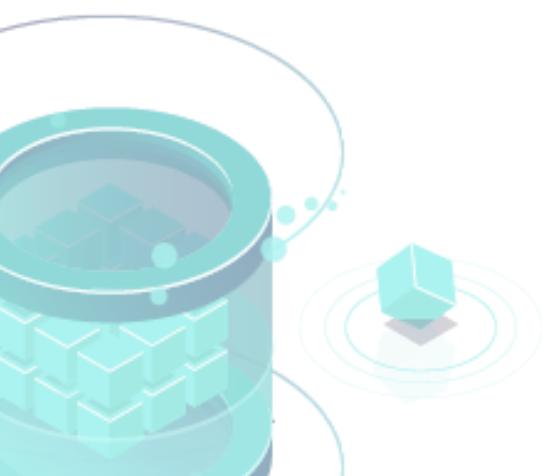


Azure Data Factory 2022

ADF & Synapse Pipelines - Intro & Unterschiede
Was gibt es Neues in ADF?

Stefan Kirner



Stefan Kirner



- › PASS Chapter Lead Karlsruhe ski@sqlpass.de & Beirat
- › Director Business Intelligence scieneers GmbH
- › Twitter: @KirnerKa



Agenda

Part 1 Azure Data Factory vs. Synapse Data Pipelines

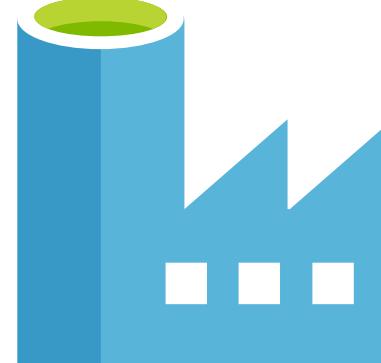
- Intro
- Features Azure Data Factory only / Synapse Pipelines only
- Lieber Azure Data Factory vs. Synapse Pipelines wenn...

Part 2 Azure Data Factory - neue Top Features 2021/2022

- Flowlets und Data Libraries
- Activity & Data Flow updates
- Application Lifecycle Management
- SAP CDC is coming

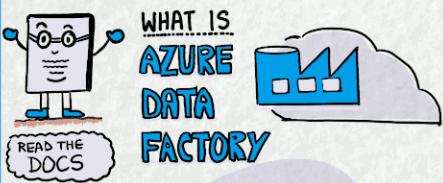
Azure Data Factory vs. Synapse Pipelines

Was ist was und wann benutze ich welches?



Was ist Azure Data Factory?

Ganz einfach...



A CLOUD-BASED DATA INTEGRATION SERVICE THAT ORCHESTRATES DATA MOVEMENT & TRANSFORMATION BETWEEN DIVERSE DATA SOURCES AT SCALE & CLOUD COMPUTE RESOURCES



7 THINGS TO KNOW ABOUT AZURE DATA FACTORY

1 ENTERPRISE READY

Data integration at cloud scale!

2 ENTERPRISE DATA READY

90+ connectors! It just works.

3 CODE-FREE TRANSFORMATION

UI driven mapping data flows

4 RUN CODE ON ANY AZURE COMPUTE

For hands-on data transformations

5 MANY SSIS PACKAGES RUN ON AZURE

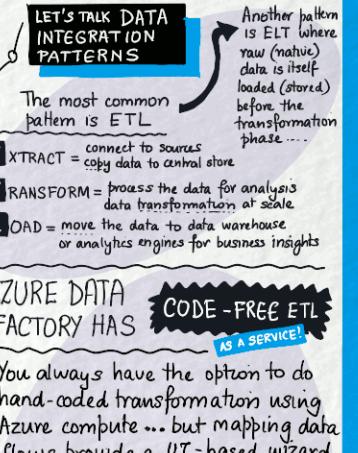
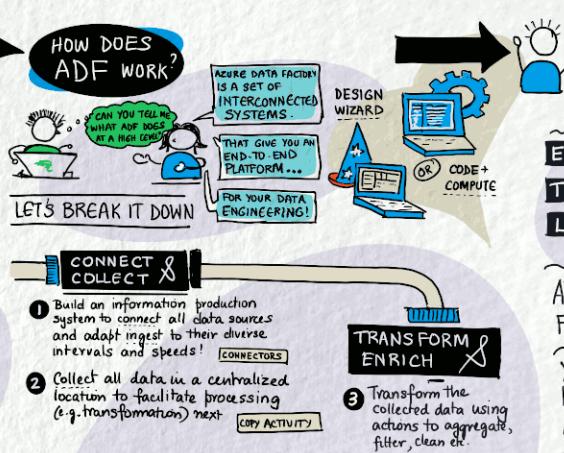
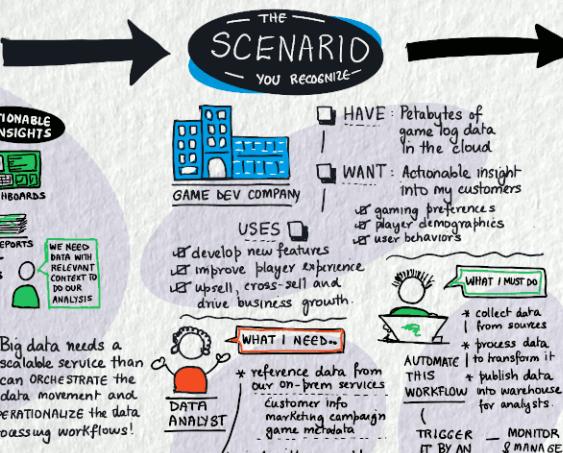
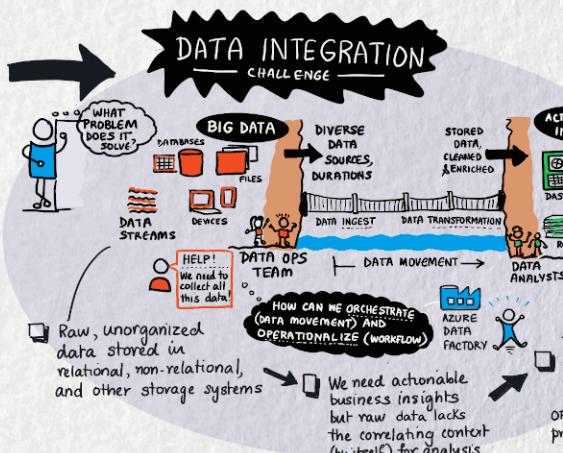
(Rehost on-prem in ADF in 3 steps)

SOURCE CONTROL AUTOMATED DEPLOYMENT

6 ADF CAN MAKE DATA OPS SEAMLESS

Managed virtual networks protect against data exfiltration, simplify your networking!

7 SECURE DATA INTEGRATION



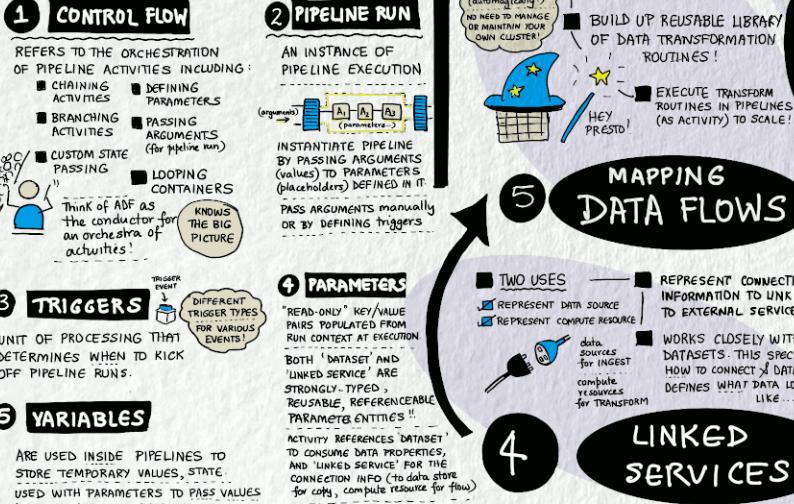
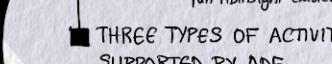
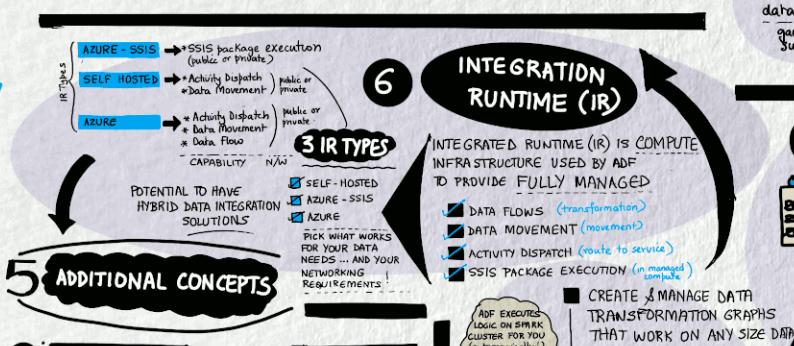
AZURE DATA FACTORY HAS CODE-FREE ETL AS A SERVICE!

You always have the option to do hand-coded transformation using Azure compute ... but mapping data flows provide a UI-based wizard to simplify your pipeline setup :

5 ASPECTS IT COVERS!

- 1 INGEST DATA
 - 2 CONTROL FLOW
 - 3 DATA FLOW
 - 4 SCHEDULE OPS
 - 5 MONITOR OPS
- * Copy data on-prem or in cloud!
 - * 90+ native connectors
 - * Serverless, auto-scaling!
 - * Design (via UI) or Create (via SDK)
 - * Utilize workflow constructs (loop, parameters, variables...) for efficiency
 - * Code-free data transforms (execute automatically in Spark)
 - * Scale out with Azure Integration Runtimes - generate data flows with SDK (code) or Designer (UI)
 - * Build/maintain ops schedules for your data pipelines
 - * Options include: Wall clock, Event-based, Tumbling windows, Sliding windows
 - * View active executions and pipeline history - details at activity/pipeline
 - * Establish alerts for key events or progress notifications

- 1 PIPELINES
 - 2 ACTIVITIES
 - 3 DATASETS
 - 4 LINKED SERVICES
 - 5 MONITOR OPS
 - 6 MAIN CONCEPTS
- A PIPELINE IS A LOGICAL GROUPING OF ACTIVITIES THAT PERFORM ONE UNIT OF WORK
 - AN ADF INSTANCE CAN HAVE ONE OR MORE PIPELINES!
 - Example: Pipeline that ingests data from Azure Blob (A1), runs Hive query on HDInsight to partition it (A2) and moves result to next store (A3)
 - Benefit: Manage correlated activities as a single entity!
 - Execution: Activities may be chained (run sequentially) or be independent (run in parallel) within pipeline!
 - AND A FEW KEY TERMS
 - JET CONTROL FLOW
 - JET PIPELINE RUN
 - JET TRIGGERS
 - JET PARAMETERS
 - JET VARIABLES
 - ADF TOOLKIT



Azure Data Factory Parts

Microsoft: „Code-free ETL as a service“

Ingest



- Multi-cloud and on-premise hybrid copy data
- 100+ native connectors
- Serverless and auto-scale
- Use wizard for quick copy jobs

Control Flow



- Design code-free data pipelines
- Generate pipelines via SDK
- Utilize workflow constructs: loops, branches, conditional execution, variables, parameters, ...

Data Flow



- Code-free data transformations that execute in Spark
- Scale-out with Azure Integration Runtimes
- Generate data flows via SDK
- Designers for data engineers and data analysts

Schedule



- Build and maintain operational schedules for your data pipelines
- Wall clock, event-based, tumbling windows, chained

Monitor



- View active executions and pipeline history
- Detail activity and data flow executions
- Establish alerts and notifications

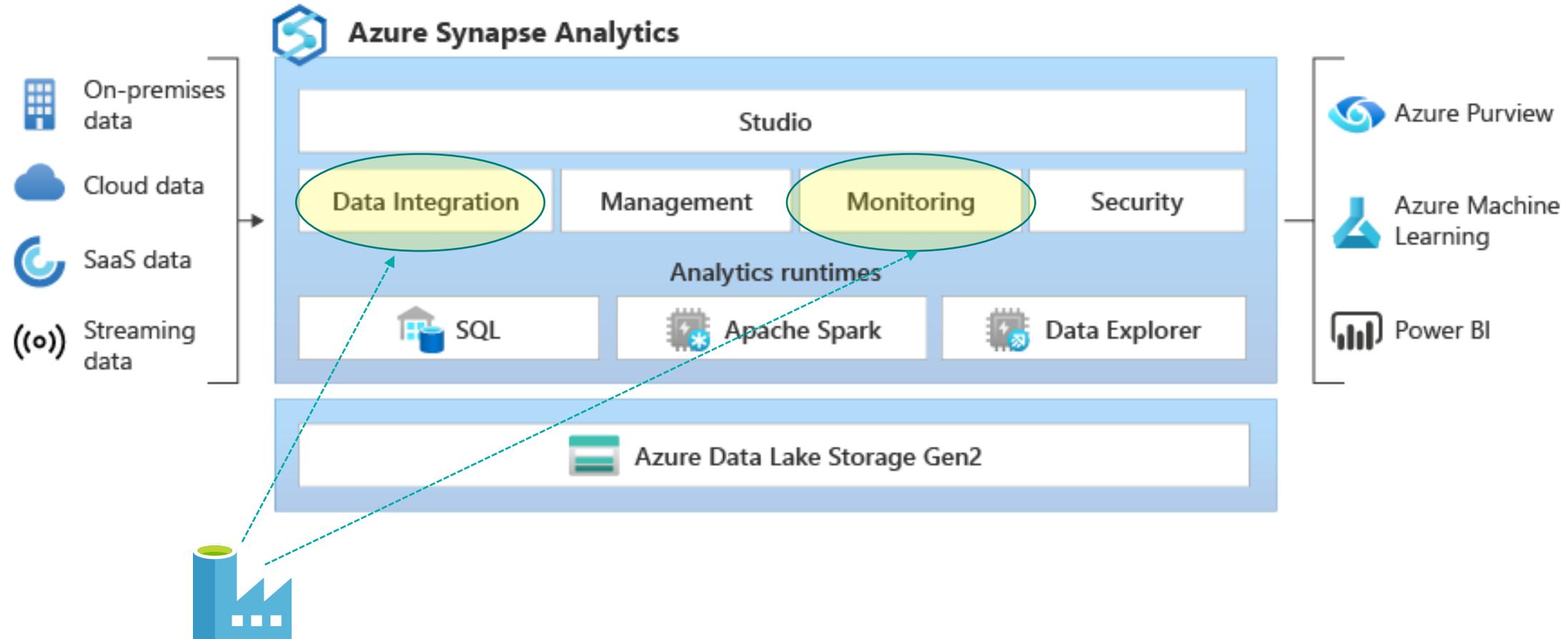


Demo

Überflug Azure Data Factory

What is Azure Synapse Analytics

Microsoft: "Azure Synapse is an enterprise analytics service"





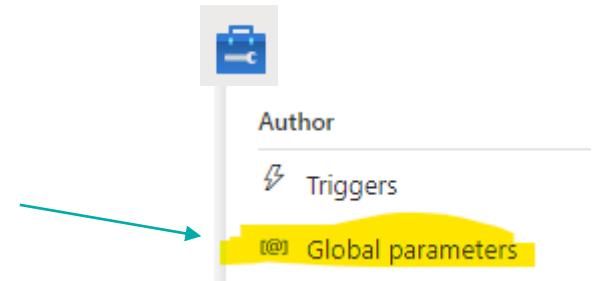
Demo

Überflug Synapse Pipelines

Azure Data Factory only stuff

Not available in Synapse

- Self-hosted Integration runtime sharing
- Support for Cross-region Integration Runtime (Data Flows)
- Power Query support aka Wrapping Data Flows
- Comfortable deployment parameters management
- *Update Snowflake source / destination*
- *Update: SSIS runtime, run SSIS packages -> in Preview for Synapse now*
- *Update: ADF Solution templates also available in Synapse*

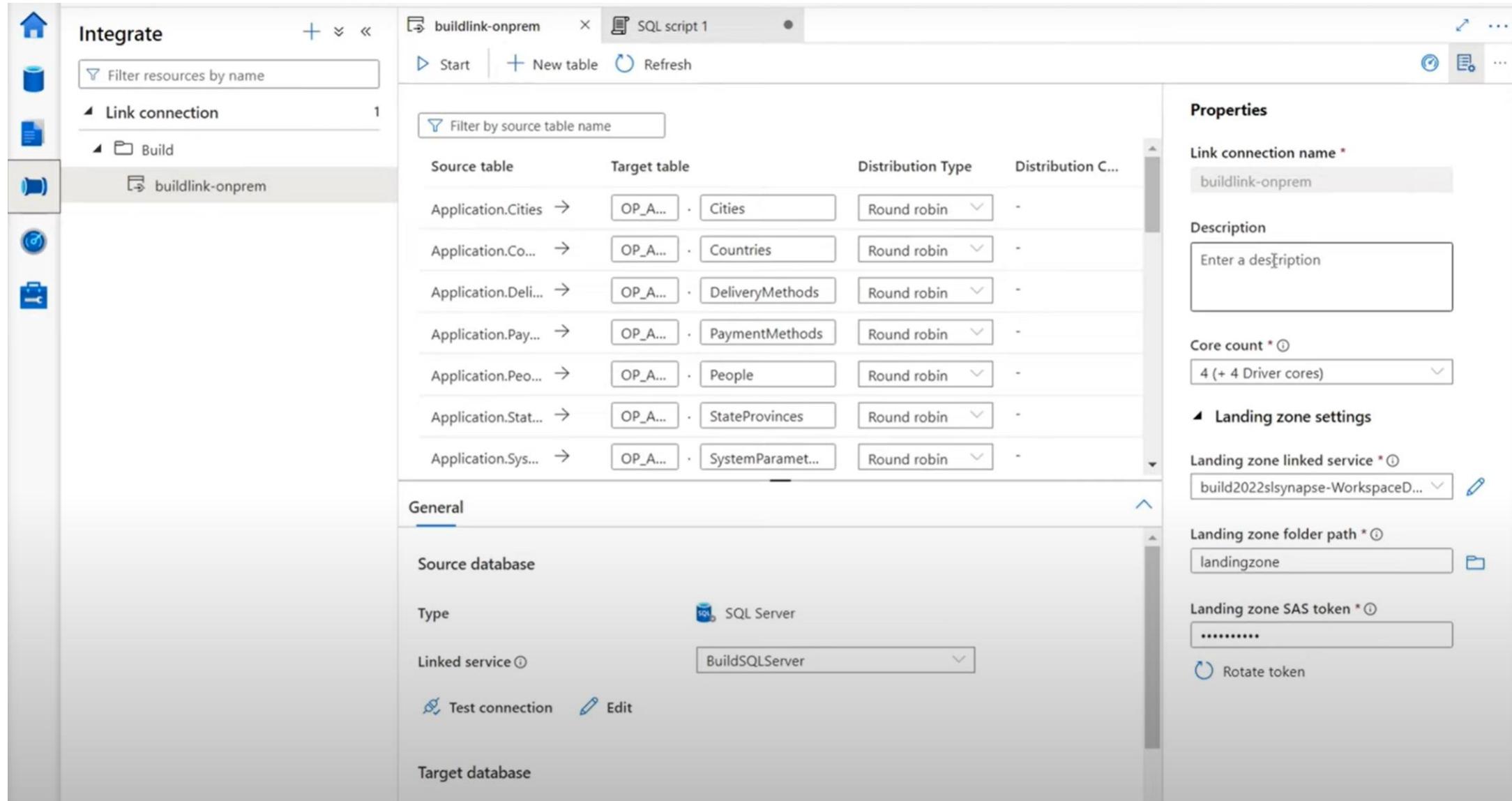


Synapse only stuff (not in ADF)

Not available in Azure Data Factory

- Spark Notebooks & Clusters
- SQL Serverless / dedicated pool (stored procedures for ETL)
- Detailed Monitoring Spark Jobs for Data Flows
- Azure Synapse Link for SQL, Dataverse, Cosmos DB

Kurzer Blick: Azure Synapse Link for SQL



The screenshot shows the Azure Synapse Studio interface with the 'Integrate' blade selected. On the left, the navigation bar includes icons for Home, Datasets, Pipelines, Triggers, and Integration. The 'Link connection' section is expanded, showing a single entry: 'buildlink-onprem'. The main workspace displays a table mapping source tables from an 'Application' database to target tables in a 'OP_A...' database, all using 'Round robin' distribution. The 'Properties' pane on the right provides configuration options for the link connection, including the name 'buildlink-onprem', a description field, core count (set to 4 + 4 Driver cores), landing zone settings (linked service 'build2022lsynapse-WorkspaceD...', folder path 'landingzone', and SAS token), and a 'Rotate token' button.

Source table	Target table	Distribution Type	Distribution C...
Application.Cities	OP_A... · Cities	Round robin	-
Application.Co...	OP_A... · Countries	Round robin	-
Application.Deli...	OP_A... · DeliveryMethods	Round robin	-
Application.Pay...	OP_A... · PaymentMethods	Round robin	-
Application.Peo...	OP_A... · People	Round robin	-
Application.Stat...	OP_A... · StateProvinces	Round robin	-
Application.Sys...	OP_A... · SystemParamet...	Round robin	-

Properties

Link connection name *
buildlink-onprem

Description
Enter a description

Core count * ⓘ
4 (+ 4 Driver cores)

Landing zone settings

Landing zone linked service * ⓘ
build2022lsynapse-WorkspaceD...

Landing zone folder path * ⓘ
landingzone

Landing zone SAS token * ⓘ

Rotate token

Announcing Synapse Link for SQL: <https://www.youtube.com/watch?v=pgusZy34-Ek>

Better use Azure Data Factory for

- Small projects (copy a bunch of tables every day)
- To keep a small footprint in infrastructure
- If having only permissions in dedicated Azure Resource Group
- As orchestration tooling for Databricks projects
- Migration from M-Code (Power BI)

Better use Synapse Data Integration if

- You plan to do data exploration on a data lake using SQL or connect Power BI
- Data platform projects including many data sources
- Spark Notebooks used as part of data integration
- Want to use the newest features
- Want to sync masses of different tables via Azure Synapse Link .. (having money for dedicated SQL pools)

Neues in der Azure Data Factory

- Flowlets und Data Libraries
- Activity & Data Flow updates
- Application Lifecycle Management
- SAP CDC is coming

Flowlets

In Mapping Data Flows

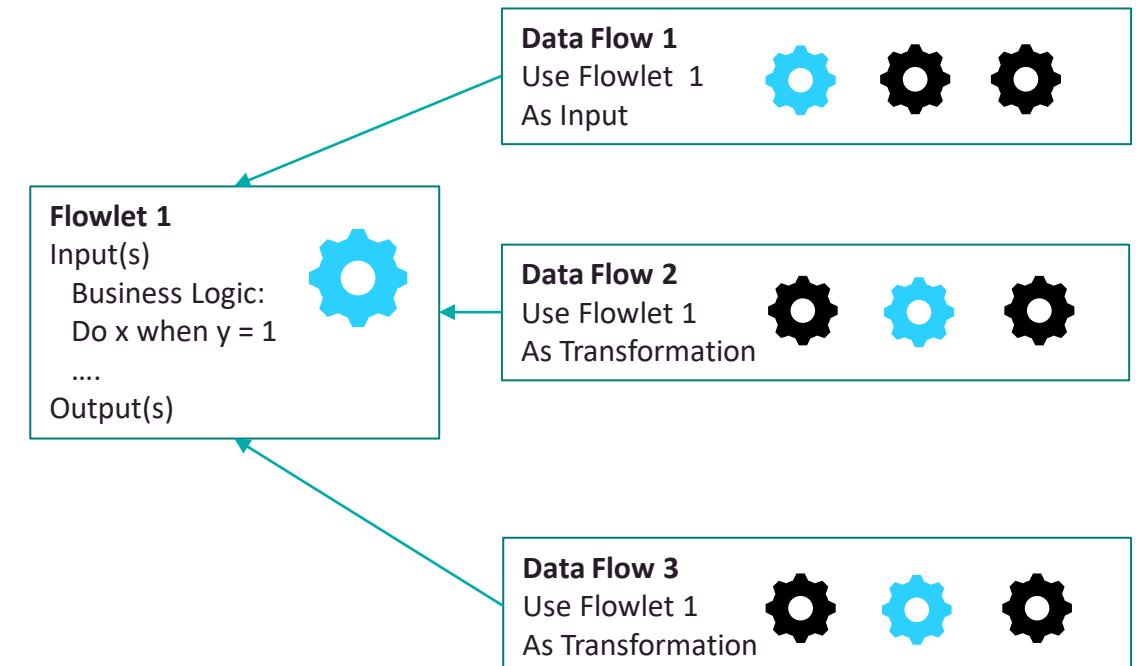
Flowlets in Mapping Data Flows

Wiederverwendbarer „Container“ von Aktivitäten

- **Don't Repeat Yourself Prinzip (DRY) :**

Durch Kapselung Duplikation von Code / Logik vermeiden

- Kann from scratch oder aus bestehenden Mapping Data Flows erstellt werden
- Einsatz als Quelle oder Transformation möglich
- Datentypen und Felder möglichst generisch und klein halten
- Schnittstellen-Vertrag zwischen Flowlets und den Data Flows die es verwenden





Demo

Flowlets



Hands-On

Flowlets

Data Flow Libraries

User Defined Functions

Data Flow Libraries - User defined functions

Noch mehr Kapselung – aber auf einer anderen Ebene

- Anlegen von Funktionen in der Expression Language die in allen Mapping Data Flows verwendet werden können
- Parametrisierung möglich
- Wiederverwendbarkeit von komplizierten Codeschnipseln
- DRY Prinzip

New data flow Library

Name *	DemoUDF
Description	Description

Functions

+ New	Name	Return Type	Body
No functions in data flow library			

New data flow function

Name *	Name
Arguments	+ Add

This function has no arguments.

Body *	Enter expression...	ANY
--------	---------------------	-----

Data Flow Libraries - User defined functions

Definition von User Defined Functions mit Parametern

- Parameter als Arguments hinzufügen
- Namen sind vorgegeben als i1, i2 ...
- Im bekannten Expression Editor entwickeln
- Einbinden dann im Data Flow im Expression Editor z.B. bei Derived Column unter Data flow library functions

Edit data flow function

Name *	genderFinder	
Arguments	+ Add	
Name	Type	
i1	string	[Delete]

Body *

```
SELECT CASE WHEN i1 = 'M' THEN 'Male' WHEN i1 = 'F' THEN 'Female' ELSE 'Unknown' END AS Gender
```

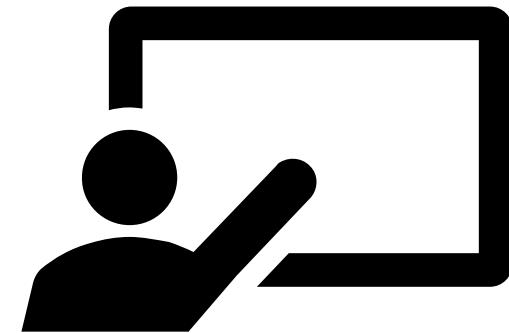
Expression elements

[All](#)
[Functions](#)
[Input schema](#)
[Parameters](#)
[Cached lookup](#)
[Data flow library functions \(preview\)](#)
[Locals](#)

Expression values

[Filter by keyword](#)

abc genderFinder(abc string)



Demo

Data Flow Libraries – User Defined Functions

Activity & Data Flow updates

Script Activity

Besser mit Datenbanken können

- Kann einiges mehr als Lookup und Stored Procedure Activity (siehe rechts)
- Ein oder mehrere Statements
- Auch Rückgabe von Result sets sowie Abfrage-Meldungen

	Script Activity	Lookup Activity	SProc Activity
<i>Supported data source</i>	Database (SQL family, Snowflake, Oracle)	all data sources	SQL family
<i>Supported Operations</i>	Read / Modify	Read*	Modify
<i>Multiple query support</i>	Yes	No**	Yes***
<i>Query parameter support</i>	Input / Output	Not supported	Input
<i>Output Result set support</i>	One or more	One	No
<i>Output query logs (PRINT)</i>	Yes	No	No
<i>Integrated CICD (ADF)</i>	Yes	Yes	No



Demo

Script Task

Assert Transformation

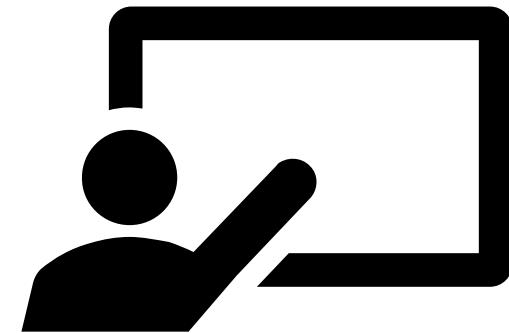
Datenqualität inhaltlich checken im Data Flow

- Erstellen von Regeln die auf die Spalten im Flow angewendet werden
- Checken auf
 - Wertebereiche (Expect true)
 - Eindeutigkeit (Expect unique)
 - Vergleich zwischen Streams (Expect exists)
- Markieren der Zeilen oder Abbruch möglich

<input type="checkbox"/> Assert type	Assert Id	Assert description	Filter	Expression	Ignore nulls
<input type="checkbox"/> Expect true	assert1	ANY	ANY	ANY	<input type="checkbox"/>

A dropdown menu is open under the 'Assert type' column, showing the following options:

- Expect true (selected)
- Expect true
- Expect unique
- Expect exists



Demo

Assert Transformation

Rest Connector für Mapping Data Flows

Als Source und Sink einsetzbar

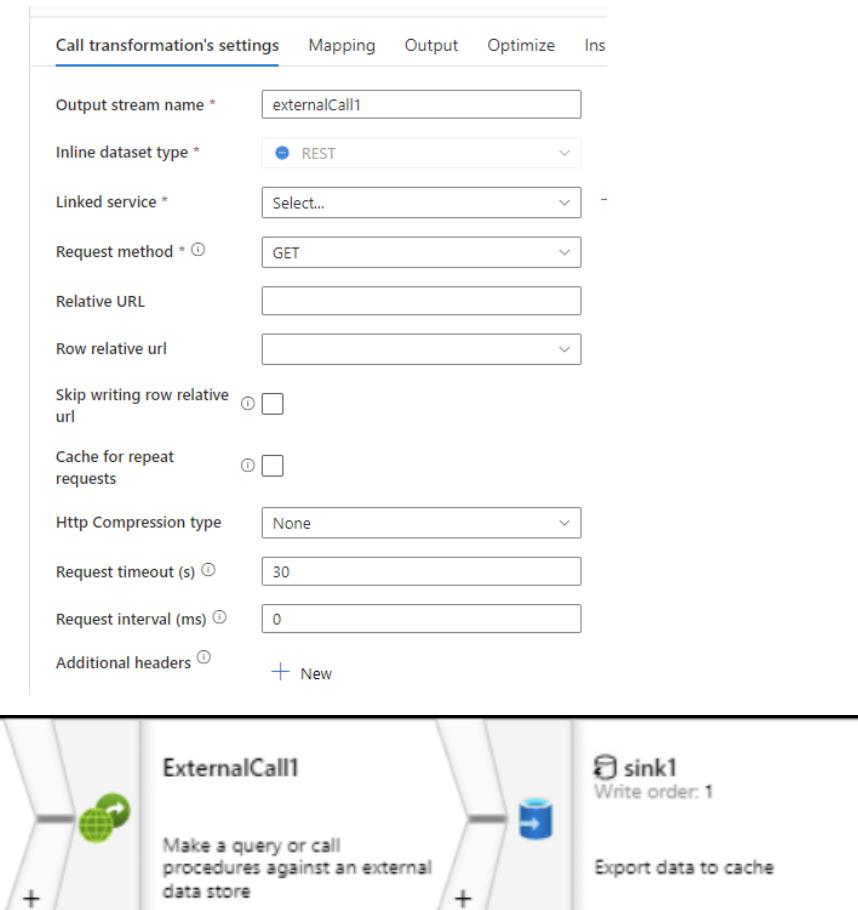
- requestMethod HTTP, Get & Post erlaubt
- integration datasets und inline datasets
- Relativ breiter Pagination Support
- Bei Verwendung als Sink ein Alter Row davor schalten um die Action zu bestimmen

Sink	Settings	Mapping	Optimize	Inspect	Data preview
	<u>Insert method</u>	PUT			
	<u>Delete method</u>	DELETE			
	<u>Upsert method</u>	PUT			
	<u>Update method</u>	PATCH			
	<u>Row relative url</u>	{ } periods			
	<u>Skip writing row relative url</u>	<input type="checkbox"/>			
	<u>Http Compression type</u>	None			
	<u>Batch size</u> ⓘ				
	<u>Request timeout (s)</u> ⓘ	30			
	<u>Request interval (ms)</u> ⓘ	0			
	<u>Additional headers</u> ⓘ	+ New			

External Call Transformation

Auch innerhalb eines Data Flows direkt Rest Services nutzen

- Einbinden externer Services im Zeilenmodus, z.B. Wetterdaten pro GPS Koordinate der Zeilen im Data Flow
- Wird noch erweitert auf Stored Procedures



External Call Transformation

Am Beispiel Adressvalidierung

PostalCode
98011
98011
75201
85004
H1Y 2H5
98004



Column	Expression
<input type="checkbox"/> ziplookup	"/{PostalCode}/degrees"

External Call Transformation

Konfiguration



Call transformation's settings

Output stream name *	ZipLookup
Inline dataset type *	REST
Linked service *	RestService6
Request method *	GET
Relative URL	<input type="text"/>
Row relative url	abc ziplookup
Skip writing row relative url	<input checked="" type="checkbox"/>
Include header in output	<input checked="" type="checkbox"/>
Http Compression type	None
Request timeout (s) *	30

Call transformation's settings

Options	<input type="checkbox"/> Skip duplicate input columns				
	<input type="checkbox"/> Skip duplicate output columns				
Input columns *	<input checked="" type="checkbox"/> Auto mapping <input type="button" value="Reset"/> <input type="button" value="Add mapping"/> <input type="button" value="Delete"/> <table border="1"> <tr> <td><input type="checkbox"/> createURL's column</td> <td><input type="button" value="Name as"/></td> </tr> <tr> <td><input type="checkbox"/> abc ziplookup</td> <td><input type="button" value="ziplookup"/></td> </tr> </table>	<input type="checkbox"/> createURL's column	<input type="button" value="Name as"/>	<input type="checkbox"/> abc ziplookup	<input type="button" value="ziplookup"/>
<input type="checkbox"/> createURL's column	<input type="button" value="Name as"/>				
<input type="checkbox"/> abc ziplookup	<input type="button" value="ziplookup"/>				

Call transformation's settings

Import projection	Schema options
-------------------	----------------

Output

Response body definition
`(acceptable_city_names as string[], area_codes as string... {})`

External Call Transformation

Unterstützte Rest Request Methoden

Call transformation's settings Mapping Output Optimize 

Output stream name * ZipLookup

Inline dataset type * REST

Linked service * RestService6

Request method * GET 

Relative URL

Row relative url

Skip writing row relative url 

GET
POST
PATCH
PUT
DELETE

Achtung wird auch beim Preview
ausgeführt (z.B. Delete)

Power Query

Aka Wrangling Data Flows

- Umbenannt, jetzt als „Power Query“ in ADF
- Wird Stück für Stück vervollständigt, Umstellung von Power BI Power Query bedeutet ja Spark hier im Backend
- Fühlt sich für mich trotzdem immer wie ein Fremdkörper an weil auch das Handling von Datasets anders ist usw.

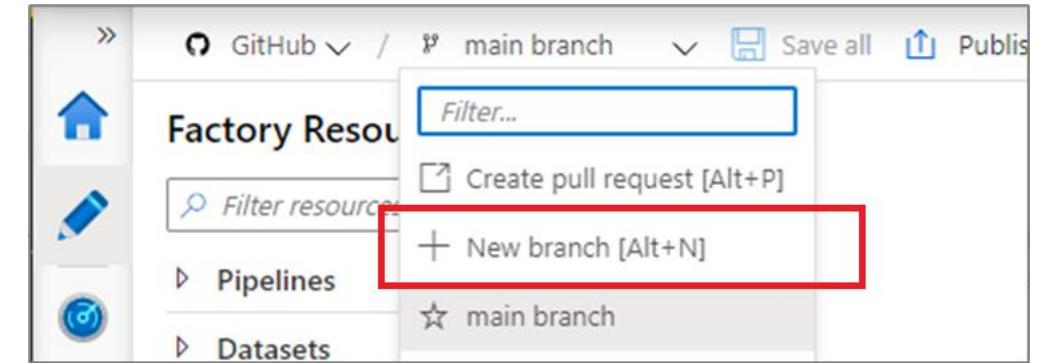
Application Lifecycle Management News

Build to survive

Azure Data Factory Github Integration

Limits entfallen & Branch of Branch

- Jetzt auch mehr als 1000 Data Factory Resourcen pro resource type möglich (gemeint sind datasets, pipelines & co.)
- Hilft bei der Arbeit mit Code Versioning & Deployments auch in Zusammenhang mit der GitHub API
- Es kann jetzt auch ein Branch basierend auf einem anderen Branch erstellt werden



Create a new branch

Branch name *

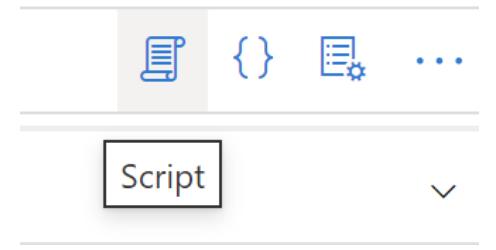
Base on *

 collab1 branch

Intro Data Flow Script (DFS)

Programmiersprache hinter den Mapping Data Flows

- Data Flow Script auch Teil des Json Codes (ganz unten)
- Entsteht aber beides automatisch durch die grafische Oberfläche
- Editierbar
- Inhalte übertragbar in andere Data Flows / Projekte
- Verwenden in Powershell cmdlets
- Einfaches Find & Replace im Text
- Programmatisch Data Flows erzeugen aus Metadaten
- Einfacheres Finden von Abweichungen zwischen Versionen – vereinfacht auch in Json seit der Einführung von scriptline
- Viele Beispiele von MS oder Blog-Artikeln können einfach integriert werden



```
1 source(output(
2     medallion as string,
3     hack_license as string,
4     vendor_id as string,
5     rate_code as string,
6     store_and_fwd_flag as string,
7     pickup_datetime as string,
8     dropoff_datetime as string,
9     passenger_count as short,
L0     trip_time_in_secs as long,
L1     trip_distance as double,
L2     pickup_longitude as double,
L3     pickup_latitude as double,
L4     dropoff_longitude as double,
L5     dropoff_latitude as double
L6 ),
L7     allowSchemaDrift: false,
L8     validateSchema: false,
L9     ignoreNoFilesFound: false) ~> TripData
source(output(
    medallion as string,
    { hack_license} as string,
    { vendor_id} as string,
```

Intro Data Flow Script (DFS)

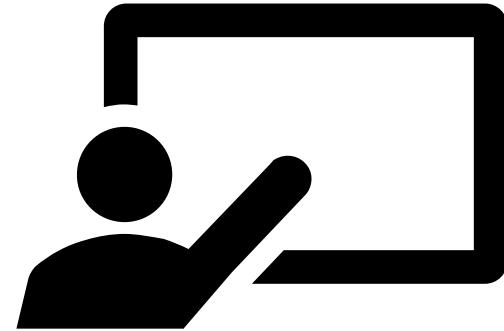
Scriptline in Json Code Definition

- JSON Definition in ADF für Data Flow früher: c

```
|   |   |
|   |   |   }
|   |   ],
|   |   "script": "source(output(\n\t\tmedallion as string,\n\t\tthack_license as string,\n\t\tvendor_id as string,\n\t\trate_code as string,\n\t\tstore_and_fwd_fl
```

- JSON Definition in ADF für Data Flow seit Einführung scriptline:

```
],
"scriptLines": [
    "source(output(",
    "    {Emp ID} as string," ,
    "    {Name Prefix} as string," ,
    "    {First Name} as string," ,
    "    {Middle Initial} as string," ,
    "    {Last Name} as string," ,
    "    Gender as string," ,
    "    {E Mail} as string," ,
    "    {Father's Name} as string," ,
    "    {Mother's Name} as string," ,
    "    {Mother's Maiden Name} as string," ,
    "    {Date of Birth} as string," ,
    "    {Time of Birth} as string," ,
```



Demo

Data Flow Script

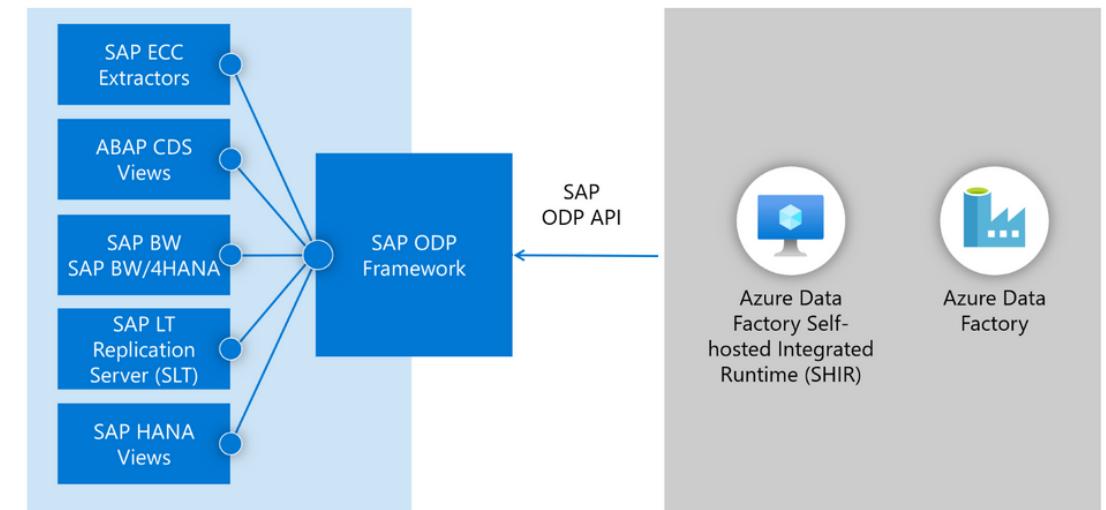
SAP CDC is coming!

I want it all (not) – incremental sourcing

SAP CDC

Verwenden von SAP eigener Change Tracking Mechanismen für Delta Loads

- new SAP CDC connector leveraging SAP ODP framework
- can connect to all SAP systems that support ODP, such as R/3, ECC, S/4HANA, BW, and BW/4HANA
- directly at the application layer or indirectly using SAP Landscape Transformation (SLT) replication server as a proxy
- fully or incrementally extract SAP data
- physical tables & logical objects such as ABAP Core Data Services (CDS) views
- will be released for public preview on June 30, 2022





Hands-On

Auf geht's Leute!