Kamil Nowiński

PASS
SQLSATURDAY
SLOVENIA | 08 DEC 2018

# Azure Data Factory:
# Data Flows – first blood

# Thank you to our AWESOME sponsors!

# About me

## Kamil Nowinski

Microsoft CERTIFIED
Solutions Associate
SQL Server 2012

MCSE Data Platform

MVP Microsoft® Most Valuable Professional

Microsoft Data Platform **MVP**
Speaker, blogger, data enthusiast
Senior Data Engineer at ASOS (www.asos.com)
13+ yrs experience as DEV/DBA
Member of the Data Community PL (SQLDay)
Project member of „SCD Merge Wizard"
Founder of blog SQLPlayer (www.SQLplayer.net)

SQL Server Certificates:
MCITP, MCP, MCTS, MCSA, MCSE Data Platform,
MCSE Data Management & Analytics
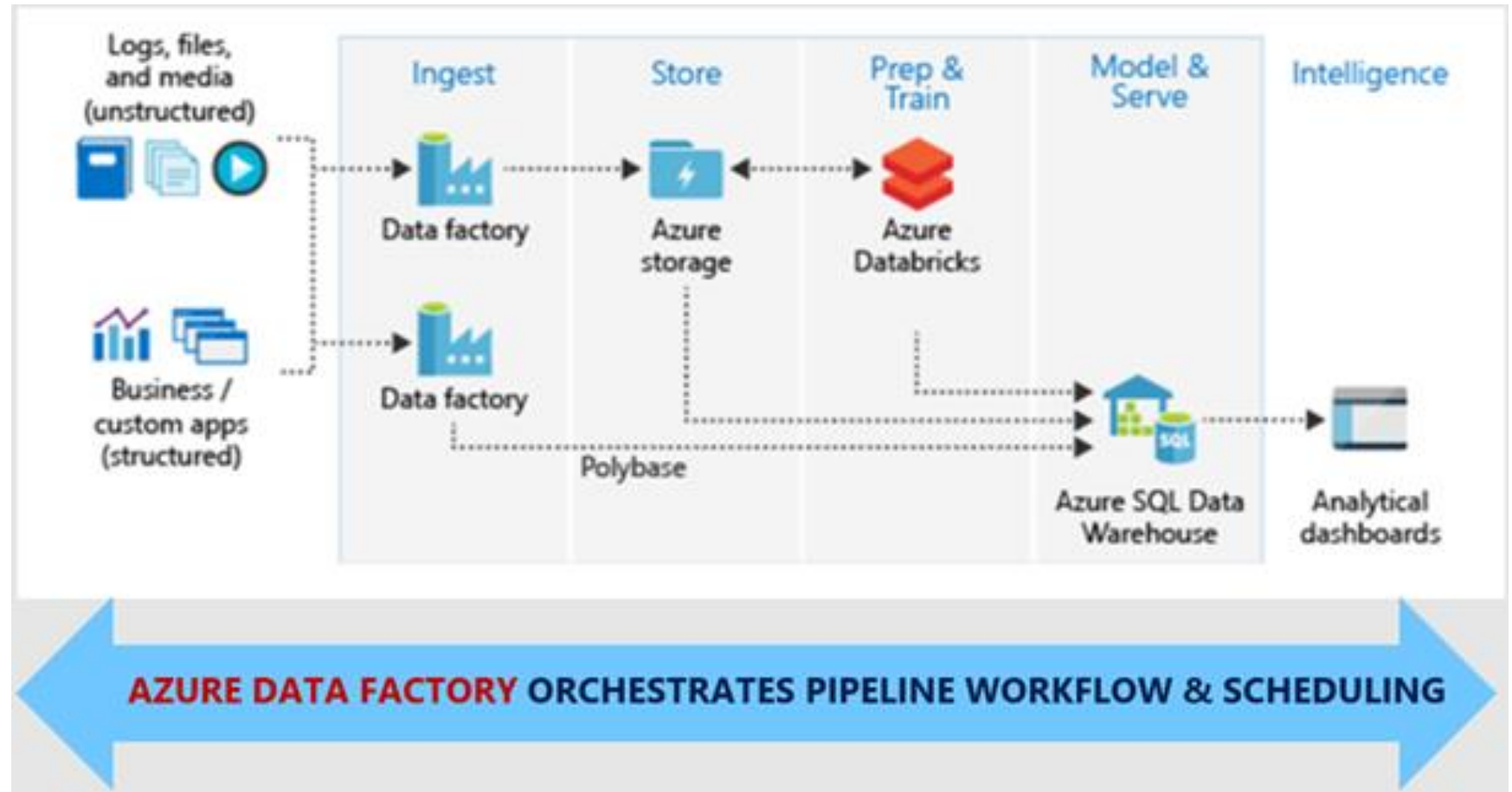Moreover: Bicycle, Running, Digital photography
@NowinskiK, @SQLPlayer

redgate FRIENDS

PASS
SQLSATURDAY
SLOVENIA | 08 DEC 2018

# (Almost) BRAND NEW BLOG

INTERVIEWS

SQLPlayer

Play with data & have fun!

www.SQLPlayer.net

# PODCAST – Interviews with…



www.SQLPlayer.net

# What the Azure Data Factory is?



Logs, files, and media (unstructured)

Business / custom apps (structured)

Ingest — Data factory

Store — Azure storage

Prep & Train — Azure Databricks

Model & Serve — Azure SQL Data Warehouse

Intelligence — Analytical dashboards

Polybase

**AZURE DATA FACTORY ORCHESTRATES PIPELINE WORKFLOW & SCHEDULING**

# Access all your data:

- 75+ connectors & growing
- Azure IR available in 20 regions
- Hybrid connectivity using self-hosted IR: on-prem & VNet

| Azure (13) | | Database (24) | | File Storage (5) | NoSQL (3) | Services and Apps (28) | | Generic (4) |
|---|---|---|---|---|---|---|---|---|
| Blob Storage | Amazon Redshift | Netezza | Amazon S3 | Cassandra | Amazon MWS | Office 365 * | HTTP |
| Cosmos DB (MongoDB API) * | DB2 | Oracle | File System | Couchbase | CDS for Apps | Paypal | OData |
| Cosmos DB (SQL API) | Drill | Phoenix | FTP | MongoDB | Concur | QuickBooks | ODBC |
| Data Lake Storage Gen1 | Google BigQuery | PostgreSQL | HDFS | | Dynamics 365 | Salesforce | REST * |
| Data Lake Storage Gen2 | Greenplum | Presto | SFTP | | Dynamics CRM | Salesforce Marketing Cloud | |
| DB for MySQL | HBase | SAP BW | | | GE Historian | Salesforce Service Cloud | |
| DB for PostgreSQL | Hive | SAP HANA | | | Google AdWords | SAP C4C | |
| File Storage | Impala | Spark | | | HubSpot | SAP ECC | |
| Kusto * | Informix | SQL Server | | | Jira | ServiceNow | |
| Search Index | MariaDB | Sybase | | | Magento | Shopify | |
| SQL DB | Microsoft Access | Teradata | | | Marketo | Square | |
| SQL DW | MySQL | Vertica | | | Oracle Eloqua | Web table | |
| Table Storage | | | | | Oracle Responsys | Xero | |
| | | | | | Oracle Service Cloud | Zoho | |

Supported as Source and Sink

Supported as Source only

Supported as Sink only

# ADF Key Concepts

# How to create Azure Data Factory with Data Flow?



Home > Data factories > New data factory

**Data factories**
ASOS.com Ltd

+ Add     Edit columns     ••• More

Filter by name...

NAME ↑↓

bigfactory555

BigFactoryCDM

BigFactoryDF

BigFactoryPP

BigPlayer

**New data factory**

\* Name ⓘ
MyDataFactory666 ✓

\* Subscription
Visual Studio Enterprise ⌄

\* Resource Group ⓘ
○ Create new   ● Use existing
rg-datafactory ⌄

Version ⓘ
V2 with data flow (preview) ⌄

\* Location ⓘ
Southeast Asia ⌄

Include data flow sample
☐

ADF Data Flow is currently in private preview. The normal Azure SLAs do not apply to use of this preview feature and all support must route through this email address:
adfdataflowext@microsoft.com

Visual Data Transformations with **Data Flow**

# What the hell Data Flows are?



| Data Sources | Staging | Transformations | Destination |
|---|---|---|---|
| • Explicit user action<br>• User places data source(s) on design surface, from toolbox<br>• Select explicit sources | • Implicit/Explicit<br>• Data Lake staging area as default<br>• User does not need to configure this manually<br>• Advanced feature to set staging area options<br>• File Formats / Types (Parquet, JSON, txt, CSV ...) | • Explicit user action<br>• User places transformations on design surface, from toolbox<br>• User must set properties for transformation steps and step connectors | • Explicit user action<br>• User chooses destination connector(s)<br>• User sets connector property options |

Transformations: Sort, Merge, Join, Lookup ...

# Authoring of Azure Data Factory (v2) – what's new?

# Code-free Data Transformation At Scale

Does not require understanding of Spark, Big Data Execution Engines, Clusters, Scala ...

Focus on building business logic and data transformation

Data cleansing

Aggregation

Data conversions

Data prep

Data exploration

ETL Data Loading into DW



... **not**

# Authoring of Azure Data Factory (v2)

# Simple Copy Flow

# Guided experience to build data flows

# Data Flow: Components = Actions *

**Multiple inputs/outputs**

- New Branch
- Join
- Conditional Split
- Union
- Lookup

**Schema modifier**

- Derived Column
- Aggregate
- Surrogate Key
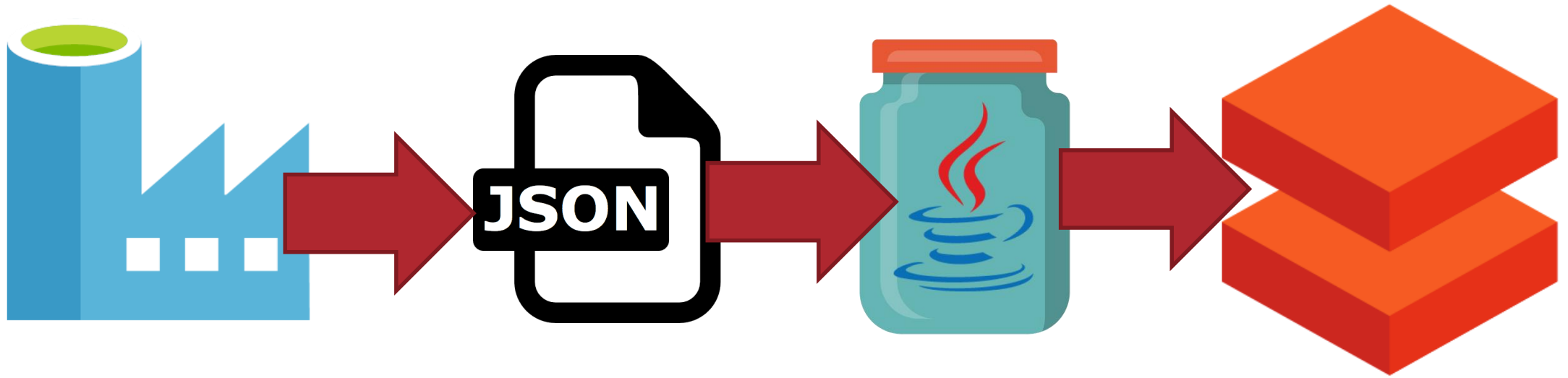
**Row modifier**

- Exists
- Select
- Filter
- Sort

**Custom**

- Extend

**Destination**

- Sink

\* With some small exceptions

# What is going on behind the scenes?



JAR

Azure
Databricks

# Azure Databricks version

# DEMO TIME

# Debug mode provides row-level context and visible results in inspector pane

# Data Flow Execution Plan

# Data Flow Data Lineage

# Resources

- Microsoft Azure Data Factory – [Tutorials & API Reference](#)
- Azure Data Factory [Overview](#)
- Azure Data Factory – [Data integration service](#)
- ADF Data Flow's [documentation](#)
- ADF Data Flow's [videos](#)
- SQLPlayer blog:
  - [Azure Data Factory v2 and its available components in Data Flows](#)
  - Follow this tag on SQLPlayer blog: [ADFDF](#)

# Q&A

# Thank you!    Hvala!    Dziękuję!

📧 [kamil@nowinski.net](mailto:kamil@nowinski.net)

🐦 @NowinskiK     @SQLPlayer

🔗 SQLPlayer.net

**Kamil Nowinski**
Microsoft Data Platform MVP
MCSE Data Platform & MCSE Data Management and Analytics

# Thank you to our AWESOME sponsors!