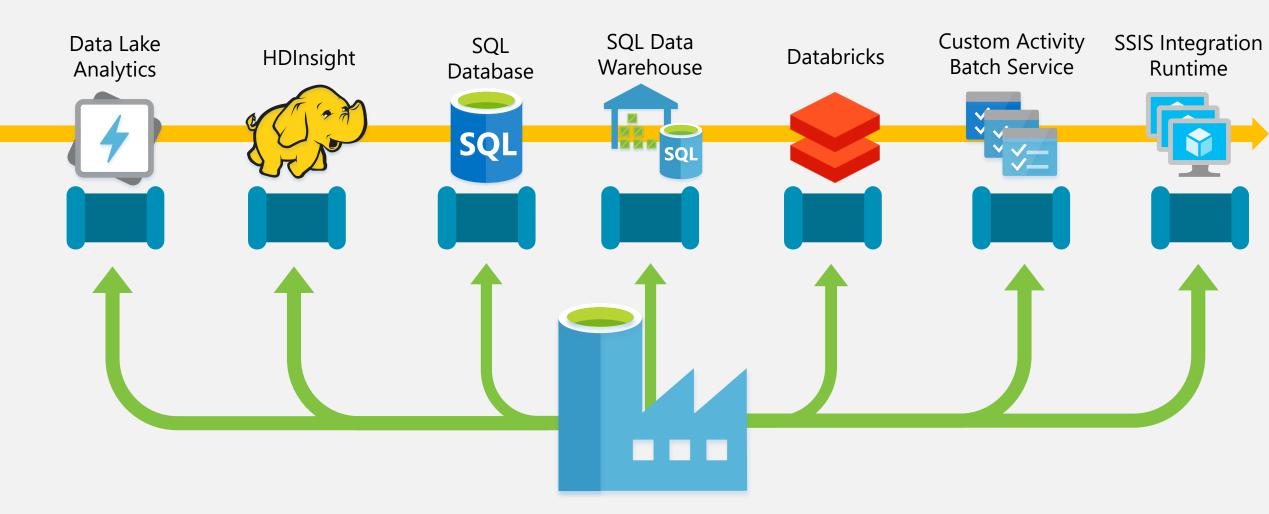






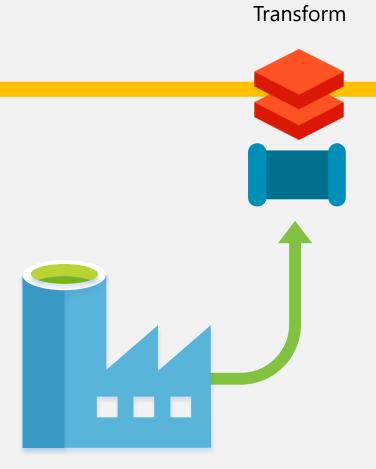


Data Transformation in Azure



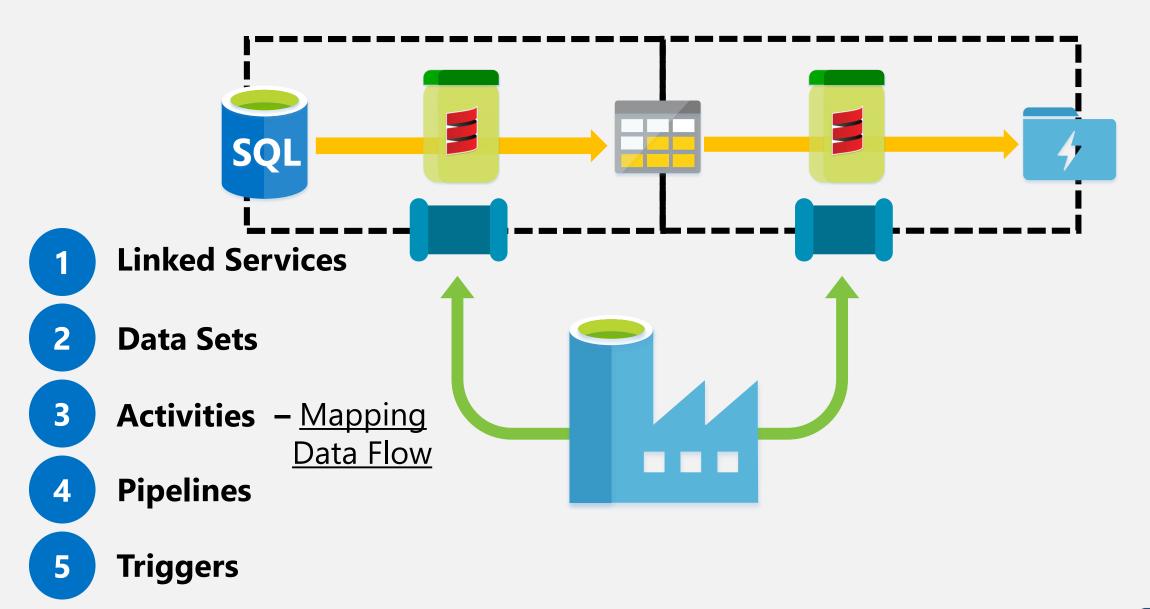


Data Transformation in **Azure**





Data Factory Control Flow Components

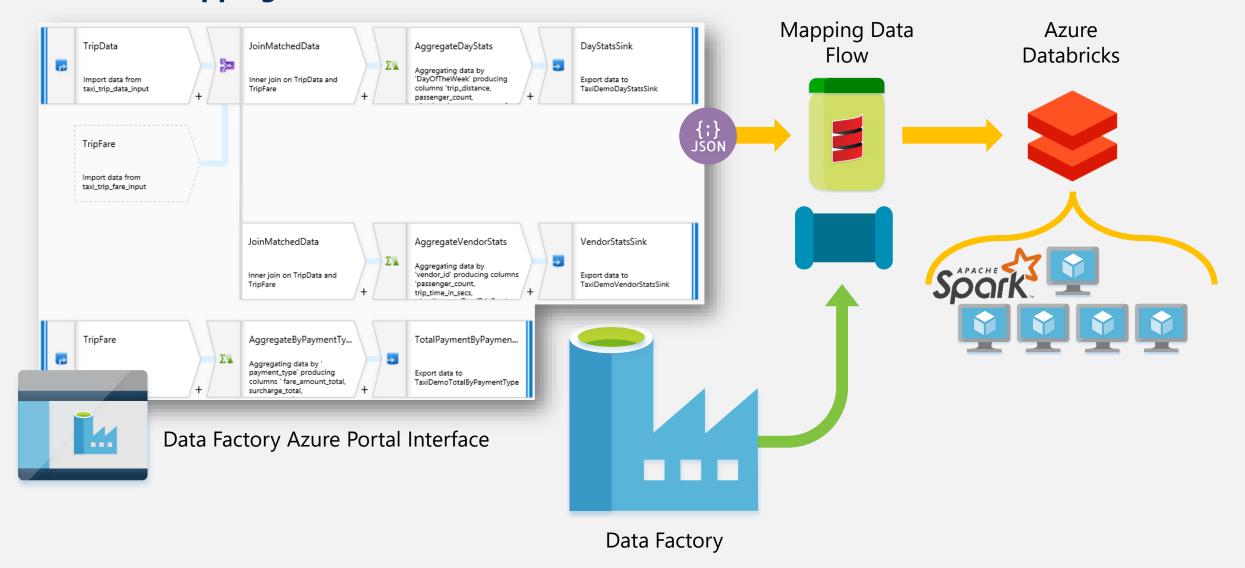




Mapping Data Flows

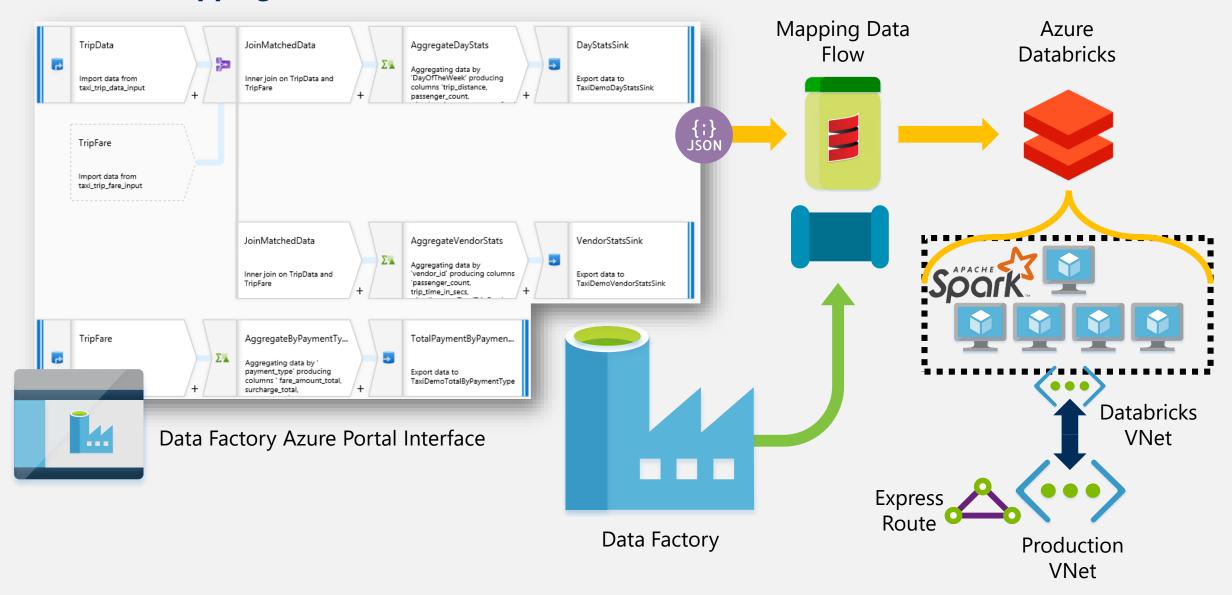


What is a Mapping Data Flow?



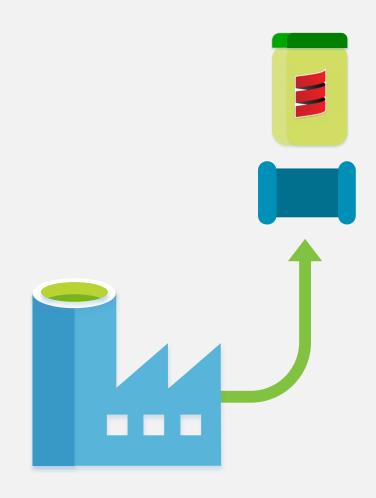


What is a Mapping Data Flow?

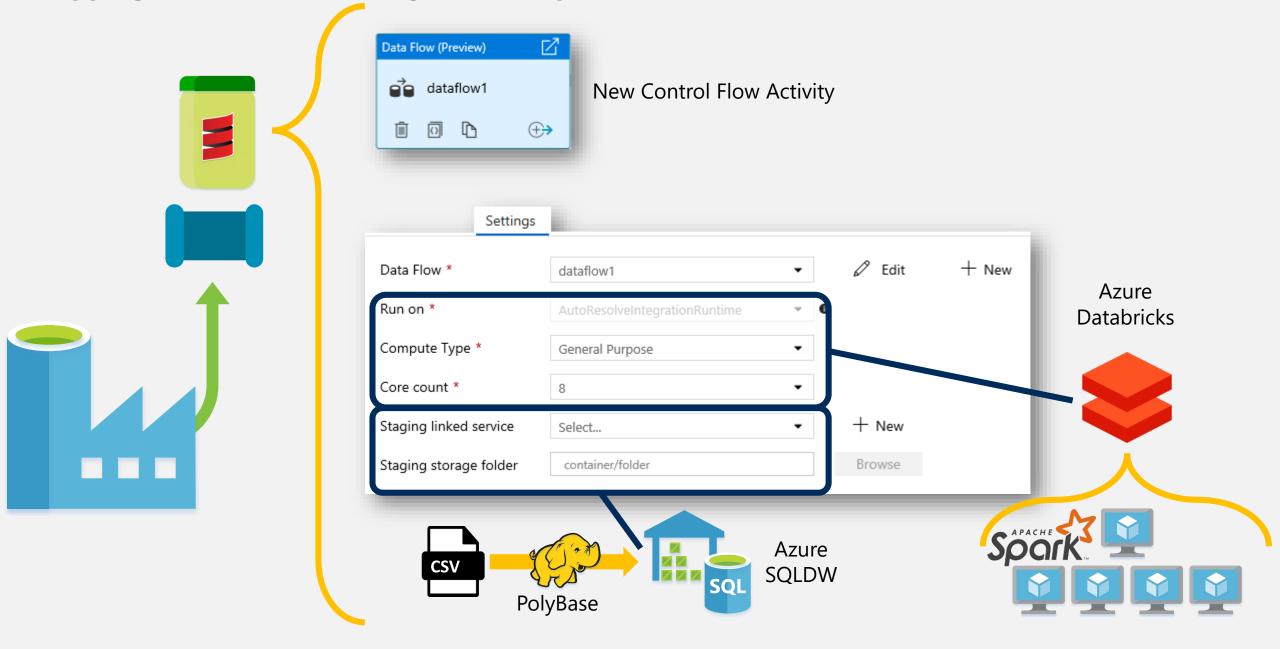


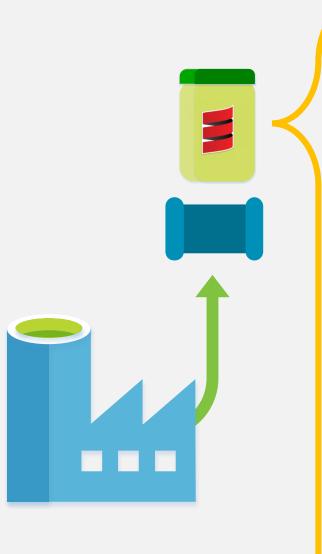


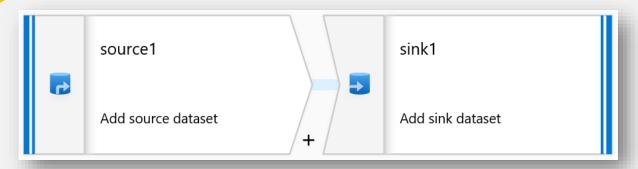
Mapping Data Flows



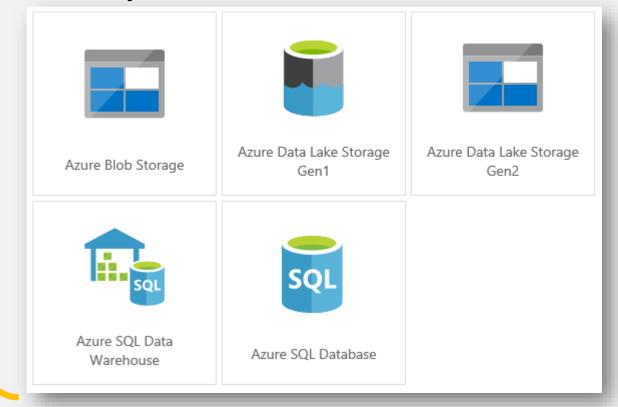




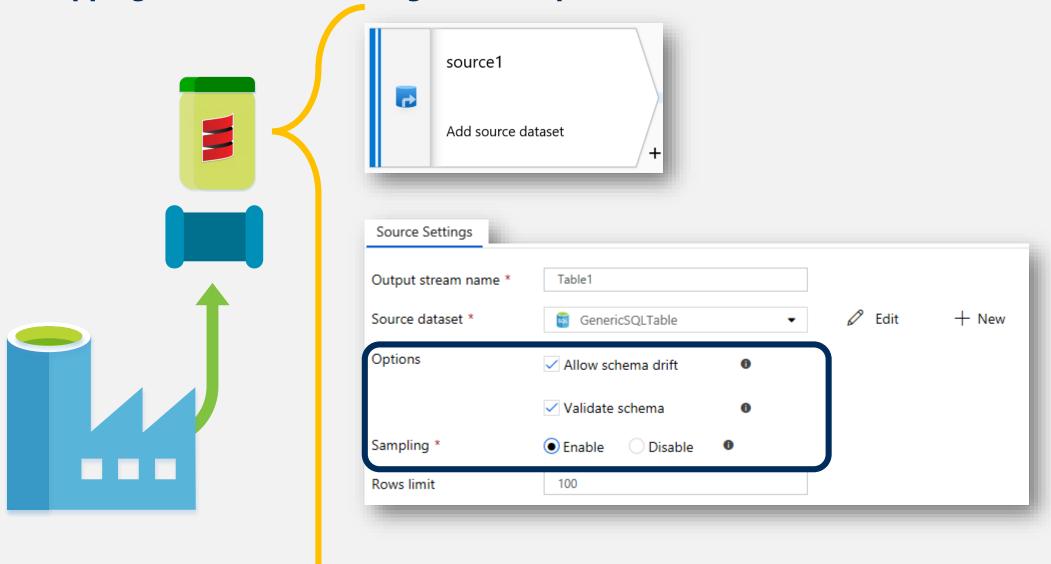




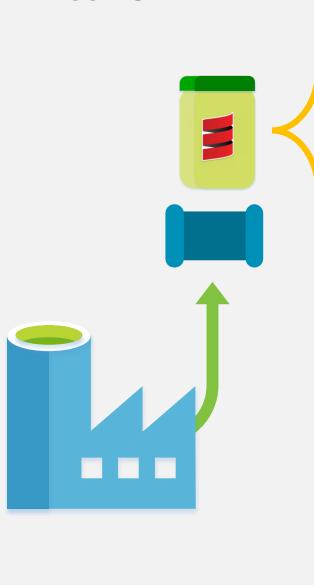
Currently Available:



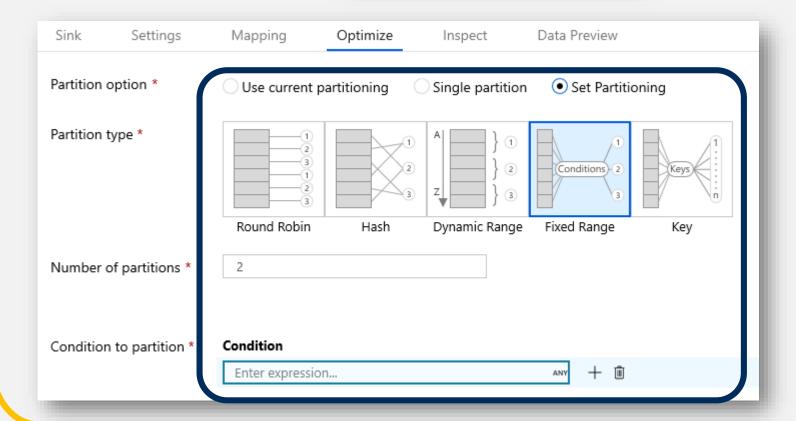






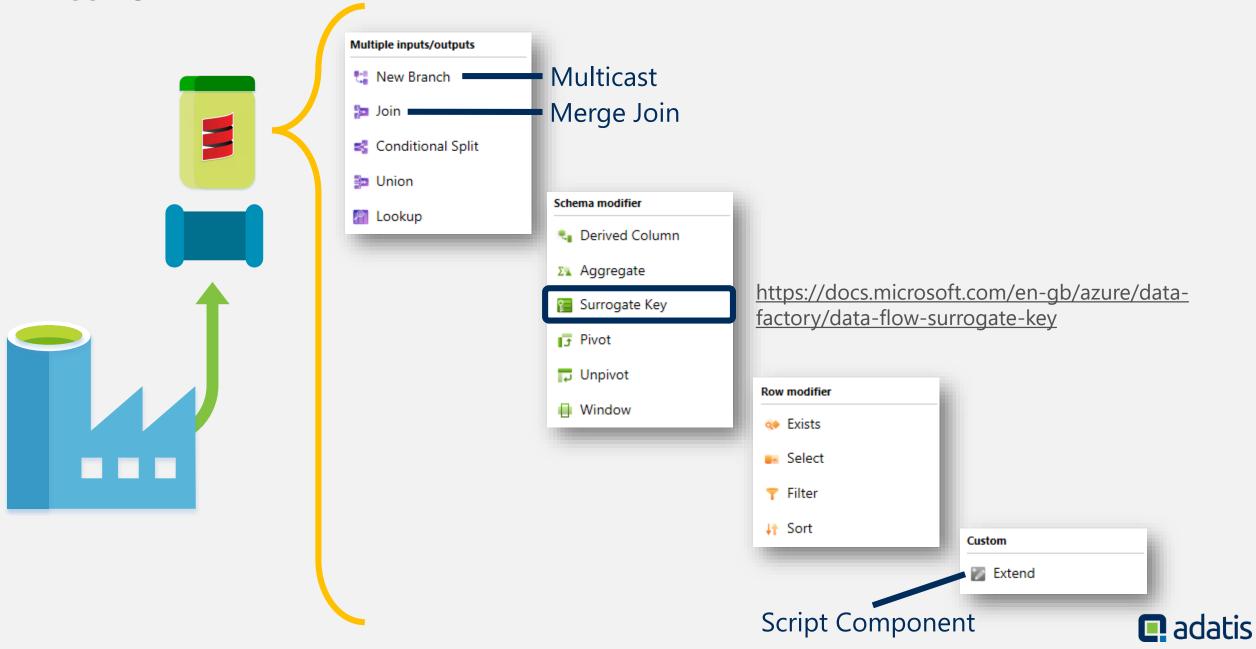




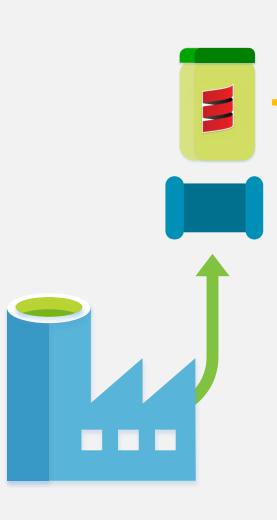


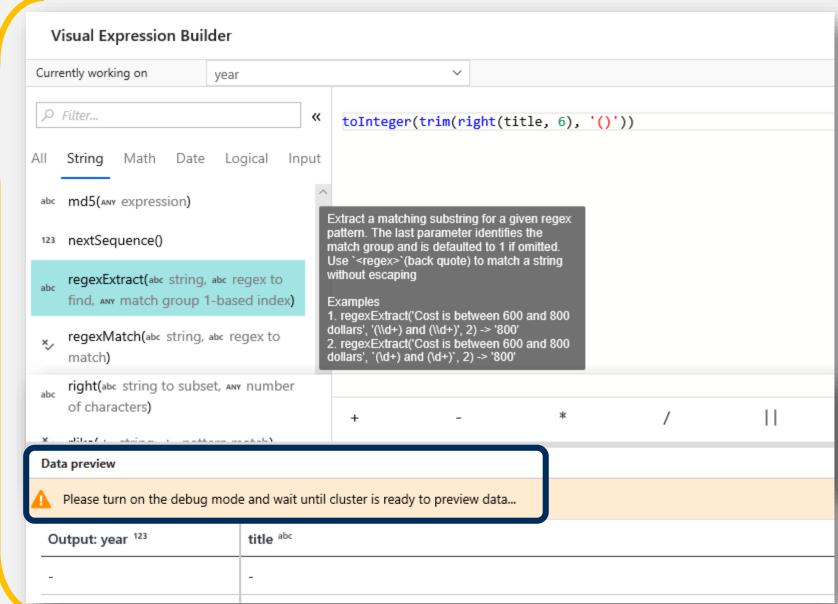


Mapping Data Flows – Transformations



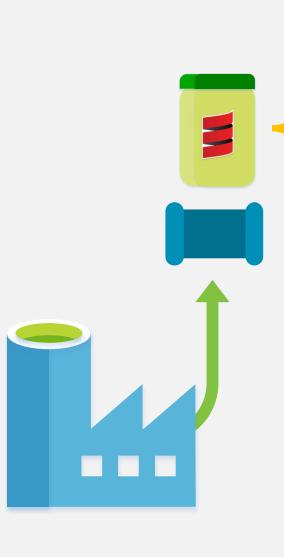
Mapping Data Flows – Expression Builder

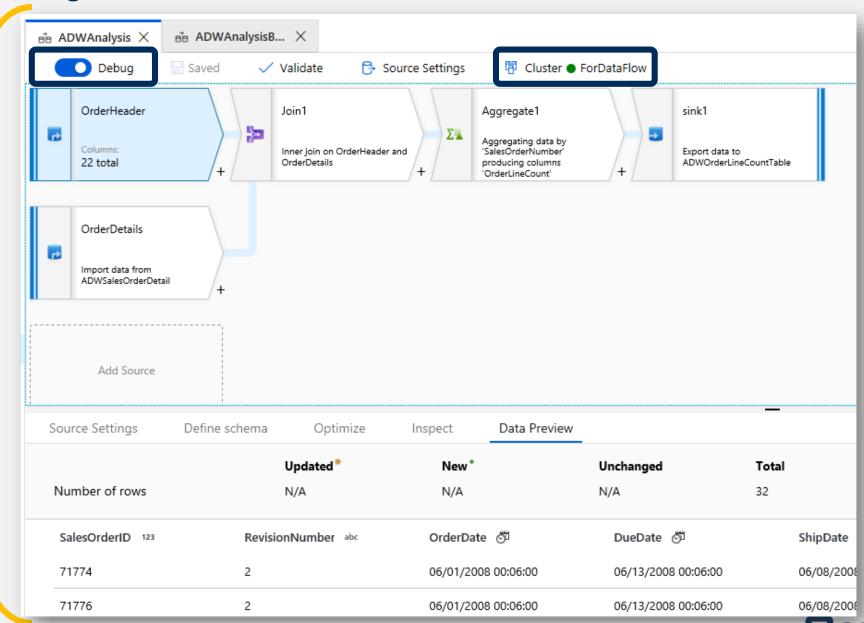




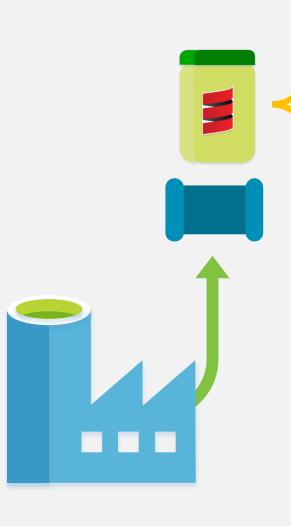


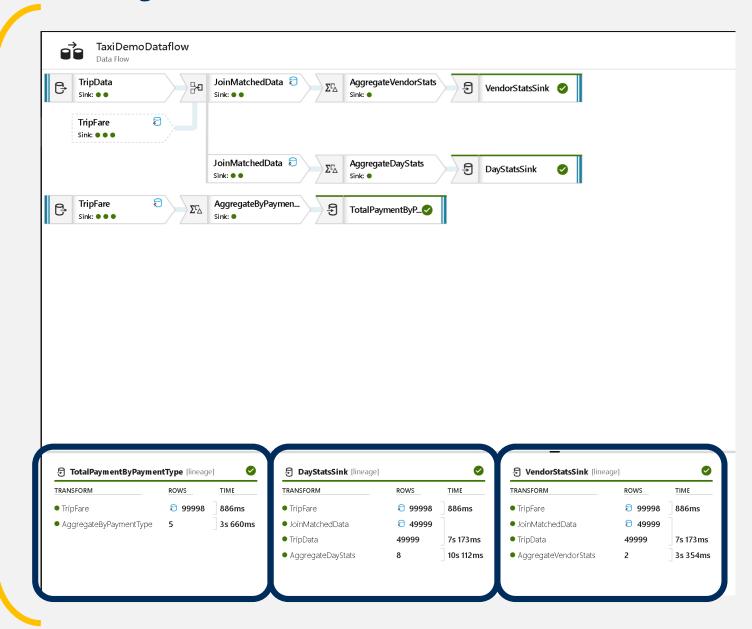
Mapping Data Flows – Debug Mode





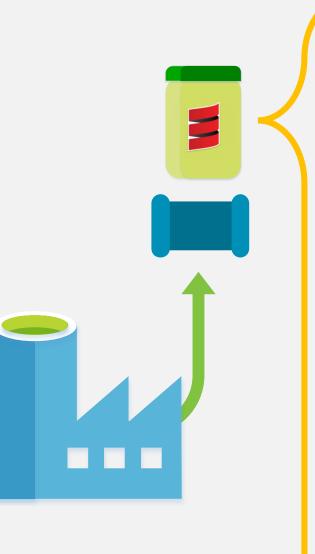
Mapping Data Flows – Monitoring







Mapping Data Flows



Activity

https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-overview

Source & Sink

https://docs.microsoft.com/en-qb/azure/data-factory/concepts-data-flow-schema-drift

Transformations

https://docs.microsoft.com/en-gb/azure/data-factory/data-flow-aggregate

Expression Builder

https://docs.microsoft.com/en-gb/azure/data-factory/data-flow-expression-functions

Debug Mode

https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-debug-mode

Monitoring

6

https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-monitoring

https://github.com/kromerm/adfdataflowdocs/tree/master/videos



Demo



Demo Summary

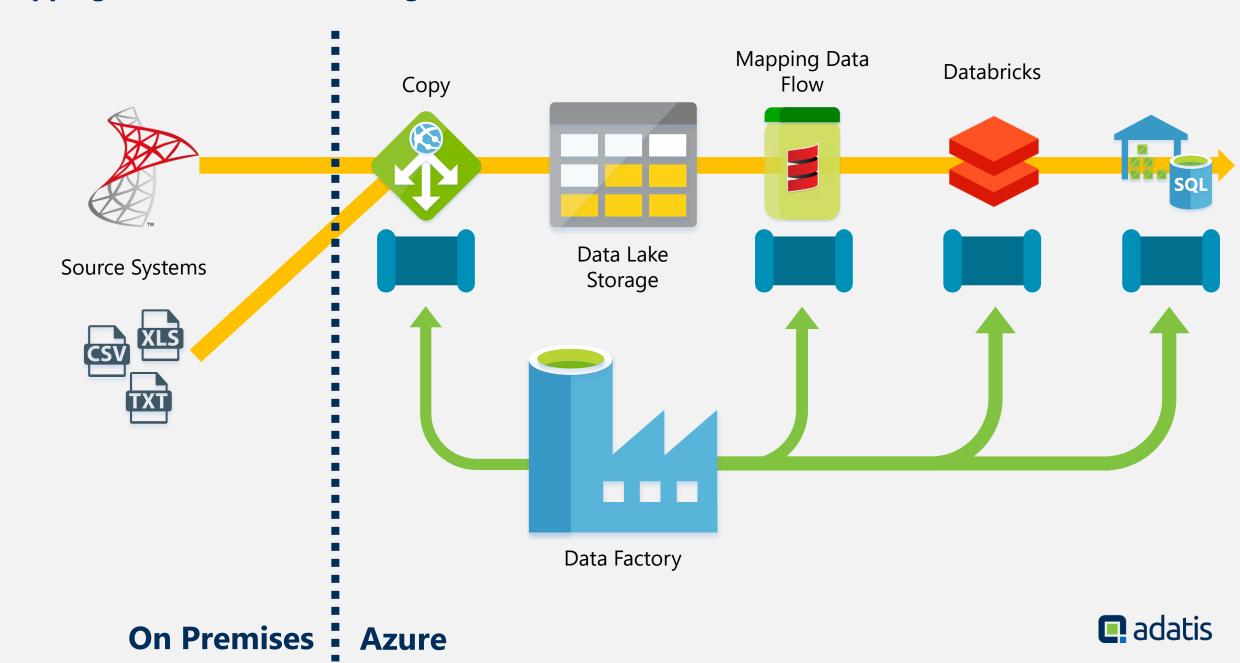
Transformation Method		Graphical Development	Scales Out	Scales Up	Cloud Native Tech
SQL	T-SQL (SQLDB)	×	*		*
	SSIS		*		*
	Scala (Databricks)	*			
	Mapping Data Flow				



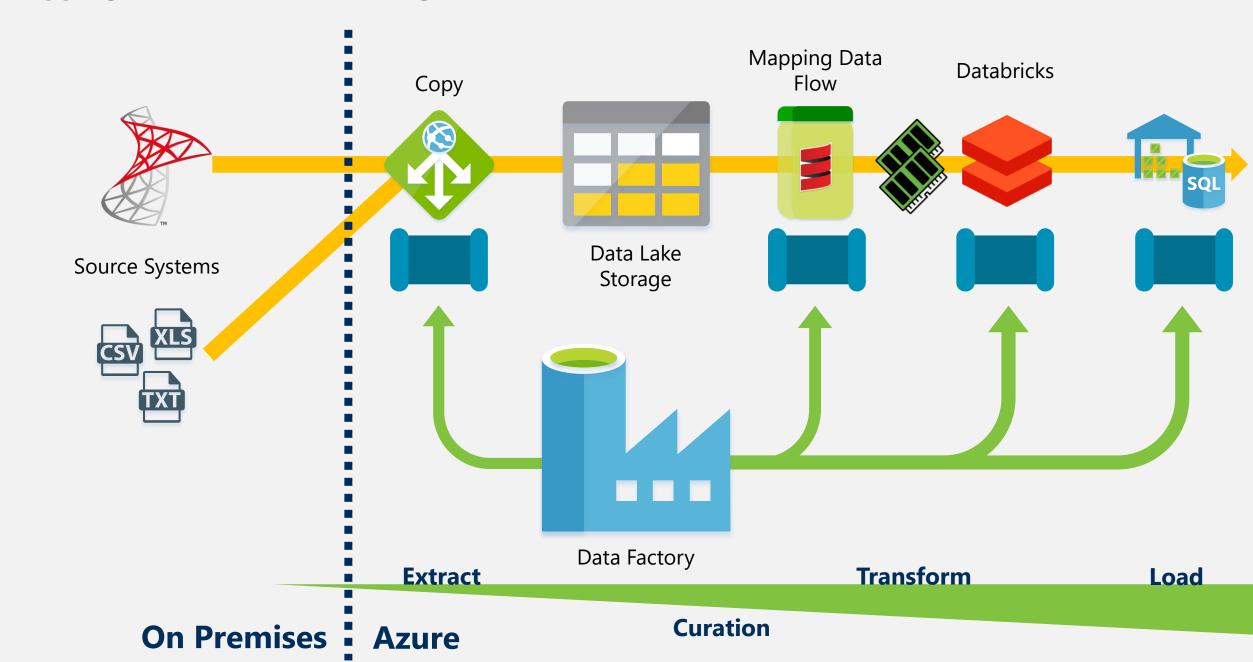
Design Patterns



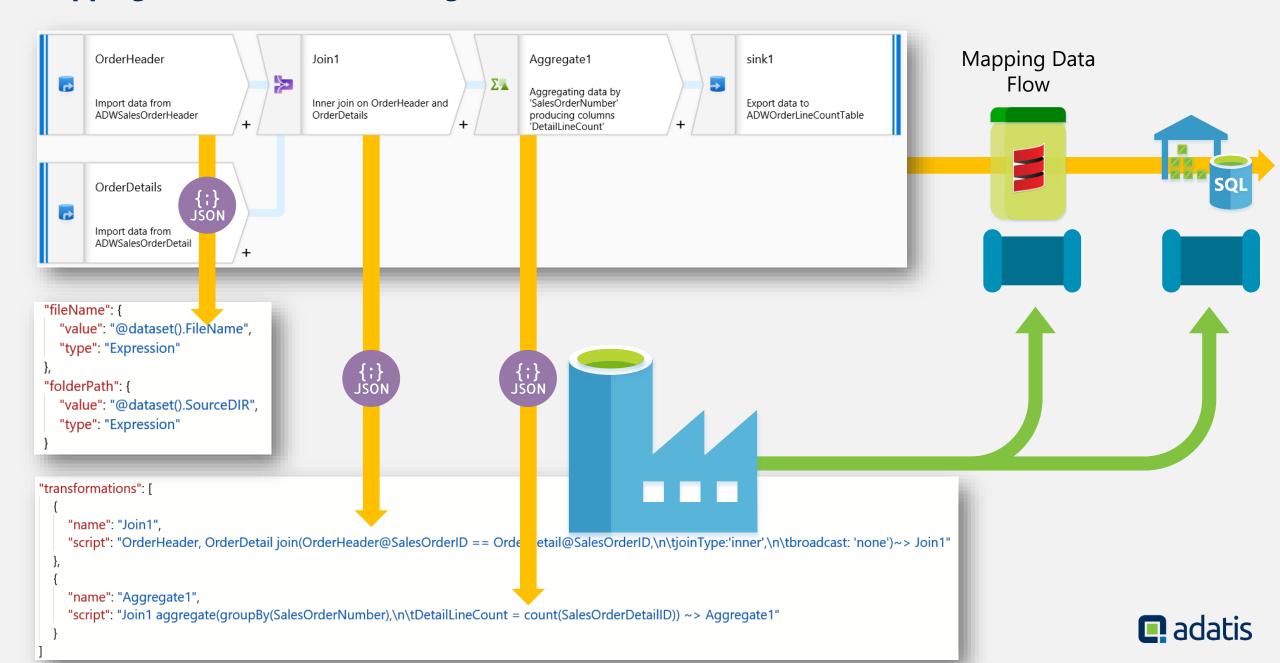
Mapping Data Flow Future Design Patterns ???



Mapping Data Flow Future Design Patterns ???



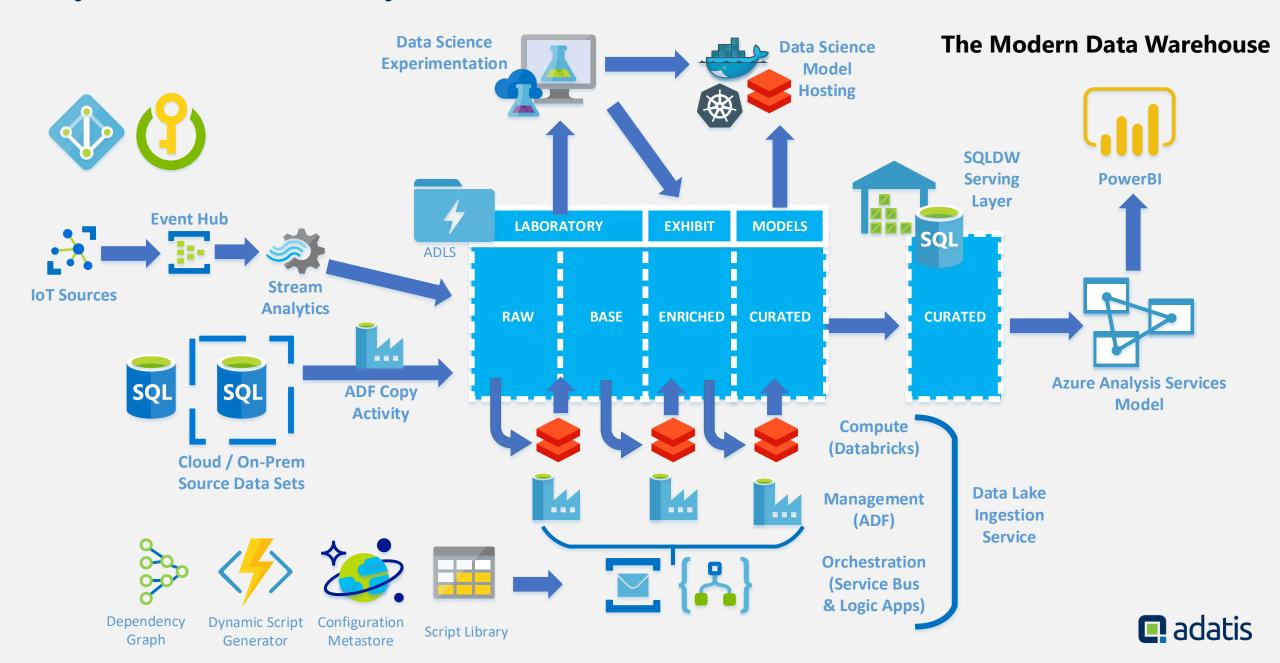
Mapping Data Flow Future Design Patterns ???



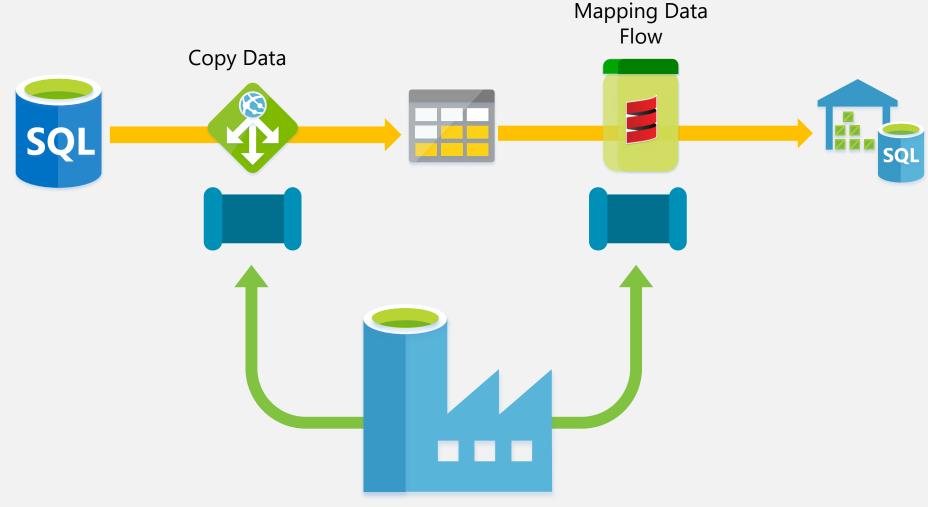
Conclusion



Why use Azure Data Factory?



What is Azure Data Factory?

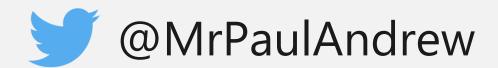


Orchestrator of our solution <u>Control Flow</u> operations. Orchestrator of our solution <u>Data Flow</u> transformations.

... using cloud native technology in Azure and now with an easy developer interface for both.

Thanks for Listening

Paul Andrew







paul@mrpaulandrew.com **Email:**

mrpaulandrew.com Blog:

GitHub: github.com/mrpaulandrew

