# AZURE DATA FACTORY – DATA FLOWS

BRINGING THE "T" IN ELT

# WHO AM I?

- Practice Manager – Pragmatic Works – www.pragmaticworks.com
- Former CIO
- Blog, Speak, Record, Tweet
  - Twitter: @bizdataviz
  - LinkedIn: in/cseferlis
  - http://blog.bizdataviz.com
- MBA from UMass
- Ski, Bike, Hike, Run and drag the family

along…

# WHY ARE WE HERE?

- Modern Data Warehouse (Platform)

- Azure Data Orchestration

- Preview New Azure Data Factory Features

- Full Circle ELT with Azure Data Factory

# AZURE DATA FACTORY… QUICKLY

- Native Cloud Data Orchestration Tool

- Started with v1, but lacked many features

- v2 goes GA 6/2018 – work has begun on Data Flows

- Data Flows and other new features currently in Limited Preview

# LIMITED PREVIEW FEATURES

- Debug Mode
- Data Inspection
- Transformations
- GitHub Support
- Detailed Flow Inspection
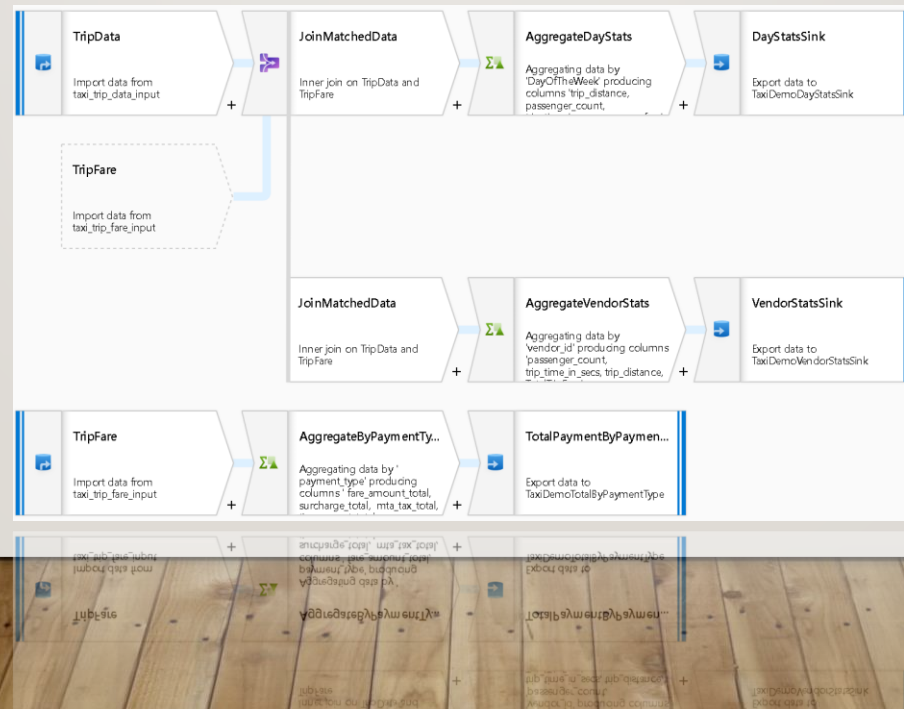- All-New Expression Builder

# VISUAL DATA FLOW AUTHORING

- Transform Data, At Scale, in the Cloud, Zero-Code
    - Cloud-first, scale-out ELT
    - Code-free dataflow pipelines

- Serverless scale-out transformation execution engine

- Maximum Productivity for Data Engineers
    - Does NOT require understanding of Spark / Scala / Python / Java

- Resilient Data Transformation Flows
    - Built for big data scenarios with unstructured data requirements
    - Operationalize with Data Factory scheduling, control flow and monitoring

# CODE-FREE DATA TRANSFORMATION AT SCALE

t: @bizdataviz
In/cseferlis

- Does not require understanding of Spark, Big Data Execution Engines, Clusters, Scala, Python …

- Focus on building business logic and data transformation

  - Data cleansing
  - Aggregation
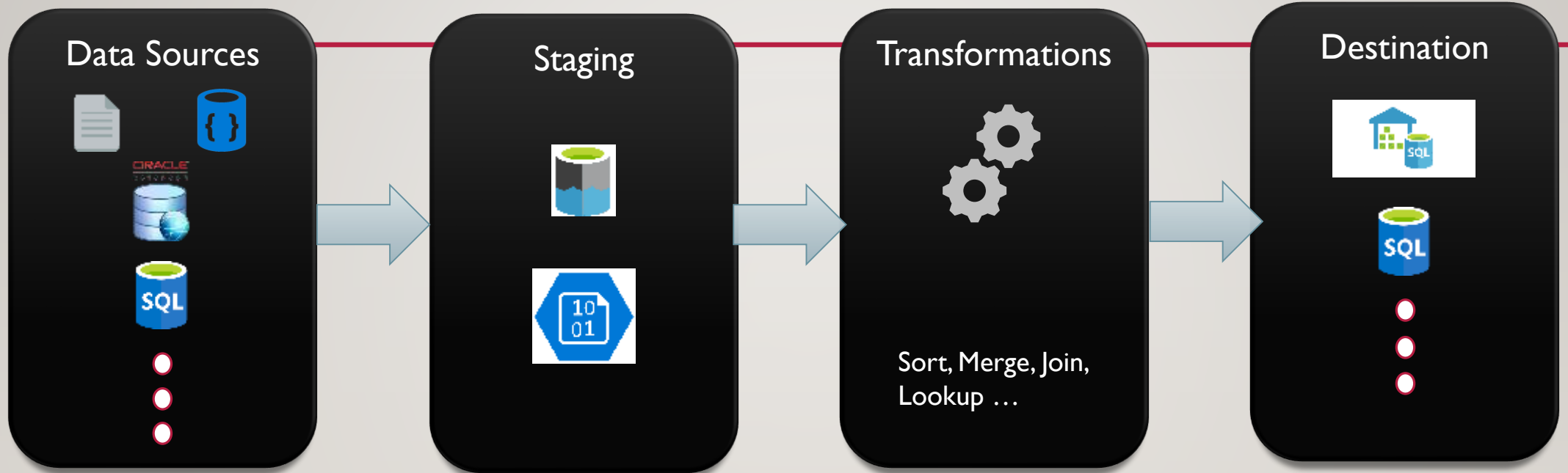  - Data conversions
  - Data prep
  - Data exploration



… not …

# ADF DATA FLOW WORKSTREAM

**Data Sources**

**Staging**

**Transformations**

Sort, Merge, Join, Lookup …

**Destination**

- Explicit user action
- User places data source(s) on design surface, from toolbox
- Select explicit sources

- Implicit/Explicit
- Data Lake staging area as default
- User does not need to configure this manually

- Explicit user action
- User places transformations on design surface, from toolbox
- User must set properties for transformation steps and step connectors

- Explicit user action
- User chooses destination connector(s)
- User sets connector property options

# DATA FLOW LIMITED PREVIEW

- Azure SLAs are NA for preview services (private or public preview) until GA of the service.

- Limited Preview Support
  - Handled directly with the Azure Engineering team via adfdataflowext@microsoft.com.

- Sign-up for ADF Data Flow service
  - http://aka.ms/dataflowpreview
  - Microsoft Azure must whitelist your subscription ID to turn on the feature for you

- Public Preview Support
  - Normal Azure customer service channels

t: @bizdataviz
In/cseferlis

# COMMON ELT SCENARIOS HANDLED WITH ADF – DATA FLOWS

# SIMPLE COPY FLOW

- ADF Data Flow is a guided construction process

- Begin by defining the Datasets for your Source and Sink

- Add Transformations to each node in your data flow

- Or simply copy from source to sink with no transformation

- Map columns and fields along the way

# SLOWLY CHANGING DIMENSION

- Common DW pattern to manage changing attributes to dimension members

- Graphically build code-free SCD ETL pattern to load your data warehouse

- Connect directly to Azure SQL DB and Azure SQL DW

- Use Lookup, Surrogate Key, Derived Column and Select transforms

# LOAD STAR SCHEMA DW

t: @bizdataviz
In/cseferlis



- Classic ETL pattern is easy to build in ADF's code-free Data Flow visual data transformation environment

- Add Aggregate transforms to produce calculations that you store in your analytical database schema

- Use Join transform to combine data from multiple data sources and data streams inside your data flow

- Land your data in your Lake folders or direct to Azure SQL DW

# DATA LAKE/DATA SCIENCE

t: @bizdataviz
In/cseferlis



- ADF supports building visual data transformations against your data directly in Data Lake locations (i.e. Azure Blob Store, Azure Data Lake Store)

- Built-in handling of schema drift for frequent changes in data lake file formats, columns, and data types

- Perform data exploration and data profiling across your data lake in ADF Data Flow with interactive debug data preview

t: @bizdataviz
In/cseferlis

# NEW FEATURES IN AZURE DATA FACTORY

# INTERACTIVE EXPRESSION BUILDER – BUILD TRANSFORM EXPRESSIONS, NOT CODE

t: @bizdataviz
ln/cseferlis

# BUILD RESILIENT DATA FLOWS WITH SCHEMA DRIFT HANDLING

# HANDLING SOURCE CHANGES

- Data Engineer Defines Source and you take ALL fields from source w/flexible schema

# HANDLING SOURCE CHANGES

- Data Engineer derives columns using template expression patterns based on name and type matching. No need to define static field names.
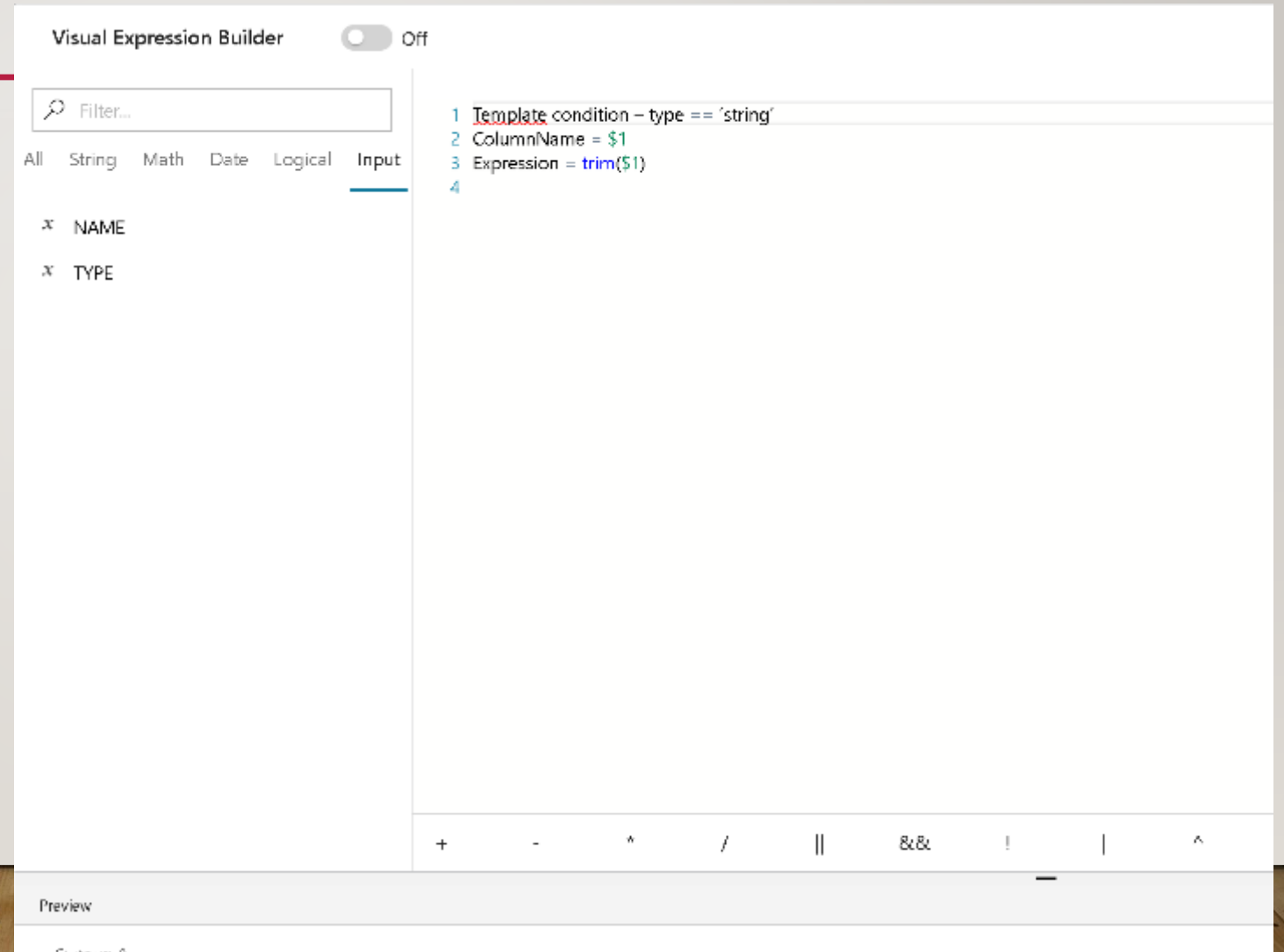
# HANDLING SOURCE CHANGES

- Data Engineer derives columns using template expression based on name and type matching. No need to define static field names.

# HANDLING SOURCE CHANGES

- Data Engineer derives columns using template expression based on name and type matching

# HANDLING SOURCE CHANGES

- Sink all incoming fields along with new derived field

# IMPORTANT LINKS:

t: @bizdataviz
In/cseferlis

- Sign-up for ADF Data Flow service:
  - [http://aka.ms/dataflowpreview](http://aka.ms/dataflowpreview)
  - Microsoft Azure must whitelist your subscription ID to turn on the feature for you
- GitHub Repository for documentation:
  - https://github.com/kromerm/adfdataflowdocs/

# THANK YOU!