

PLATINUM SPONSOR

STRATEGIC PARTNER

TECHNOLOGY
INNOVATION
DATA
KNOWLEDGE

GOLD SPONSORS



CLOUDS ON MARS



SILVER SPONSOR



BRONZE SPONSOR



ADF Mapping Data Flow

Tomasz Libera

Tomasz Libera

- Microsoft MVP Data Platform
- Microsoft Certified Trainer
- SQL Server Developer, Trener akademicki
 - WSZiB w Krakowie
- Prowadzi autoryzowane i autorskie szkolenia
 - TSQL | Stored Procedures | Performance Tuning
 - Integration Services | Reporting Services | Power BI
- Data Community
 - Lider Krakowskiej Grupy | Były Członek Zarządu
 - Organizator i prelegent SQLDay, SQLSaturday
- Pasjonat kolarstwa górskiego i maratonów MTB
 - XBOX Bike Team
- tomasz.libera@datacommunity.pl | blog.libera.net.pl



Agenda

- Azure Data Factory Overview
- Mapping Data Flow
- Transformations
- Demo

Azure Data Factory

- The Azure Data Factory service is a fully managed service for composing data storage, processing, and movement services into streamlined, scalable, and reliable data production pipelines

Azure Data Factory



- Productive
 - Build automated data integration solutions with visual drag-&-drop UI. Move data seamlessly from over **80 sources** without writing code.



- Trusted
 - Data movement using Azure Data Factory has been certified by HIPAA/HITECH, ISO/IEC 27001, ISO/IEC 27018, and CSA STAR.



- Hybrid
 - Build data integration pipelines that span on-premises and cloud. Easily lift your SQL Server Integration Services (SSIS) packages to Azure.



- Scalable
 - Build serverless cloud-based data integration with no infrastructure to manage. Take advantage of elastic capabilities to scale out with your customer growth.

Azure Data Factory



- Visual drag-&-drop UI
 - Maximize productivity by getting pipelines up and running quickly. Use the code-free drag-&-drop interface to build, deploy, monitor, and manage your data integration



- SSIS package execution in Azure
 - Easily execute and schedule your SQL Server Integration Services (SSIS) packages in managed execution environment



- Comprehensive control flow
 - Looping, branching, conditional constructs, on-demand executions, and flexible scheduling



- Multiple Language Support
 - Use the visual interface or write your own code in Python, .NET, or ARM to build pipelines using your existing skills



- Code-free data movement
 - Improve your TCO with 80+ natively supported connectors including Azure data services

New ADF service

Home > New > New data factory

New data factory

* Name ⓘ
MyAzureDataFactory ✓

* Subscription
Microsoft Azure Sponsorship ▼

* Resource Group ⓘ
☐ Create new ☒ Use existing
MyRG ▼

Version ⓘ
V2 (Preview) ▼

* Location ⓘ
East US ▼

Home > MyAzureDataFactory > Settings

MyAzureDataFactory
Data factory

✦ ✕

🗑 Delete

Essentials ^

Resource group
MyRG

Type
Data factory (V2)

Location
EastUS

Subscription name
Microsoft Azure Sponsorship

Provisioning state
Succeeded

Subscription id
14275120-0b96-4eef-98fa-b54e9272fee4

Getting started
Quick start

All settings →

Quick links

📖 Documentation

✍ Author & Monitor

Monitoring

PipelineRuns

100

Settings
MyAzureDataFactory

🗑 ✕

🔍 Filter settings

SUPPORT + TROUBLESHOOTING

🔧 Diagnose and solve problems >

📋 Activity log >

🛎 New support request >

GENERAL

☰ Properties >

GETTING STARTED


☁ Quick start >

RESOURCE MANAGEMENT

🏷 Tags >

🔒 Locks >

👤 Users >



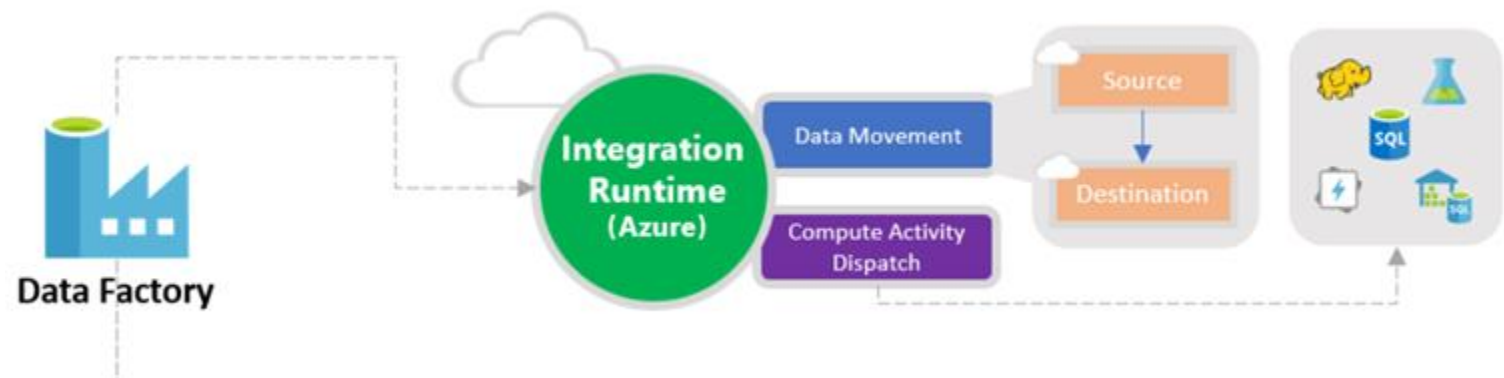
Data Community

Components

- **Pipeline** – key concept in ADF; workflow of activities
- **Activity** - processing step in a pipeline (copy data, run integration job)
- **Datasets** – data structures
 - point to the data you want to use as inputs or outputs
- **Linked service** – connection string
 - Data store - Azure/ on-premise/ Oracle/ Azure blob storage account
 - Compute resource - execution of an activity – HDInsight Hive/ Pig/ MapReduce/ Streaming
- Triggers - unit of processing
- Parameters

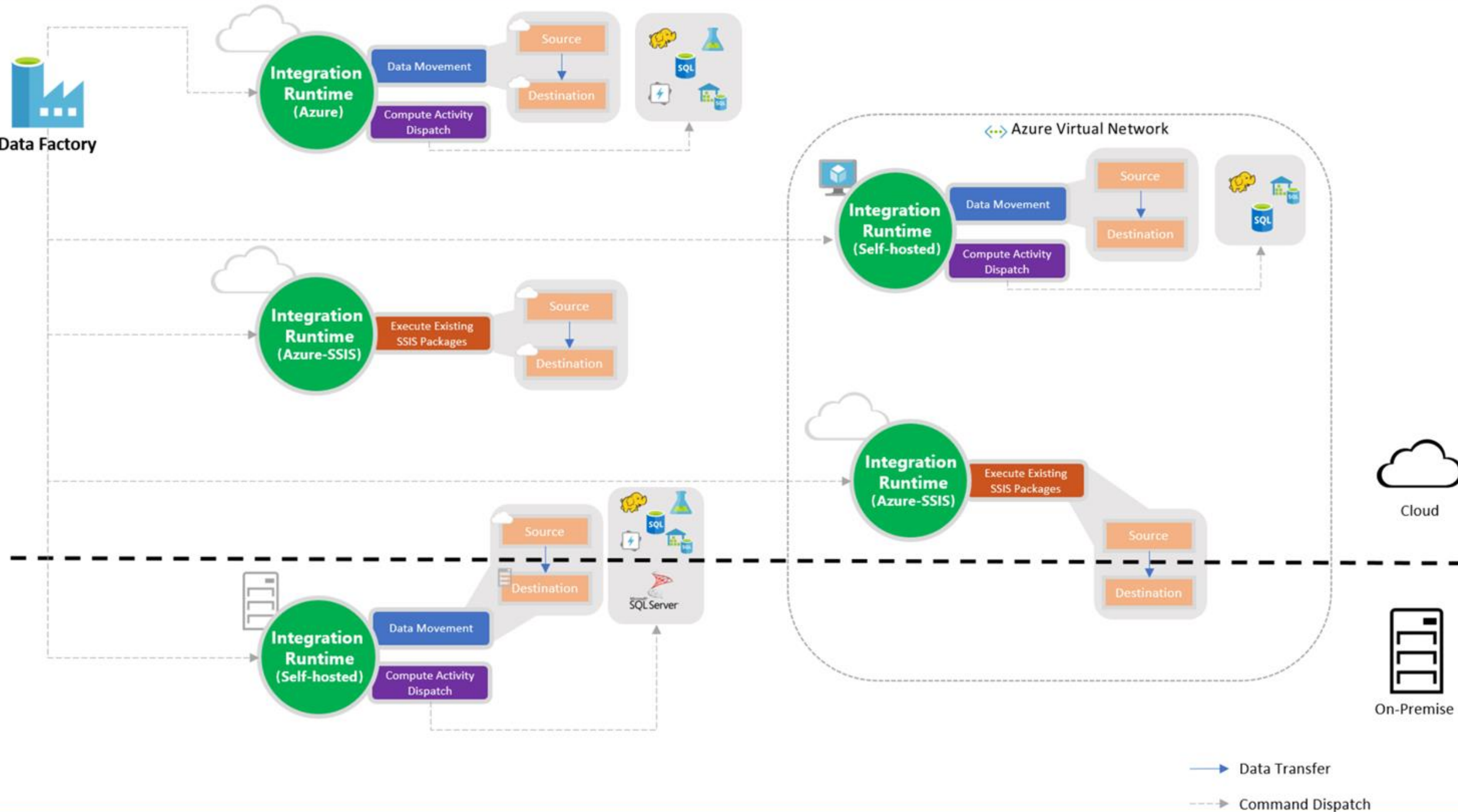
Integration Runtime (IR)

- Data movement
- Activity dispatch - monitor transformation activities
- SSIS package execution: Natively execute SSIS packages
- Types:
 - Azure
 - Self-hosted
 - Azure-SSIS





Data Factory



Azure Data Factory

ADF MAPPING DATA FLOW

Mapping Data Flow













- Transformation capabilities for Azure Data Factory.
Behind the scene ADF JSON code (generated by drag&drop interface) is converted into Scala programming language and compiled and executed in Azure Databricks which is automatically scale-out as needed.
- Does not require understanding of Spark, Big Data Executions Engines, Clusters, Scala, Python, Java... (Zero-Code)
- Focus on building business logic and data transformations
 - Data cleansing
 - Aggregation
 - Data conversion
 - Data preparation
 - Data exploration












Data Flows vs SSIS



SQL Player

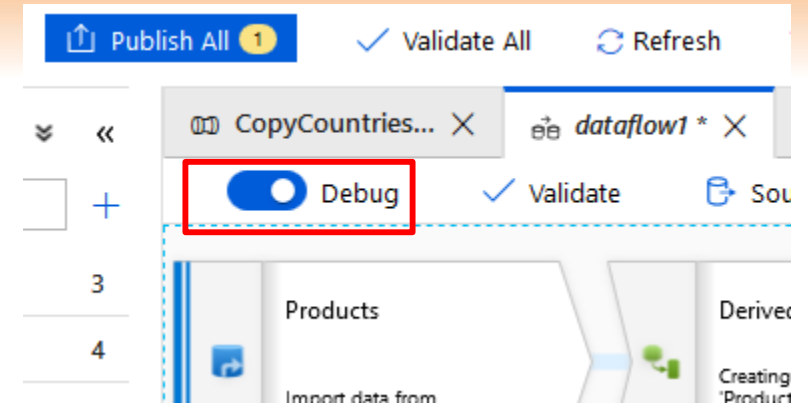
Play with data & have fun!

Operation / Activity	Description	SSIS equivalent	SQL Server equivalent
 New branch	Create a new flow branch with the same data	 Multicast (+icon)	<pre>1 SELECT INTO 2 SELECT OUTPUT</pre>
 Join	Join data from two streams based on a condition	 Merge join	<pre>1 INNER/LEFT/RIGHT JOIN, 2 CROSS/FULL OUTER JOIN</pre>
 Conditional Split	Route data into different streams based on conditions	 Conditional Split	<pre>SELECT INTO WHERE condition1 SELECT INTO WHERE condition2 CASE ... WHEN</pre>
 Union	Collect data from multiple streams	 Union All	<pre>SELECT colla UNION (ALL) SELECT collb</pre>
 Lookup	Lookup additional data from another stream	 Lookup	<i>Subselect, function,</i> <pre>LEFT/RIGHT JOIN</pre>
 Derived Column	Compute new columns based on the existing once	 Derived Column	<pre>SELECT Column1 * 1.09 as NewColumn</pre>

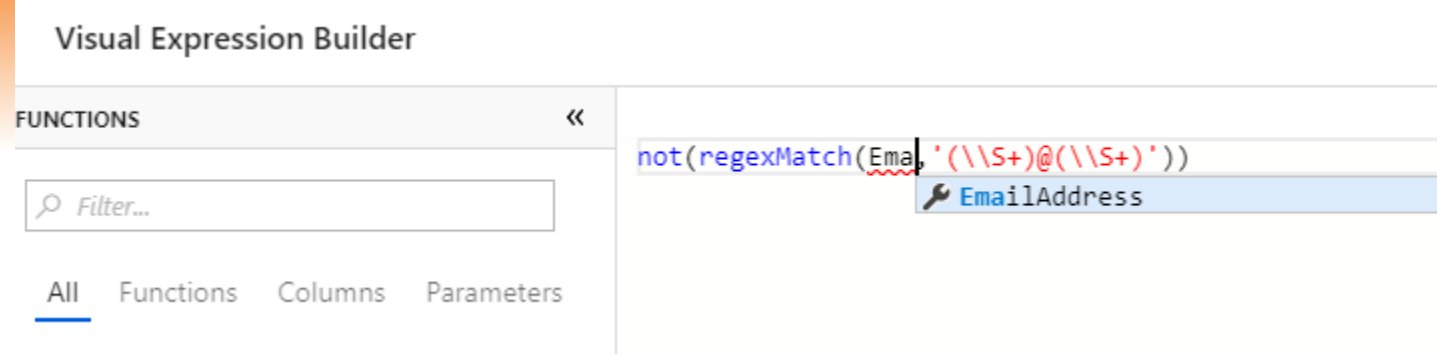
Exists			Choose columns to flow to the next stream	OUTPUT in components, mapping columns	<pre>SELECT Column1, Column4 FROM Table</pre>
Select			Filter rows in the stream based on a condition	 Conditional Split	<pre>SELECT * FROM Table WHERE [Column] LIKE '%pattern%'</pre>
Sort			Order data in the stream based on column(s)	 Sort	<pre>SELECT * FROM Table ORDER BY [Column] ASC</pre>
Extend			Use any custom logic from an external library	 Script Component	<i>SQL CLR</i>
Source			Source for your data flow. Obligatory first element of every Data Flow in ADF.	 OLE DB Source and more ...	<pre>SELECT * FROM SourceTable</pre>
Sink			Destination for your data flow	 OLE DB Destination and more...	<pre>INSERT INTO TargetTable</pre>

Debug Mode

- Interactively build data flow with a running Azure Databricks interactive cluster
- Session will close once you turn debug off in ADF.
- Hourly charges incurred by Azure Databricks during the time that the debug session turned on
- <https://github.com/kromerm/adfdataflowdocs/blob/master/Concepts/ADF-data-flow-debug-mode.md>



Expression Builder



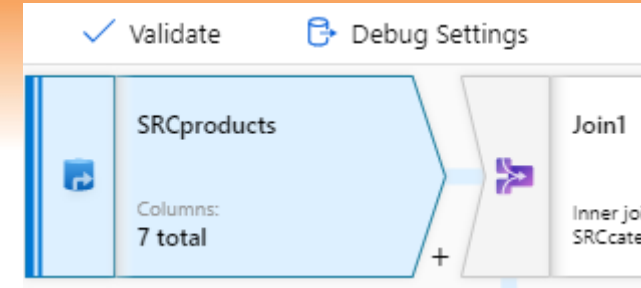
- Some transformations use expressions which could be created in Expression Builder.
- Expressions use columns, fields, variables, parameters and functions.
- Auto-complete feature reads from the entire Azure Data Factory Data Flow object model with syntax checking and highlighting.
- Using Debug mode, live in-progress preview data results and real-time live debugging is enabled.

Data preview	
Output: 	EmailAddress ^{abc}
	rucy0adventure-works.com
	rosmarie0adventure-works.com
	dominic0adventure-works.com
	kathleen0@adventure-works.com
	

Transformations

Source

- Data flow can include more than one source.
- Every data flow requires at least one source transformation.
- You can join those sources together with a join transformation or a union transformation.
- Dataset defines the shape and location of the data you want to write to or read from.
- You can use wildcards and file lists in your source to work with more than one file at a time.

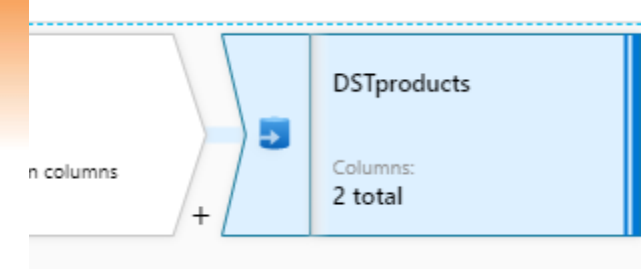


The screenshot shows the 'Source Options' tab of a data transformation tool. The 'Output stream name' is 'SRCproducts'. The 'Source dataset' is 'DS_SQLDbNoTables'. The 'Options' section has 'Allow schema drift' checked and 'Validate schema' unchecked. The 'Sampling' section has 'Disable' selected. The 'Input' section has 'Query' selected. The 'Query' text area contains a SQL statement: 'SELECT ProductID, Name, Color, StandardCost, ListPrice, Size, ProductCategoryID FROM SalesLT.Product'. There is an 'Import schema' button. The 'Batch size' field is empty.

Transformations

Sink

- Data flow can include more than one sink (destination).
- File name options
 - Default
 - Pattern
 - Per partition
 - As data in column
- Database options
 - Update method
 - Recreate table
 - Truncate table
 - Batch size
 - Enable staging



The screenshot displays the Databricks UI for configuring a Sink. It shows two panels, one for 'DSTbikes' and one for 'DSTproducts', both with the 'Settings' tab selected.

DSTbikes Settings:

- Output stream name: DSTbikes
- Incoming stream: BikesAndNotBikes@Bikes
- Sink dataset: DS_DelimitedText
- Options:
 - ☒ Allow schema drift
 - ☐ Validate schema
- File name option: ☐ Default ☐ Pattern ☐ Per partition ☐ As data in column ☒ Output to single file
- File name: bikes.csv

DSTproducts Settings:

- Output stream name: DSTproducts
- Incoming stream: Sort1
- Sink dataset: DS_SQLDb_ProductsCnt
- Options:
 - ☒ Allow schema drift
 - ☐ Validate schema
- Update method:
 - ☒ Allow insert
 - ☐ Allow delete
 - ☐ Allow upsert
 - ☐ Allow update
- Table action: ☐ None ☐ Recreate table ☒ Truncate table
- Batch size: [Empty field]

Transformations

Alter Row

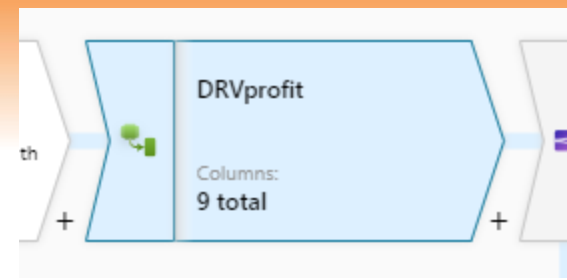
- To set insert, delete, update, and upsert policies on rows.
- Based on expressions (conditions) rows will be inserted, updated, deleted, or upsert
- Alter Row can produce both DDL & DML actions against database.

The screenshot shows a data pipeline with four stages: SRCcustomer (Import data from DS_SQLDB_Customer), DeleteIncorrectEmails (Columns: 15 total), DerivedColumn1 (Creating/updating the columns 'CustomerId, NameStyle, Title, FirstName, MiddleName, LastName, Suffix'), and DSTcustomer (Export data to DS_SQLDB_Customer). The 'Alter Row Settings' tab is active, showing the following configuration:

- Output stream name: DeleteIncorrectEmails
- Incoming stream: SRCcustomer
- Alter row conditions:
 - Delete if: `not(regexMatch(EmailAddress, '(\S+)(\S+)'))`
 - Update if: `year(ModifiedDate) = 2005`

Transformations

Derived Column



- Generate new columns or to modify existing fields.
- Expression Builder window to build the expression for the derived columns using expression functions.

Derived Column's Settings Optimize Inspect Data Preview

Output stream name *

Incoming stream * [Select1](#)

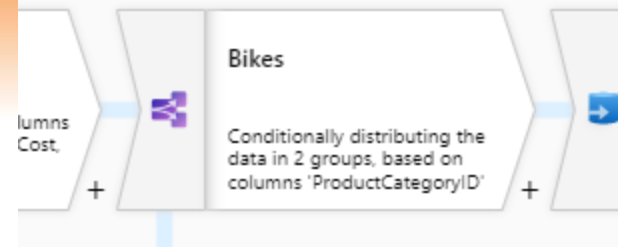
Columns *

Visual Expression Builder [Expression reference doc](#)

COLUMN SCHEMA	FUNCTIONS	EXPRESSION
Profit	<div><input type="text" value="Filter..."/></div> <div><div>All Functions Columns Parameters</div><div><div>123 ProductID</div><div>abc Name</div></div></div>	<div><code>toDecimal(ListPrice-Cos, 10,2)</code></div> <div><div>Cost</div><div>cos</div><div>cosh</div><div>concatWS</div></div>

Transformations

Conditional Split



- Generate new columns or to modify existing fields.
- Expression Builder window to build the expression for the derived columns using expression functions.

Conditional Split Settings Optimize Inspect Data Preview

Output stream name *

Incoming stream * [DRVprofit](#)

Split on ☒ First matching condition ☐ All matching conditions

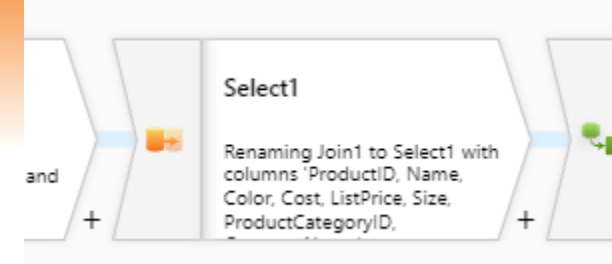
Split condition

STREAM NAMES	CONDITION
<input type="text" value="Bikes"/>	<input type="text" value="and(ProductCategoryID >= 5, ProductCategoryID <= 7)"/>
<input type="text" value="BikesNot"/>	Rows that do not meet any condition will use this output stream

Transformations

Select

- Column selectivity (reducing number of columns)
- Alias columns and stream names



Select Settings Optimize Inspect Data Preview

Output stream name *

Incoming stream * [Join1](#)

Input columns * ☒ Auto Mapping ⓘ ☒ Re-map ☐ Delete 8 items: All inputs mapped

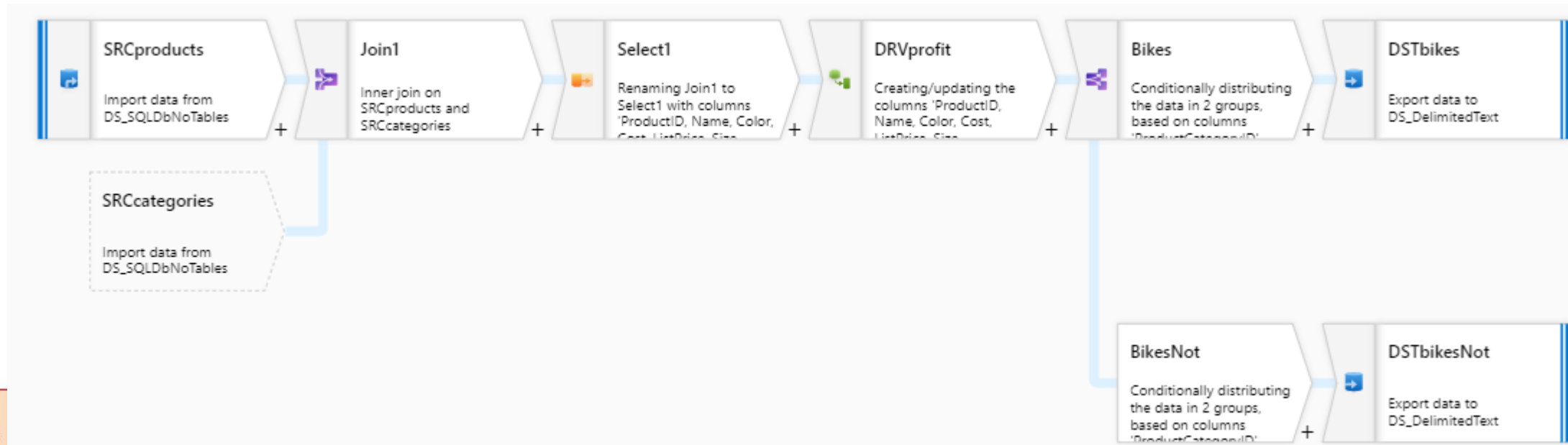
<input type="checkbox"/>	Join1's column		Name as
<input type="checkbox"/>	123 ProductID	→	<input type="text" value="ProductID"/>
<input type="checkbox"/>	abc SRCproducts@Name	→	<input type="text" value="Name"/>
<input type="checkbox"/>	abc Color	→	<input type="text" value="Color"/>
<input type="checkbox"/>	e ^x StandardCost	→	<input type="text" value="Cost"/>

DEMO

A. Load products from table, join with category names and split into bikes and other products.

B. Load products.csv, aggregate – grouping by category name, sort and save results into new table.

C. Sync data in table using Alter Row (delete incorrect email).



PLATINUM SPONSOR

STRATEGIC PARTNER

TECHNOLOGY
INNOVATION
DATA
KNOWLEDGE

GOLD SPONSORS



CLOUDS ON MARS



SILVER SPONSOR



BRONZE SPONSOR



DZIĘKUJĘ ZA UWAGĘ

tomasz.libera@datacommunity.pl
@tomasz_libera

bit.ly/sqlday2019_adf