February 6th-7th

# nic
## 20/20 VISION

Oslo Spektrum

# Understanding Azure Data Factory
## The What, When, and Why

**Cathrine Wilhelmsen**
NIC · February 6th, 2019

# Understanding Azure Data Factory

What is at the core of every Business Intelligence, Data Science, and Machine Learning project?

Data.

You need data to understand what has happened in the past, to predict what may happen in the future, to discover patterns and anomalies, and to gain the insight necessary for making faster and better decisions.

But before you can do any of those things, you need to collect, store, transform, integrate, and prepare your data. Azure Data Factory (ADF) is a service that enables you to quickly and efficiently create automated data pipelines – without having to write any code!

In this session, we will go through the fundamentals of Azure Data Factory and see how easy it is to build solutions that can work with all your data on-premises and in the cloud. We will explore some key features such as Mapping Data Flows for visual data transformations and Wrangling Data Flows for visual data preparation, as well as how to schedule and orchestrate your finished data pipelines. Throughout the session, we will discuss different use cases and scenarios, as well as when and why you should use Azure Data Factory for your projects.
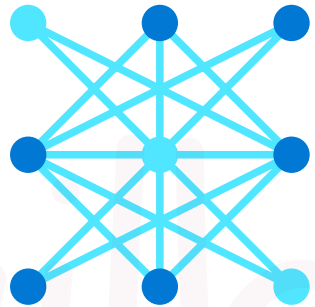
# cathrine
## WILHELMSEN

inmeta EVANGELIST MICROSOFT AZURE

MVP Microsoft Most Valuable Professional

@cathrinew

CW cathrinew.net
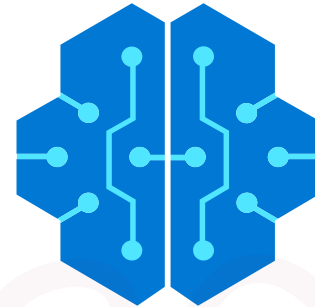
**Data Warehousing**

**Big Data and Analytics**

**Business Intelligence**

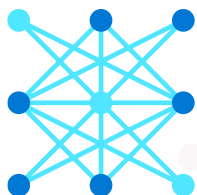**Artificial Intelligence**

**Data Science**

**Machine Learning**
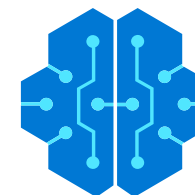
Big Data and Analytics

Data Warehousing

Business Intelligence

Artificial Intelligence

Machine Learning

Data Science

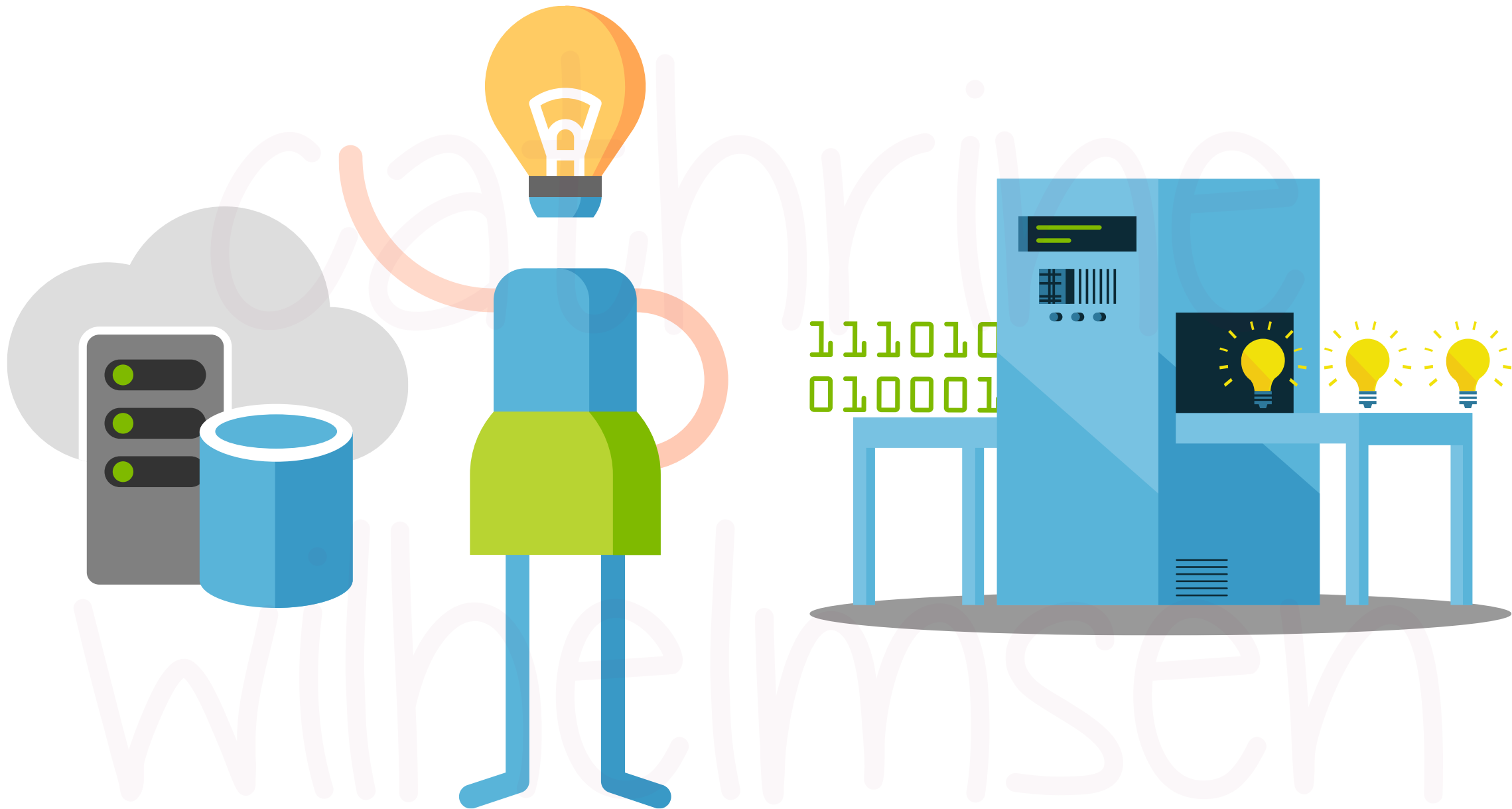What has happened?

What?

What will happen?

When did it happen?

When?

When will it happen?

Why?

Why did it happen?
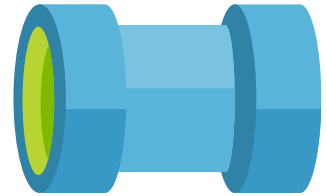
Collect
Store
Transform
Integrate
Prepare

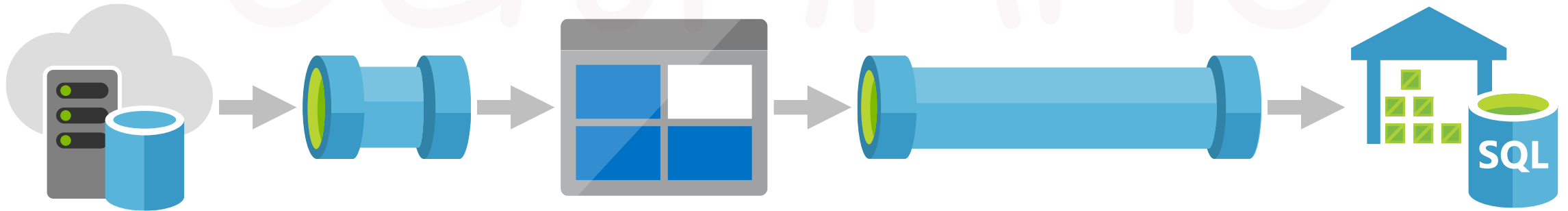# What is Azure Data Factory?

Hybrid data integration service
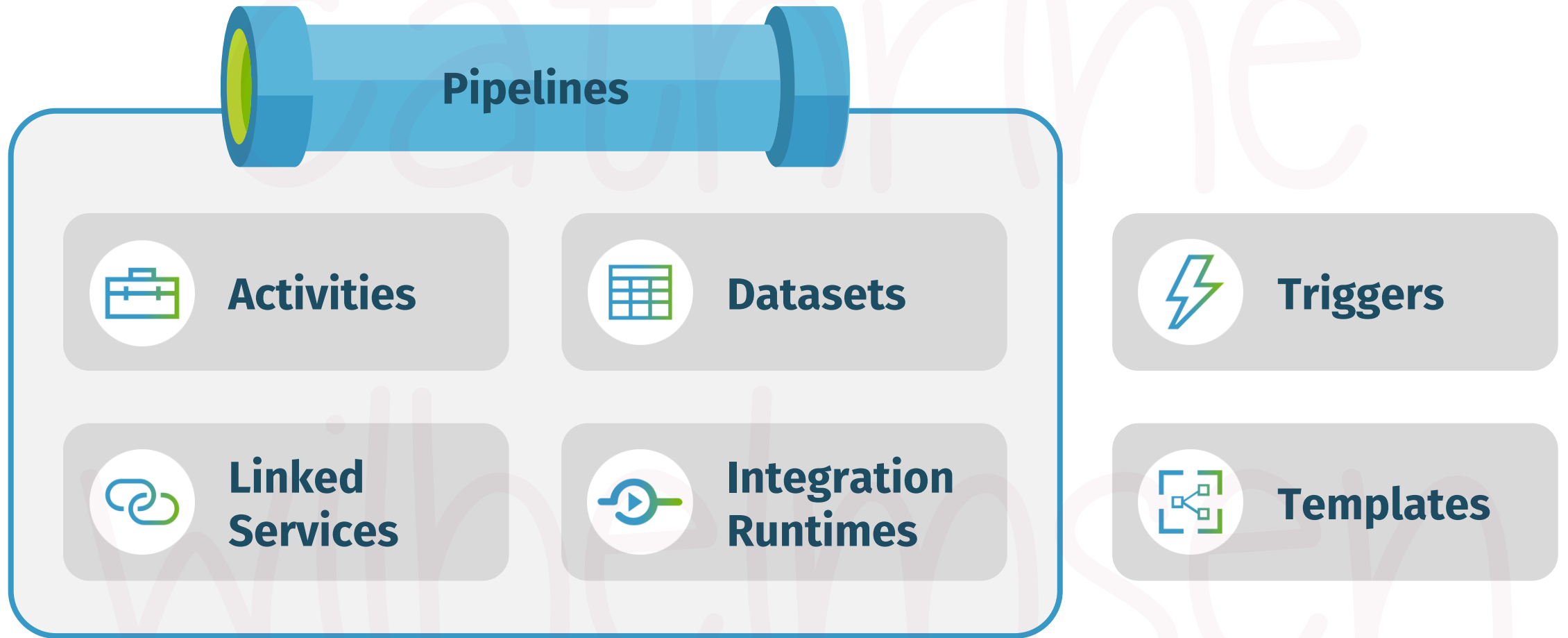
Complex and scalable pipelines

No-code ETL/ELT data flows

# What can you do in Azure Data Factory?

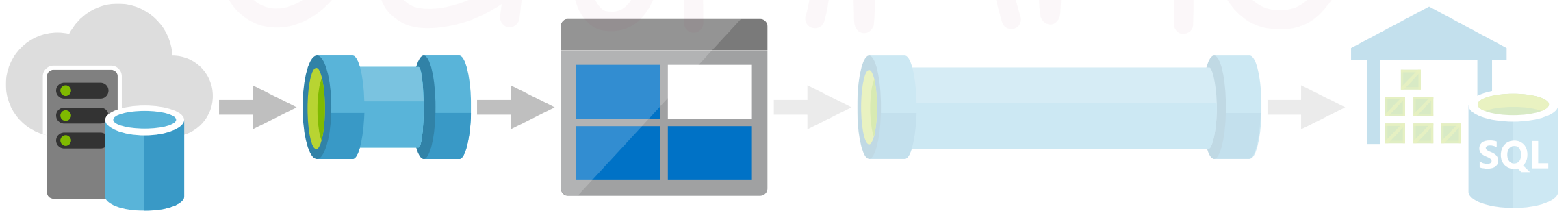**Copy Data**

**Transform Data**

# What is inside Azure Data Factory?

**Pipelines**

🧰 Activities

📊 Datasets

⚡ Triggers

🔗 Linked Services

▶️ Integration Runtimes

📐 Templates

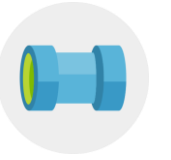DEMO

# Let's look inside
# Azure Data Factory!

# What can you do in Azure Data Factory?
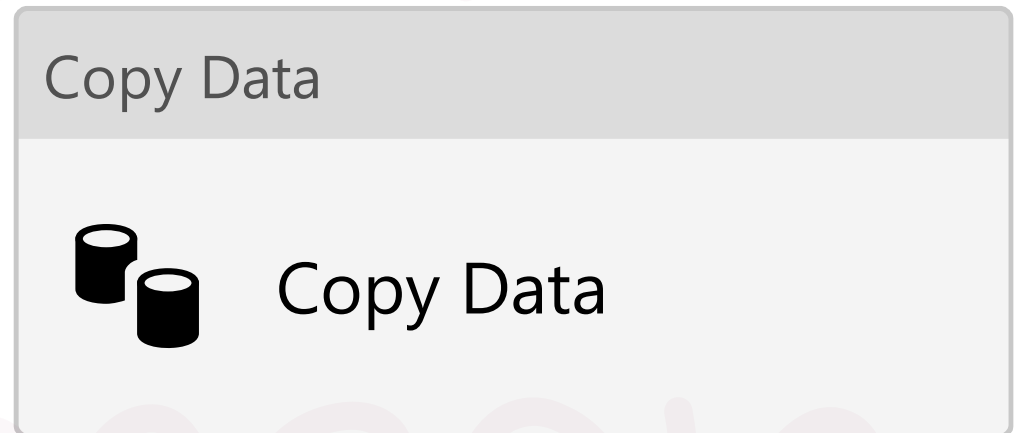


**Copy Data**

**Transform Data**

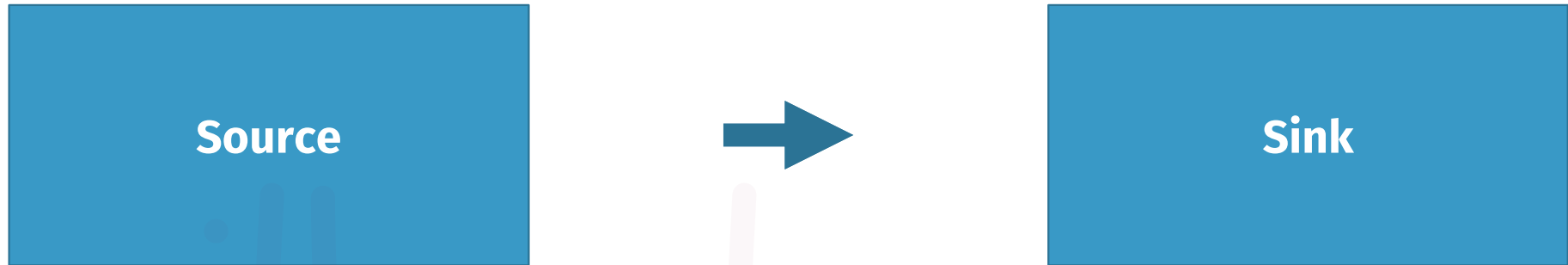# What is the Copy Data Activity?

The *core* activity *

Supports 80+ connectors

Copy from *Source* to *Sink*
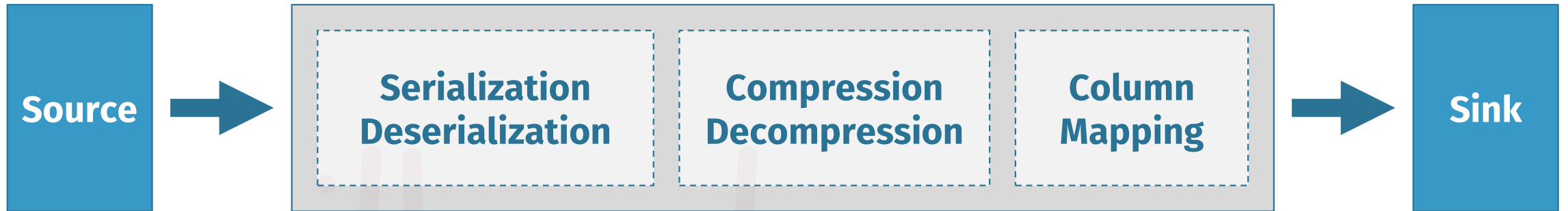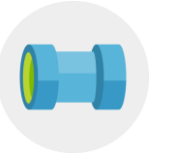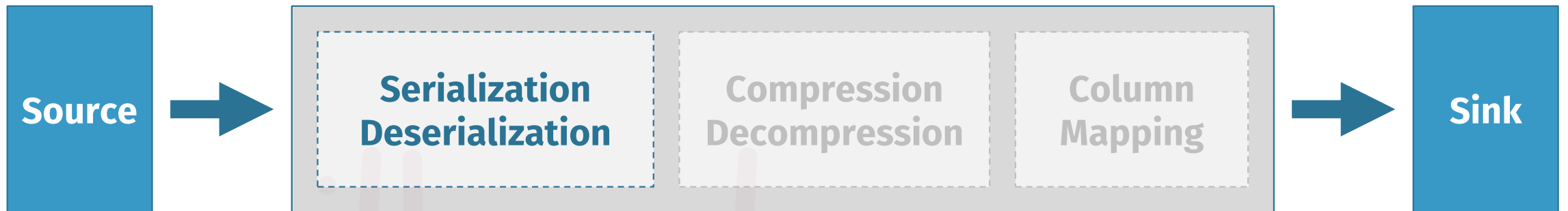
**Copy Data**

Copy Data

**\*** *Cathrine's opinion :)*

# Copy Data Process: Binary Files
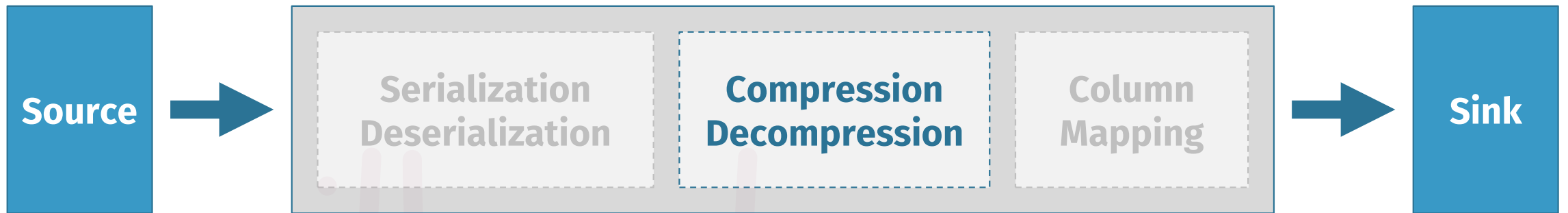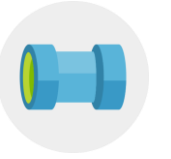
# Copy Data Process: Complex Files

| Source | → | Serialization Deserialization | Compression Decompression | Column Mapping | → | Sink |

# Copy Data Process: Complex Files

Source → Serialization Deserialization | Compression Decompression | Column Mapping → Sink

Convert file formats

# Copy Data Process: Complex Files

Source → | Serialization Deserialization | **Compression Decompression** | Column Mapping | → Sink

Zip or unzip files

# Copy Data Process: Complex Files

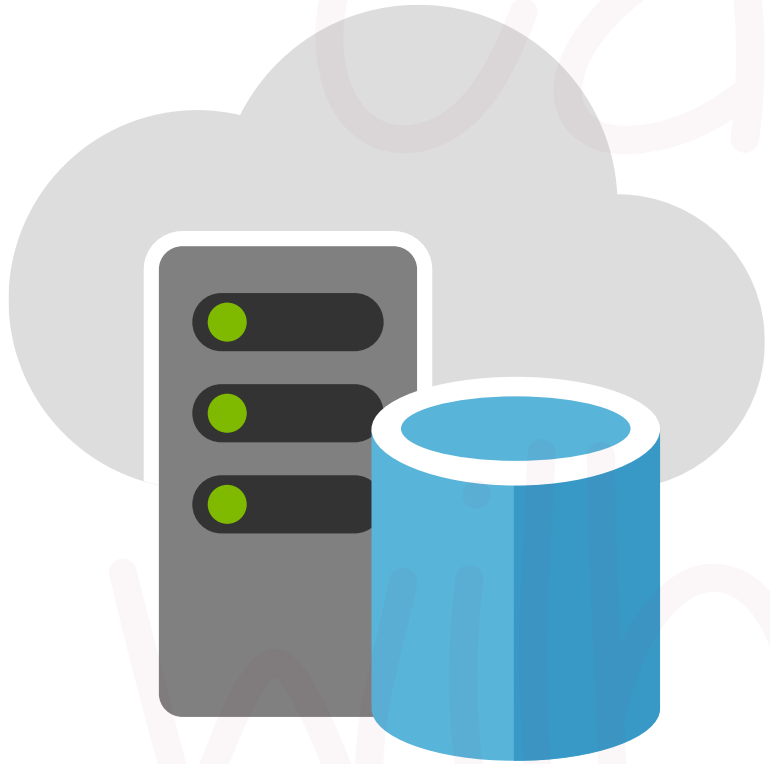| Source | → | Serialization Deserialization | Compression Decompression | Column Mapping | → | Sink |

Map columns implicitly or explicitly

DEMO

# Let's copy some data!

what if my systems are on-premises?

# Hybrid Azure Data Factory

Azure excels at cloud data integrations, but can also work with your on-premises systems!

# What are Integration Runtimes?

Azure Integration Runtime

Self-Hosted Integration Runtime

# Azure Integration Runtime

Restrict to specific Azure regions

- Data does not leave that specific region

Fully managed compute infrastructure

- Scale up by specifying Data Integration Units (DIUs)

# Self-Hosted Integration Runtime

Acts like a gateway

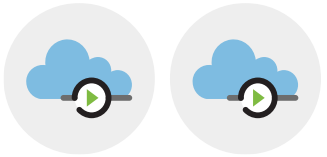- Get access to on-premises system within the network

Bring your own compute infrastructure

- Scale out by installing up to 4 nodes

# Copy Data Scenarios

Use Azure Data Factory to copy data between:
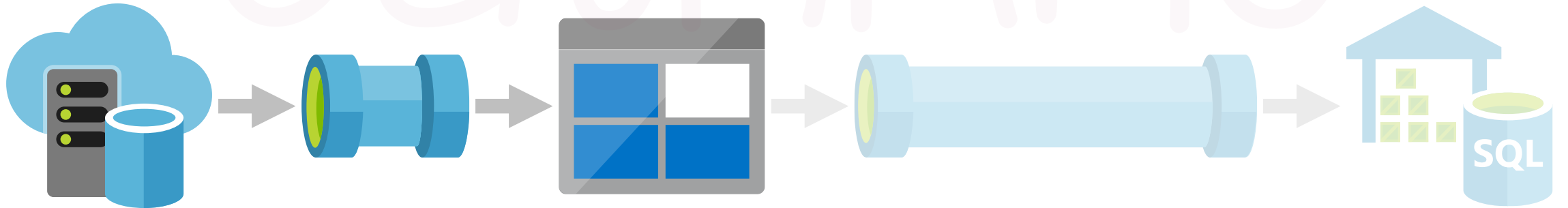
Cloud Stores

Cloud and On-premises Stores

On-premises Stores

DEMO

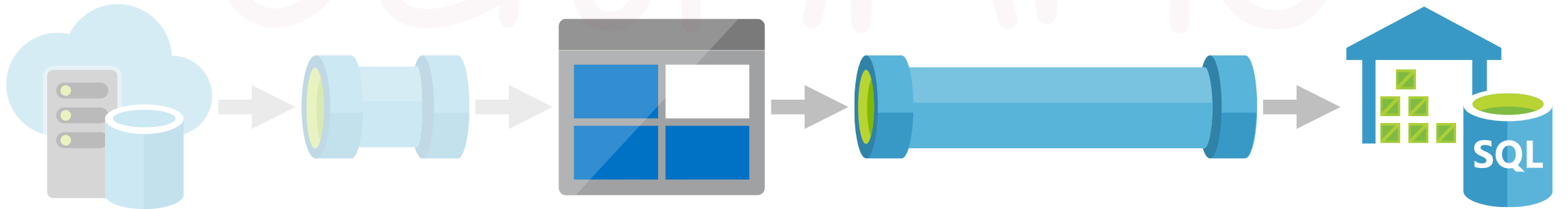# Let's connect to an *on-prem* SQL Server!

# Ok, so we can copy data...



**Copy Data**

**Transform Data**

# ...what about transforming data?

**Copy Data**

**Transform Data**

# Mapping or Wrangling
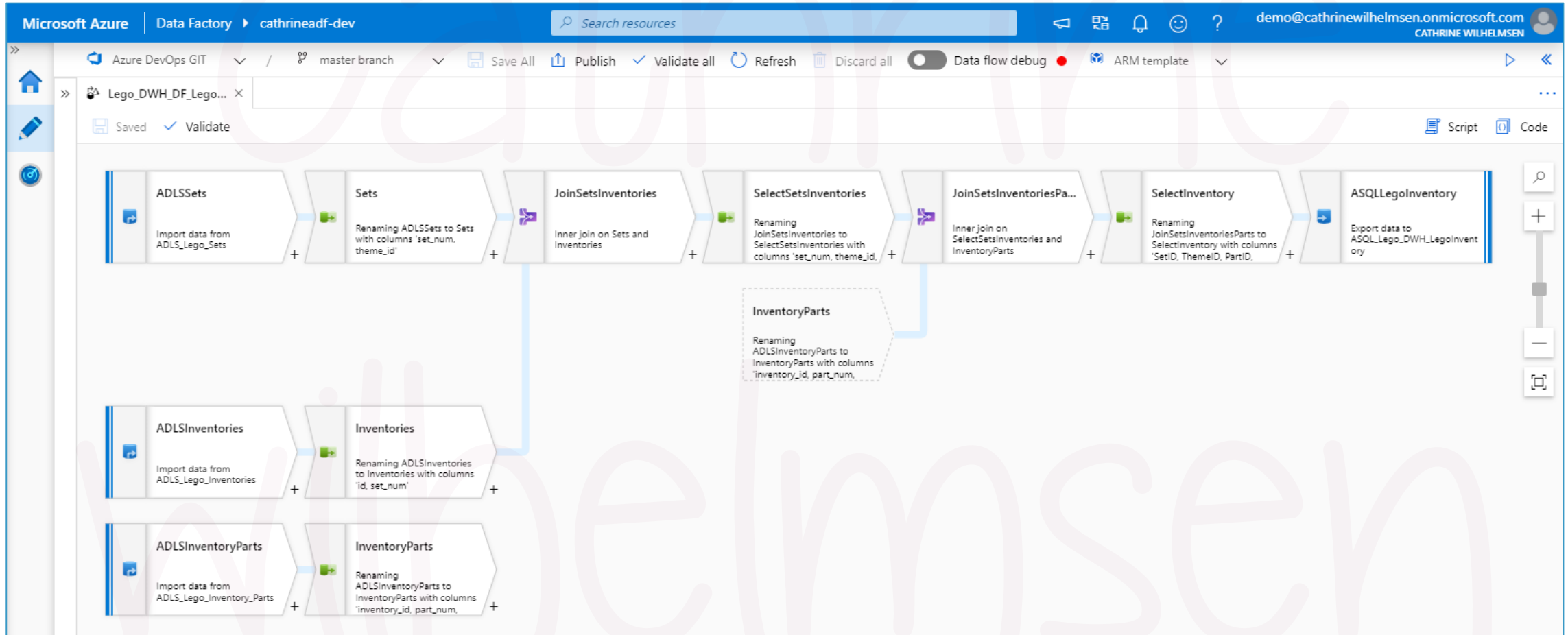
# What are Mapping Data Flows?

Data *transformation* at scale

Visual editor, no-code experience

Runs on Spark clusters

# How do Mapping Data Flows work?

# What are Wrangling Data Flows?

Data *preparation* at scale

Visual editor, no-code experience

Runs Power Query Online

# How do Wrangling Data Flows work?

**Mapping Data Flows**

Data Transformation

Similar to SSIS

**Wrangling Data Flows**

Data Preparation

Power Query Online

DEMO

# Let's transform
# some data!

how do we schedule data pipelines?

# Trigger pipelines...

On a set Schedule

In a Tumbling Window

When Event Happens

Now

# Triggers: Schedule

Execute one or more pipelines on a set schedule

- Every Wednesday at 06:00

- Last day of the month at 18:00

- Every Monday at 04:00 and Friday at 20:00

# Triggers: Tumbling Window

Execute a single pipeline for each time slice

- For every 15 minutes

- For every 1 hour

- For every 24 hours

# Triggers: Event Based

Execute one or more pipelines when event happens

- Blob is Created

- Blob is Deleted

- Blob is Created or Deleted

# Triggers: Now

Execute a single pipeline immediately

# Monitoring Triggers

Triggers save and log execution information

Information is available on the Monitor page

DEMO

# Let's schedule some pipelines!

# Azure Data Architectures

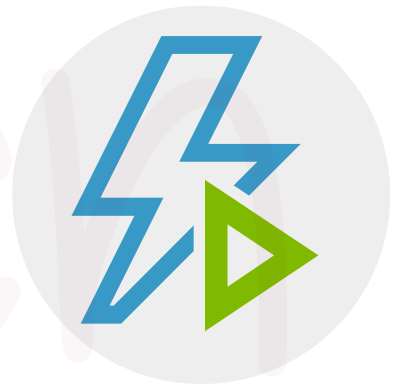# Advanced Analytics on Big Data



Logs, files, and media (unstructured)

Business/custom apps (structured)

Model and serve

Ingest — Azure Data Factory

Store — Azure Data Lake Storage

Prep and train — Azure Databricks (Python, Scala, Spark SQL, SparkR, Spark ML, SparklyR)

PolyBase

Azure Synapse Analytics

Azure Analysis Services

Power BI

Azure Cosmos DB

Web Application

1 2 3 4 5 6 7

# Real-time Analytics



Ingest      Store      Prep and train      Model and serve

Business/custom apps (structured)

Azure Data Factory — 2

Azure Data Lake Storage

PolyBase

Azure Synapse Analytics

Azure Analysis Services

Sensors and IoT (unstructured)

Clickstreams and Events (unstructured)

1

Azure HDInsight (Kafka)

3
4

Azure Databricks
(Python, Scala, Spark SQL, Spark R, Spark Structured Streaming)

5

6

Power BI

7

8

Azure Cosmos DB

Real-time Apps

https://azure.microsoft.com/en-us/solutions/architecture/real-time-analytics/

# Modern Data Warehouse



Logs, files, and
media (unstructured)

Business/custom
apps (structured)

Ingest — Store

Azure Data Factory ①

Azure Data
Lake Storage ②

Prep and train

Azure Databricks
(Python, Scala, Spark SQL,
SparkR, Spark ML, SparklyR)

PolyBase

Model and serve

Azure Synapse
Analytics ③

Azure Analysis
Services ④

Power BI ⑤

https://azure.microsoft.com/en-us/solutions/architecture/modern-data-warehouse/

© 2020 Cathrine Wilhelmsen (hi@cathrinew.net)

**Sources**

On-Premises

Cloud

SaaS

**Azure Synapse Analytics**

**Visualize**

Power BI

Good luck!

# *thank you!*

hi@cathrinew.net

@cathrinew

cathrinew.net

inmeta
EVANGELIST

MICROSOFT
AZURE