# DATA WRANGLING FOR ETL ENTHUSIASTS

MOHAMED KABIRUDDIN

CLOUD SOLUTIONS ARCHITECT (DATA)

MICROSOFT

# ETL PROCESS

LANDING

STAGING

DIMENSIONAL MODEL

DID I JUST SPEND 10 HOURS PERFECTING THAT LOOKUP?

AND STILL APPLY INDEXES TO GAIN PERFORMANCE?

Data Flow | Event Handlers | Package Explorer | Execution Results

Excel Lookup

Generate 3 states

```
using System;
using System.Data;
using Microsoft.SqlServer.Dts.Pipeline.Wrapper;
using Microsoft.SqlServer.Dts.Runtime.Wrapper;

[Microsoft.SqlServer.Dts.Pipeline.SSISScriptComponentEntryPointAttribute]
public class ScriptMain : UserComponent
{
    /// <summary>
    /// Send 3 rows down the pipeline
    /// </summary>
    public override void CreateNewOutputRows()
    {
        StateCodeOutputBuffer.AddRow();
        StateCodeOutputBuffer.StateCode = "MO";

        StateCodeOutputBuffer.AddRow();
        StateCodeOutputBuffer.StateCode = "NE";

        StateCodeOutputBuffer.AddRow();
        StateCodeOutputBuffer.StateCode = "KS";
    }
}

-- EXCEL Query
SELECT
    S.StateCode
,   S.StateName
FROM
    `Sheet1$` S
```
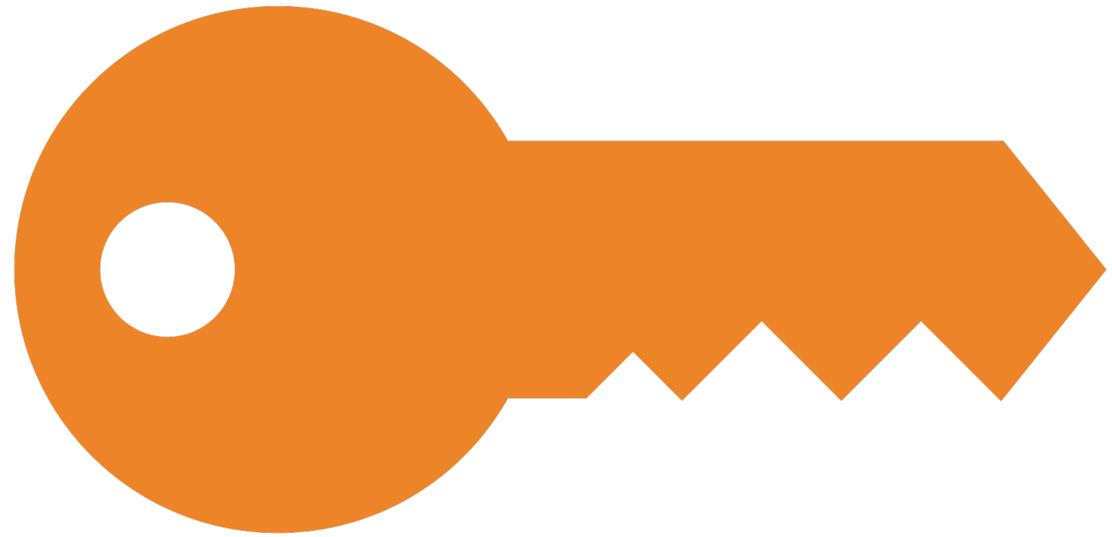
Lookup

Lookup Match Output

bit bucket

# CONSTRUCT THAT POWERFUL SURROGATE KEY

ULATIMATE GOAL:

KEEP THE
DATA WAREHOUSE
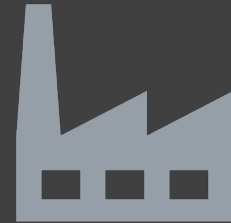
1. UPDATED
2. RELEVANT
3. OPERATIONAL

# DATA STAGES IN WRANGLING

RAW

REFINED

PRODUCTION

# DATA LAKE DESIGN CONSIDERATIONS

## Data Lake Zones

### Transient Landing Zone

Temporary storage of data to meet regulatory and quality control requirements. Limited access. May not be required depending on requirements.

### Raw Zone

Original source of data ready for consumption. Metadata publicly available but access to data still limited.

### Trusted Zone

Standardized and enriched datasets ready for consumption to those with appropriate role-based access. Metadata available to all.

### Curated/Refined Zone

Data transformed from Trusted Zone to meet specific business requirements.

### Sandbox Zone

Playground for Data Scientists for ad hoc exploratory use cases.

## Data Governance Considerations

### Security and Compliance

Access Control

Encryption
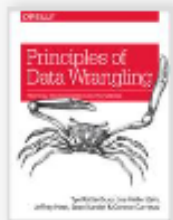
Row-Level Security

### Metadata Management

Data Quality

Metadata Management

Lifecycle Management

# ARE DATA WRANGLING AND ETL THE SAME THEN?

- Data wrangling is the process involved in transforming or preparing data for analysis

- Consider ETL to be one type of data wrangling, specifically a type of data wrangling managed and overseen by an organization's shared services or IT organization.

- Data wrangling can also be handled by business users in desktop tools like Excel, or by data scientists in coding languages like Python or R

## Principles of Data Wrangling
by Connor Carreras; Jeffrey Heer; Sean Kandel; Joseph M. Hellerstein; Tye Rattenbury
Published by O'Reilly Media, Inc., 2017

O'REILLY®

# TOOLSET FROM MICROSOFT

- Power BI
- Excel
- SSIS
- T-SQL
- U-SQL
- Polybase
- Azure Data Explorer
- Azure Data Factory
- Azure Stream Analytics
- Azure HD Insight (R & Python)
- Azure Databricks (R, Python and SparkSQL)

# MODERN DATA WAREHOUSING
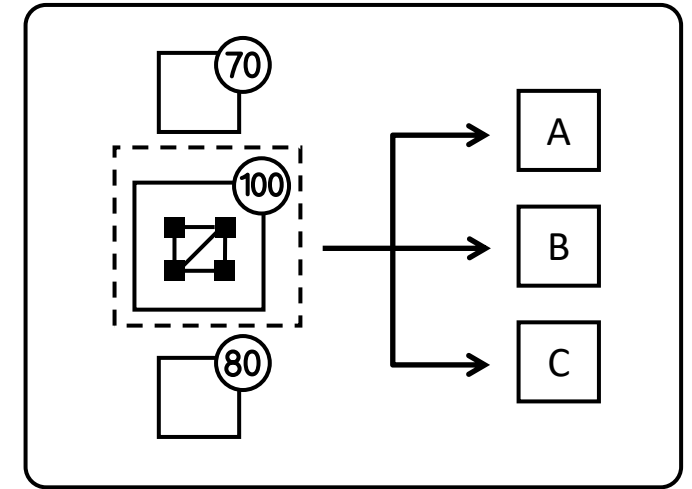
Canonical operations

## Load and ingest



Transfer and store
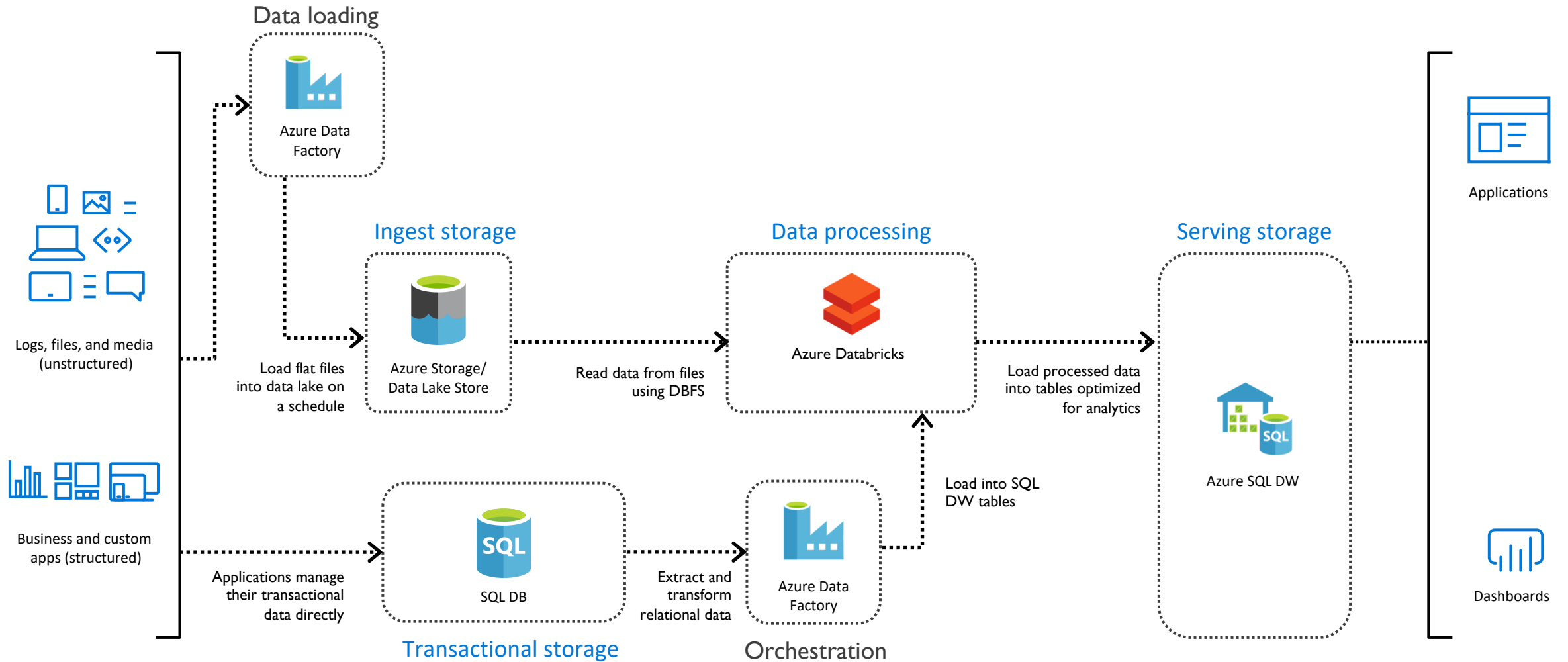
## Process



Process and clean

## Serve



Serve and analyze

# MODERN DATA WAREHOUSING PATTERN IN AZURE

Data processing with Azure Databricks

**Data loading**

Azure Data Factory

Logs, files, and media (unstructured)

Business and custom apps (structured)

**Ingest storage**

Azure Storage/ Data Lake Store

Load flat files into data lake on a schedule

**Data processing**

Azure Databricks

Read data from files using DBFS

**Serving storage**

Azure SQL DW

Load processed data into tables optimized for analytics

Applications

Dashboards

Applications manage their transactional data directly

SQL DB

**Transactional storage**

Extract and transform relational data

Azure Data Factory

**Orchestration**

Load into SQL DW tables

# New Data Flow

**Mapping Data Flow**
Code free data transformation at scale

**Wrangling Data Flow (Preview)**
Code free data preparation at scale

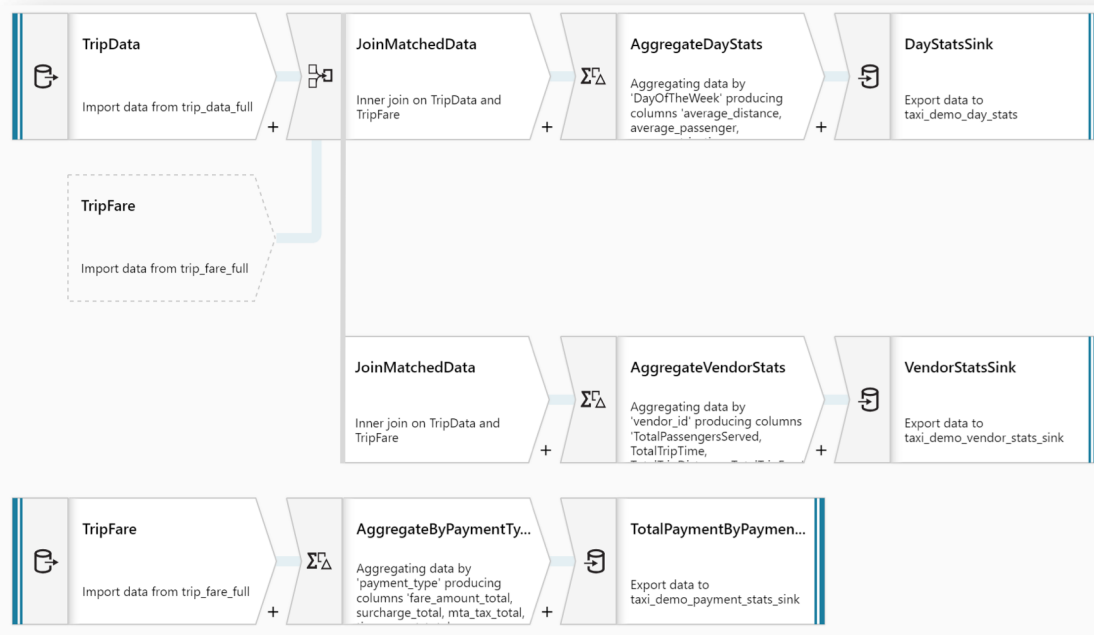# DATA FACTORY - DATAFLOWS

WRANGLING DATAFLOWS

# MAPPING DATA FLOW

No-code data transformation @ scale

Data cleansing, transformation, aggregation, conversion, etc.

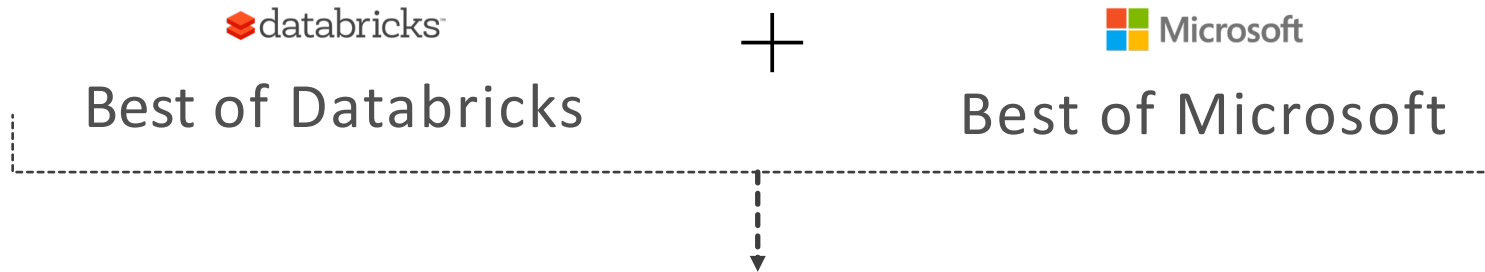Cloud scale via Spark execution

Easily build resilient data flows



... not

# AZURE DATABRICKS

A fast, easy and collaborative Apache® Spark™ based analytics platform optimized for Azure

**databricks**

**+**

**Microsoft**

## Best of Databricks

## Best of Microsoft

Designed in collaboration with the founders of Apache Spark

One-click set up; streamlined workflows

Interactive workspace that enables collaboration between data scientists, data engineers, and business analysts.

Native integration with Azure services (Power BI, SQL DW, Cosmos DB, ADLS, Azure Storage, Azure Data Factory, Azure AD, Event Hub, IoT Hub, HDInsight Kafka, SQL DB)

Enterprise-grade Azure security (Active Directory integration, compliance, enterprise-grade SLAs)

# AZURE DATABRICKS NOTEBOOKS

## Notebooks are a popular way to develop, and run, Spark Applications

Notebooks are not only for authoring Spark applications but can be *run/executed directly* on clusters

- Shift+Enter

- click the ▶ at the top right of the cell in a notebook

- Submit via Job

Fine grained permissions support so they can be *securely shared* with colleagues for collaboration

Notebooks are well-suited for prototyping, rapid development, exploration, discovery and iterative development



With Azure Databricks notebooks you have a default language but you can mix multiple languages in the same notebook:

%python   Allows you to execute python code in a notebook (even if that notebook is not python)

%sql      Allows you to execute sql code in a notebook (even if that notebook is not sql).

%r        Allows you to execute r code in a notebook (even if that notebook is not r).

%scala    Allows you to execute scala code in a notebook (even if that notebook is not scala).

%sh       Allows you to execute shell code in your notebook.

%fs       Allows you to use Databricks Utilities - dbutils filesystem commands.

%md       To include rendered markdown