**Microsoft**

# Azure Data Factory
## Visual Data Flow

*Limited Preview December 2018*

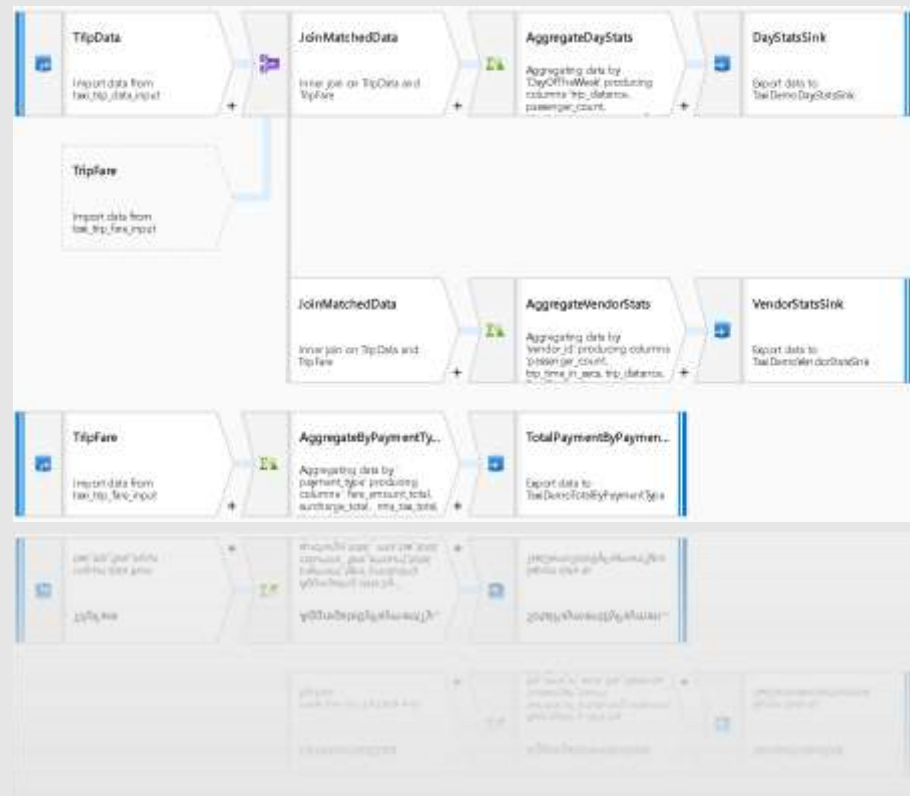# Visual Data Flow Authoring

- 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

# Visual Data Flow Key Tenets

- Visual "Data Flow Builder" / "Data Mapping"
- Extensible through scripting and expressions
- Data Flow can be embedded into ISV / SaaS apps
  - Embed UI
  - Embed Parameterize Data Flows
- A graphical UI for building data transformation routines on Spark
- Built for resiliency and operationalized environments

# Code-free Data Transformation At Scale

- 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

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

## Staging

- Implicit/Explicit
- Data Lake staging area as default
- User does not need to configure this manually
- Advanced feature to set staging area options
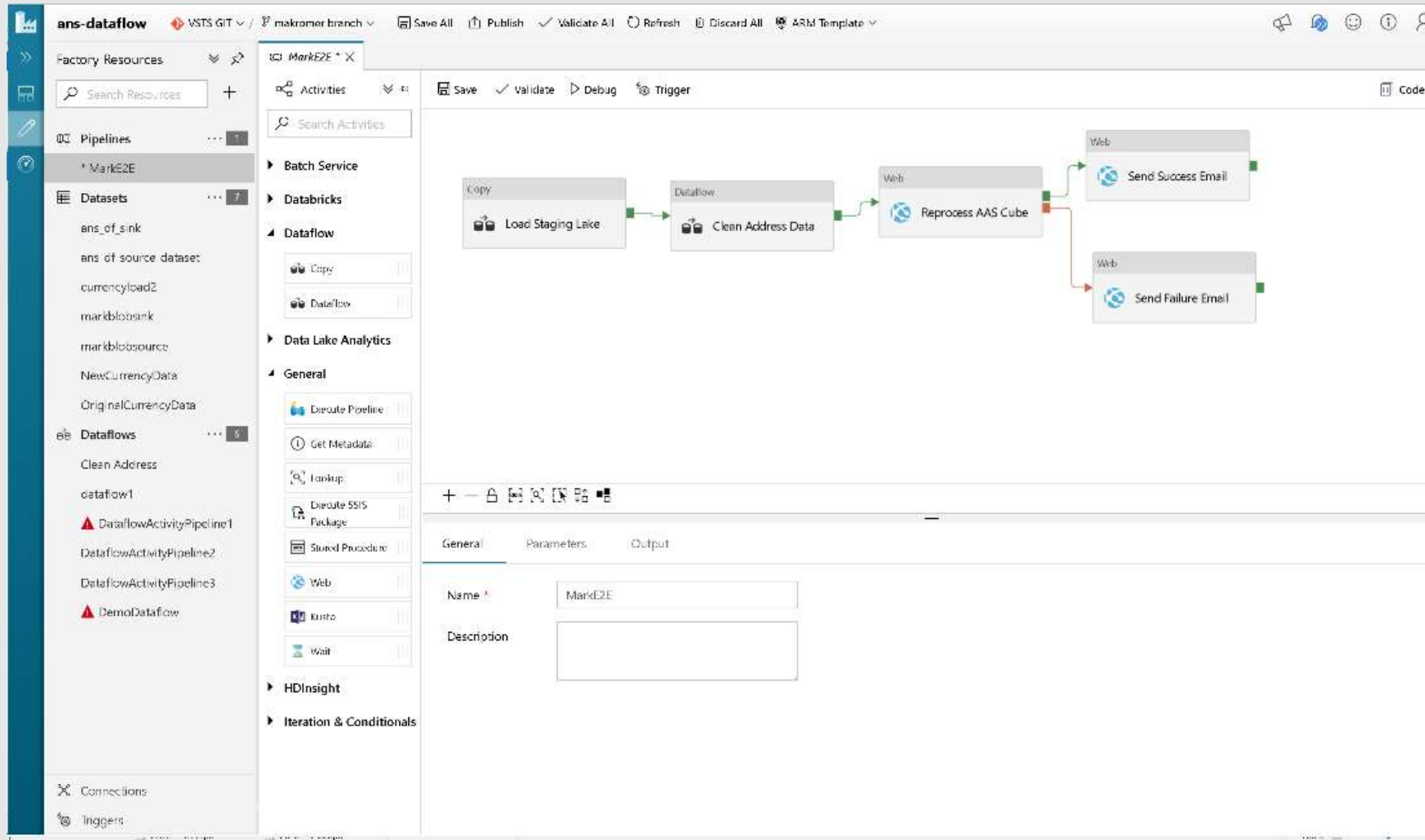- File Formats / Types (Parquet, JSON, txt, CSV ...)

## Transformations

Sort, Merge, Join, Lookup ...

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

## Destination

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

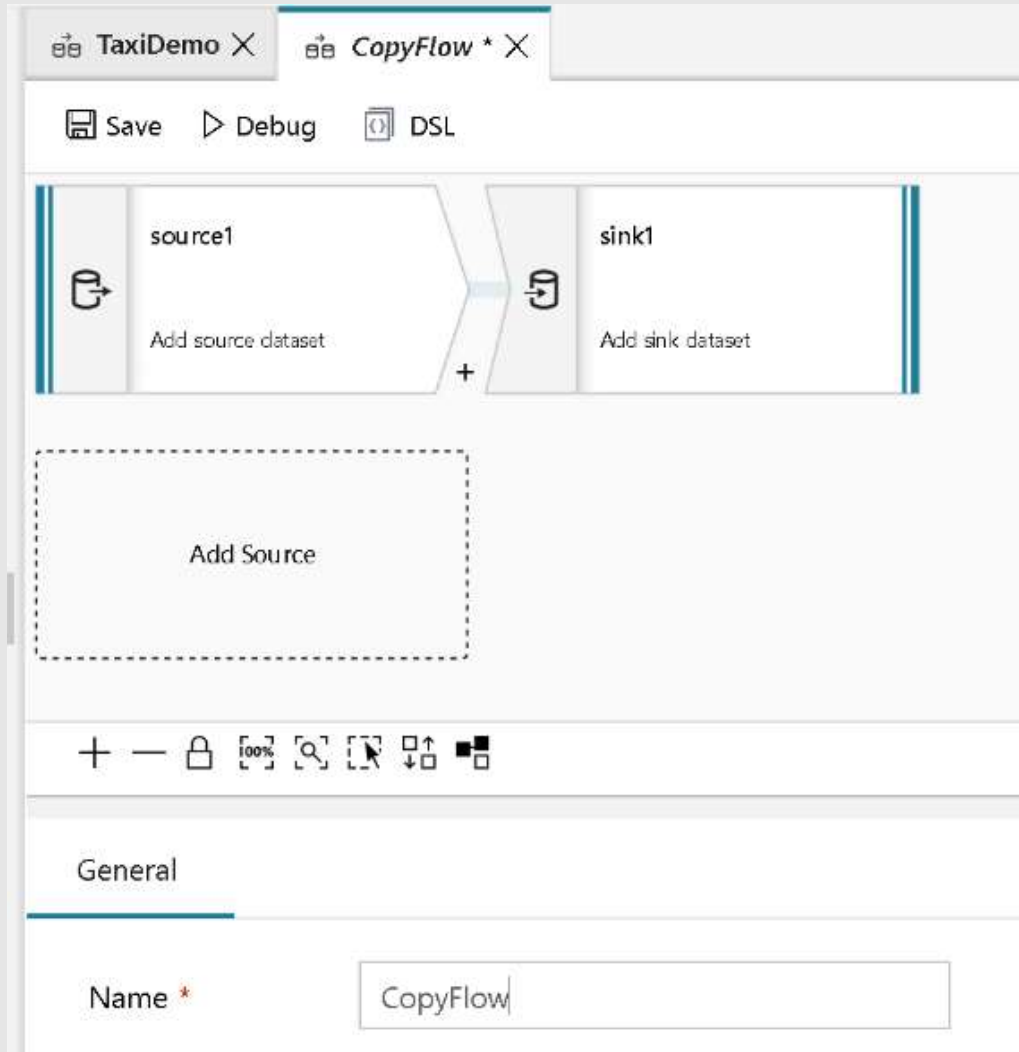# Data Flow Limited Preview Support & SLAs

- 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

# ADF Pipeline Execution of a Data Flow Activity



- Design code-free ETL workflows
- Copy data from on-prem, other clouds and Azure
- Stage data for transformation
- Build visual data transformations
- Schedule triggers for your pipeline execution
- Monitor processes and configure alerts
- All within ADF

# 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 Scenario



- 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 Scenario



- 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 Scenario



- 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 win interactive debug data preview

# Build your logical data flows adding data transformations in a guided experience

# Microsoft Azure Data Factory Continues to Extend Data Flow Library with a Rich Set of Transformations and Expression Functions

# Switch to Debug Mode and select sample data to work with for debugging

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

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

# Interactive Expression Builder – Build data transform expressions, not Spark code

# Azure Data Factory Visual Data Flow

# Deep Monitoring Introspection of Data Transformations

# Debug Data Flows with Data Preview and Data Sampling

# Build Resilient Data Flows with Schema Drift Handling

# Data Engineer Defines Source will take ALL fields from the source file with flexible schema

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

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

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

# Sink all incoming fields along with new derived field