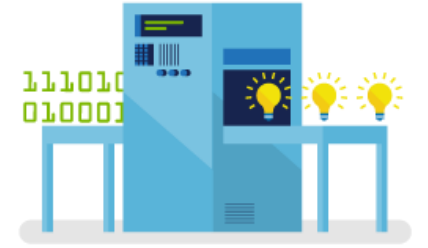# ELT using Azure Databricks and Data Factory
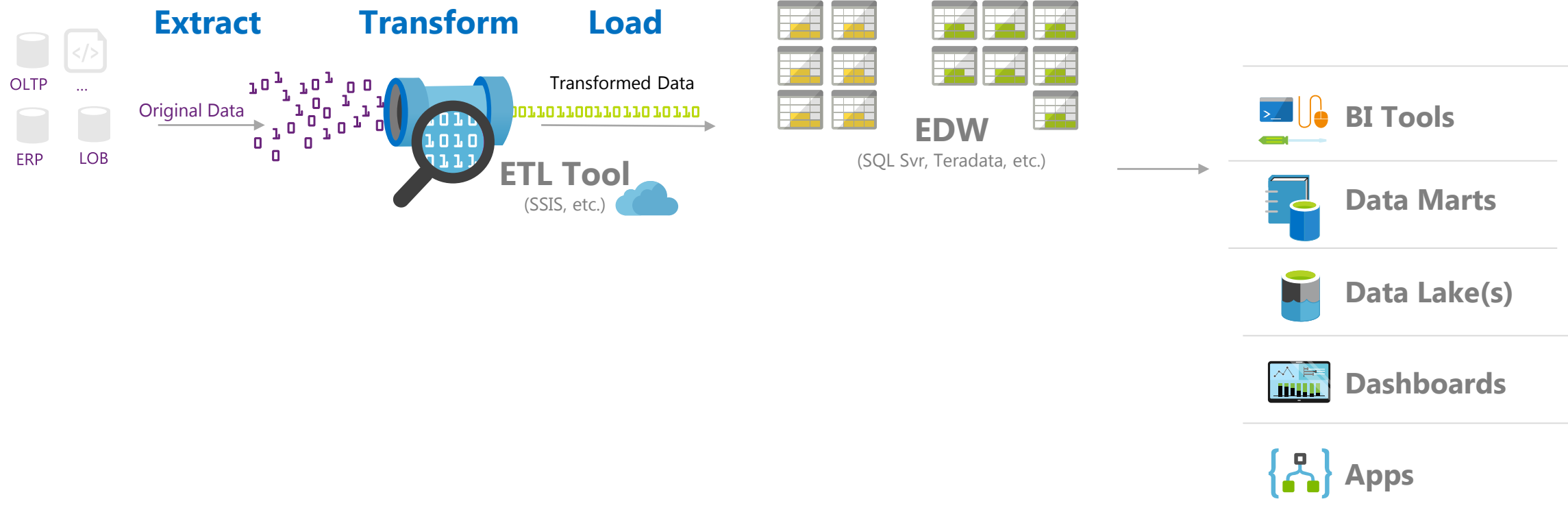
Gaurav Malhotra
Senior Program Manager-Microsoft
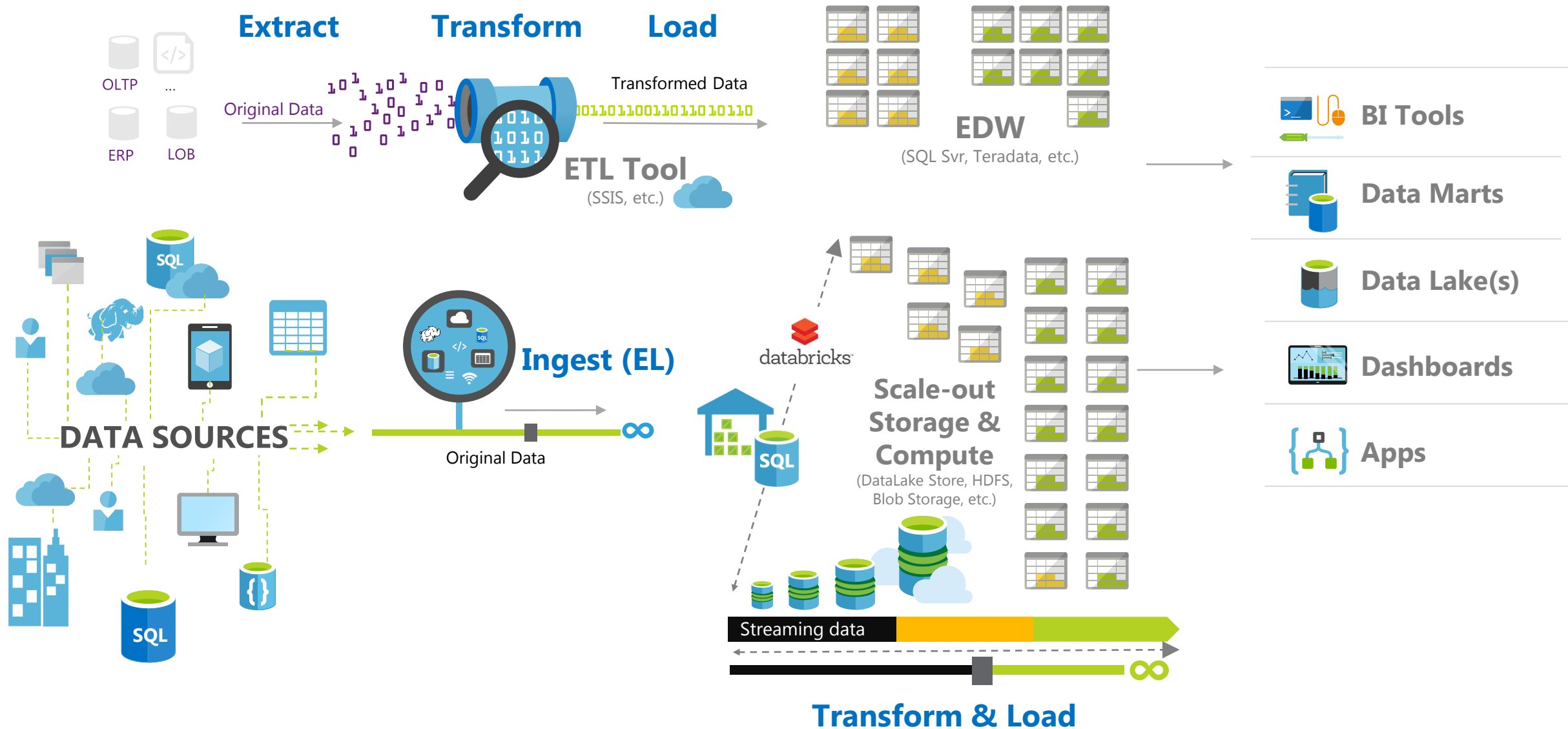
# Agenda

- Modern Data Engineering
- Azure Data Factory Overview
- Azure Databricks Overview
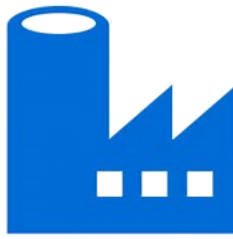- Demos
- Q & A

# Modern Data Engineering

**Extract**    **Transform**    **Load**

OLTP    ...

Original Data

ERP    LOB

Transformed Data

**ETL Tool**
(SSIS, etc.)

**EDW**
(SQL Svr, Teradata, etc.)

**BI Tools**

**Data Marts**

**Data Lake(s)**

**Dashboards**

**Apps**

**Extract**  **Transform**  **Load**

OLTP  ...

ERP  LOB

Original Data

Transformed Data

ETL Tool
(SSIS, etc.)

EDW
(SQL Svr, Teradata, etc.)

DATA SOURCES

SQL

Ingest (EL)

Original Data

databricks

Scale-out
Storage &
Compute
(DataLake Store, HDFS,
Blob Storage, etc.)

SQL

Streaming data

**Transform & Load**

BI Tools

Data Marts

Data Lake(s)

Dashboards

Apps

# Azure Data Factory
## Managed Data Integration Service

# Azure Data Factory
*Managed Data Integration Service*

## Flexible Pipeline Model

Rich pipeline orchestration
Triggers: on-demand, schedule, event

## Data Movement as a Service

Cloud, Hybrid
70+ connectors provided
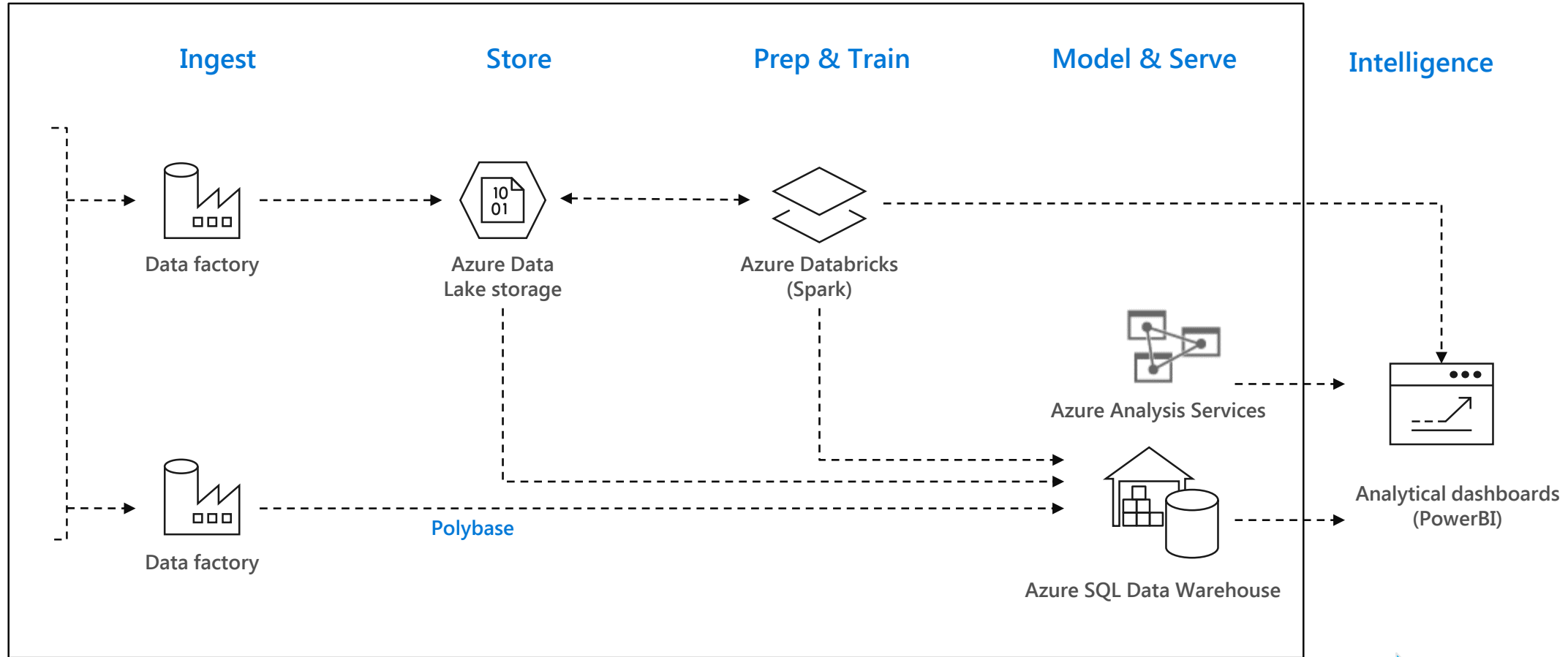
## SSIS Package Execution

In a managed cloud environment
Use familiar tools, SSMS & SSDT

## Author & Monitor

Programmability (Python, .NET, Powershell, etc.)
Visual Tools

# Modern Data Engineering for BI



Ingest | Store | Prep & Train | Model & Serve | Intelligence

Logs, files and media (unstructured)

On Prem, Cloud Apps & Data

Business / custom apps (Structured)

Data factory

Azure Data Lake storage

Azure Databricks (Spark)

Azure Analysis Services

Data factory

Polybase

Azure SQL Data Warehouse

Analytical dashboards (PowerBI)

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

# Azure Databricks
## Trusted & Reliable Platform for Data Engineering

# AZURE DATABRICKS

- Azure Databricks is a **first party** service on Azure.
  - Unlike with other clouds, it is not an Azure Marketplace or a 3$^{rd}$ party hosted service.

- Azure Databricks is integrated seamlessly with Azure services:
  - Azure Portal: Service an be launched directly from Azure Portal

  - Azure Storage Services: Directly access data in Azure Blob Storage and Azure Data Lake Store

  - Azure Active Directory: For user authentication, eliminating the need to maintain two separate sets of uses in Databricks and Azure.

  - Azure SQL DW and Azure Cosmos DB: Enables you to combine structured and unstructured data for analytics

  - Apache Kafka for HDInsight: Enables you to use Kafka as a streaming data source or sink

  - Azure Event Hub & Azure IOT Hubs: Enables you to use Event Hub and IOT Hub as a streaming data source

  - Azure Billing: You get a single bill from Azure

  - Azure Power BI: For rich data visualization

  - Azure Data Factory: ETL/ELT - See here

# Azure Databricks Role in Modern Data Warehouse



Logs, files and media (unstructured)

On Prem, Cloud Apps & Data

Business / custom apps (Structured)

**Ingest**

Data factory

Data factory

**Store**

Azure storage

Polybase

**Prep & Train**

Azure Databricks (Spark)

**Model & Serve**

Azure Analysis Services

Azure SQL Data Warehouse

**Intelligence**

Analytical dashboards (PowerBI)

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

# Infinite Scale, Lower Cost, Zero Management



1 to 1000s of Worker Nodes

Auto-scale Compute & Storage

Auto-Recovery & Upgrade

# Your Language, Your Data (Anywhere), Your Format

- SQL, Python, Scala & R Support
  - Code in your favorite language

- Source data from File System, Object stores, HDFS, Database, Pub-Sub systems & Others
  - Read and write data from/to multiple sources
  - Optimized for Azure Blob Store, ADLS, SQLDW, Event Hubs & Cosmos DB

- File Formats
  - CSV, JSON, Parquet, Text, ORC, XML & More

# Batch & Streaming Using Unified API

- Structured Streaming
  - Built on Spark SQL Engine
  - Express Streaming computation like batch computation on static data
  - Micro-batch & continuous processing support
  - Fault Tolerant, Only once computation
  - Supports
    - Late Data / Out of Order Data
    - Data de-duplication
    - Stream to Static Join
    - Stream to Stream Join

**Azure Databricks Unified Computing simplifies and accelerates Data Engineering**

# Demo-Connected Cars

# Connected car market

- The connected car market is growing

  - **45%** compound annual growth rate over 5 years
  - **10x faster** than overall car market
  - **75%** of cars shipped globally by **2020** will have necessary hardware to connect to the **Internet**

- Connected car technology is split between two approaches

  - Put the Internet connection in the car (**embedded connections**)
    - Does not require a phone data plan to operate
    - Provides access to more features and data
  - Rely on a **secondary device**

- **Embedded connections win**, because auto companies will be able to

  - Collect data on the performance of cars
  - Send updates and patches to cars remotely
  - Avoid recalls related to the car's software



**Connected-Car Shipments Forecast**
*(Global)*

■ Global Cars Shipped    ■ Shipped With Connectivity

Five-Year (2015-2020)
CAGR 45%

| | 2013E | 2014E | 2015E | 2016E | 2017E | 2018E | 2019E | 2020E |
|---|---|---|---|---|---|---|---|---|
| Global Cars Shipped | 69 | 72 | 75 | 78 | 81 | 84 | 88 | 92 |
| Shipped With Connectivity | 7 | 7 | 10 | 15 | 22 | 32 | 47 | 69 |

Millions

Source: Scotiabank, BI Intelligence Estimates

BI INTELLIGENCE

# Connected car market

**75%** of the cars shipped globally by **2020** will be built with the necessary hardware to connect to the Internet

Vehicle diagnostic

Usage-based insurance

Fleet management

Roadside assistance

Eco-driving

Engine performance remapping

Engine emission control

# Demo Architecture Diagram-Connected Cars

# Power BI dashboard

# Dataflow
## Visual Data Transformation in ADF (Private Preview)

# 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

# Azure Data Factory Visual Data Flow

# Simple Copy Flow

# Guided experience to build data flows

# 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

# Questions