12/19/2019

# MOVING DATA AROUND IN AZURE

# moving data around in Azure

In the Microsoft Data platform ecosystem in Azure, there are several possibilities to move data around. The choice can be overwhelming and it's not always clear which option should be used.

In this session, we aim to give you an overview of a couple of products in Azure for data pipelines: Azure Data Factory, Logic Apps, Integration Services et cetera. At the end, you'll have a better understanding of the various alternatives in Azure.

# Koen Verbeeck

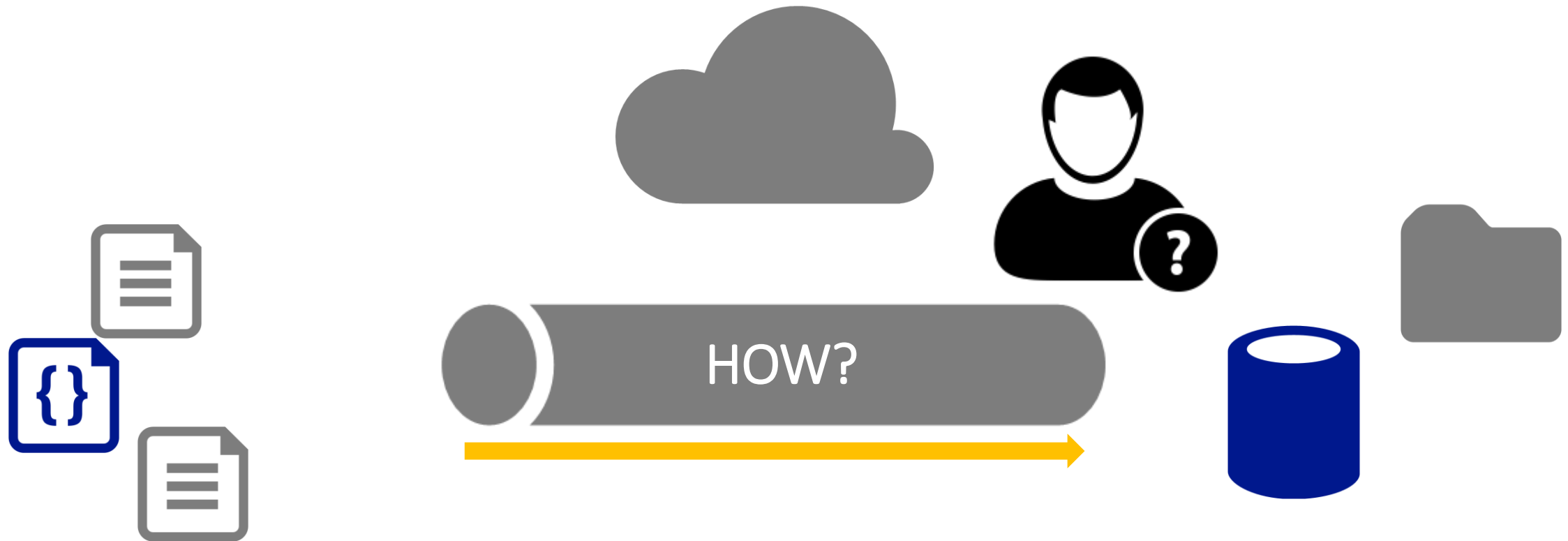## contact information

koen dot verbeeck at ae dot be

@Ko_Ver

LinkedIn

https://sqlkover.com

# the goal of this session



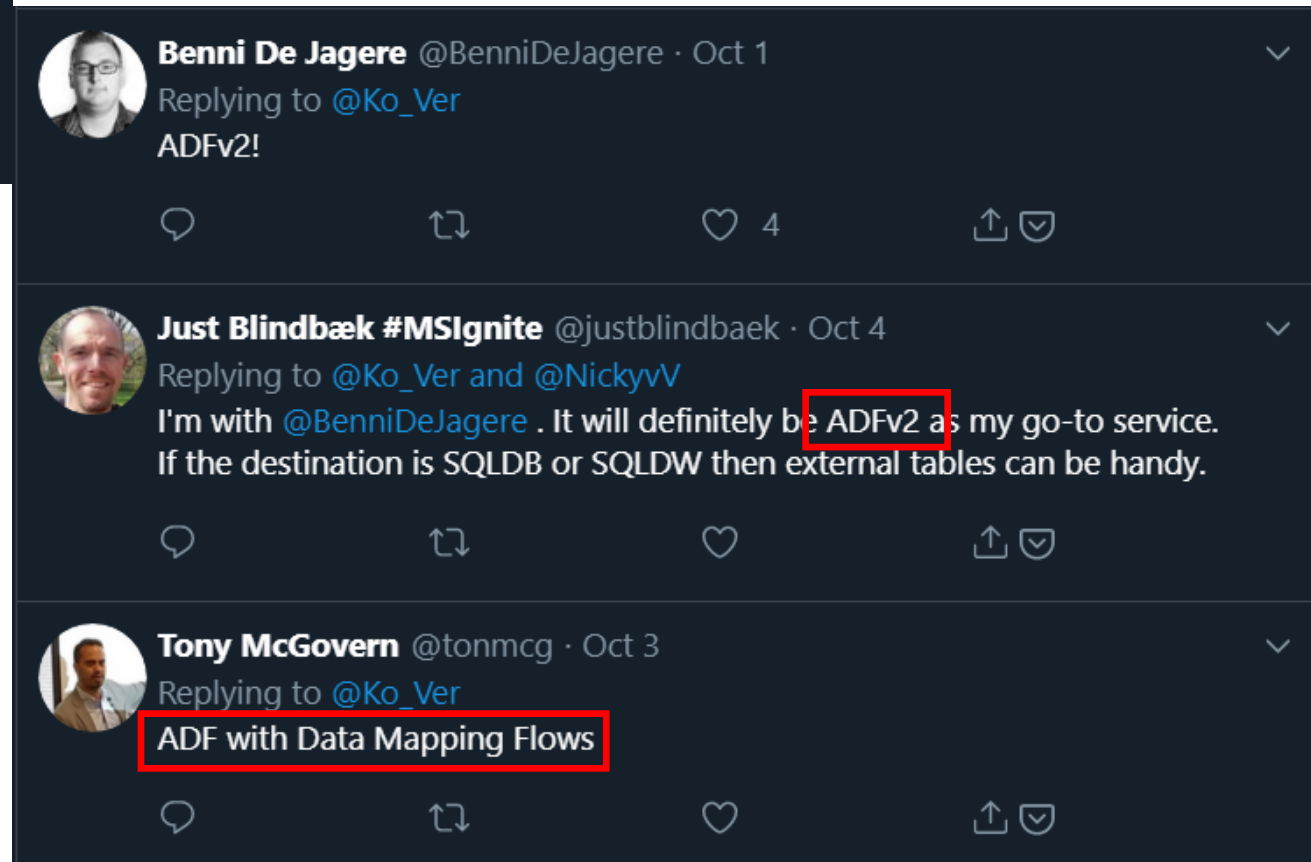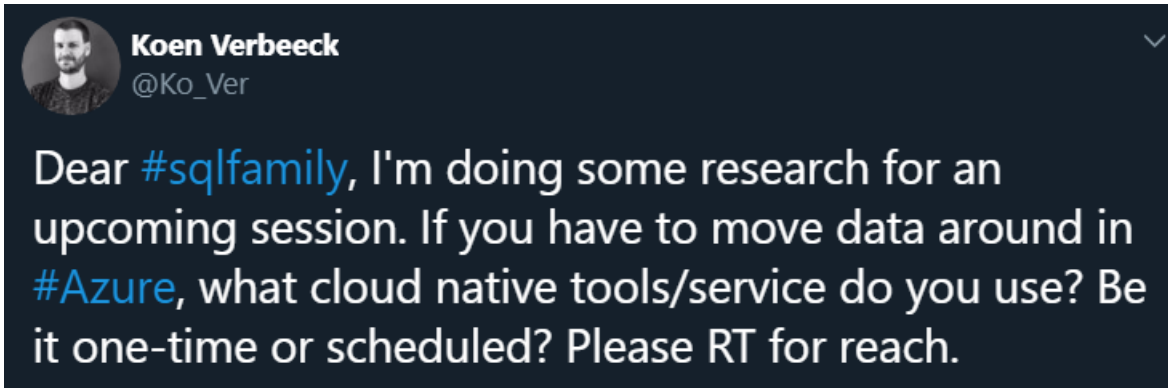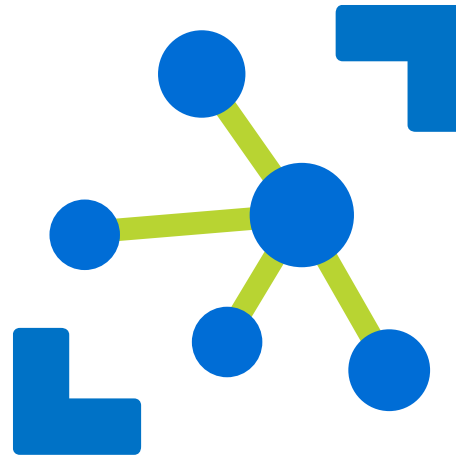HOW?

# research

what we'll
<span style="color:red">not</span>
be talking about

# specific use cases

# anything that runs on a VM

# … or some sort of backup/migration service

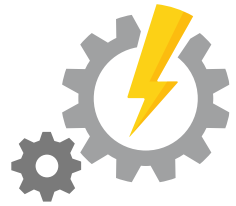also any service announced when I was making this slide deck



Azure Synapse Analytics
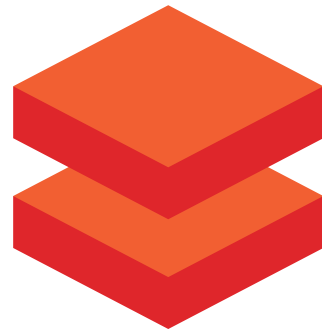(formerly SQL DW)

# what we will be talking about

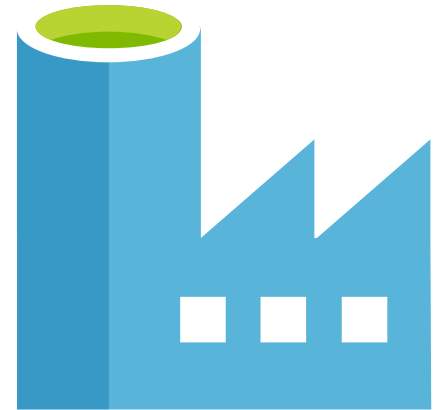# some Azure "data" services

## writing code

automation

functions

databricks

## not writing code

logic apps

data factory

azure automation

# azure automation

aka runbooks

can run Python & PowerShell

no Visual Studio integration

better suited for automating administrative tasks

# azure functions

# azure functions

serverless compute

event-driven

Visual Studio integration

C#, PowerShell, Python, Batch, Bash, JavaScript, PHP or F#

# azure functions

# DEMO TIME

azure databricks

# azure databricks

# azure databricks

# azure databricks

# azure databricks

# azure databricks

**DataFrame**

| | | |
|---|---|---|
| | | |
| | | |
| | | |
| | | |
| | | |

partition

partition

partition

**RAW**

CSV

## DataFrame
- **Schema** ← Parameter
- **Format** ← Parameter
- **Location** ← Parameter

```
df = (spark
        .read
        .schema(newSchema)
        .format(fileFormat)
        .load(dataLocation)
    )
```

pythondatacopy (Python)

⏱  ?  👤

⊹ ● my-spark-cluster | ⌄   📄 File ▾   🖼 View: Code ▾   🔒 Permissions   ⊙ Run All   ✎ Clear ▾          ⌨   📅 Schedule   💬 Comments   ▶ Runs   ↺ Revision history

| path | name | size |
|---|---|---|
| dbfs:/mnt/koenblobstorage/Top250Movies.csv | Top250Movies.csv | 5552 |

Command took 8.31 seconds -- by koen.verbeeck@outlook.com at 10/31/2019, 10:38:54 AM on my-spark-cluster

Cmd 4

```
1  df = spark.read.format("csv").option("header", "true").option("encoding", "cp1252").load("dbfs:/mnt/koenblobstorage/Top250Movies.csv")
```

▸ (1) Spark Jobs

▸ ▦ df: pyspark.sql.dataframe.DataFrame = [Movies: string]

Command took 4.58 seconds -- by koen.verbeeck@outlook.com at 10/31/2019, 10:39:22 AM on my-spark-cluster

Cmd 5

```
1  df
```

Out[2]: DataFrame[Movies: string]

Command took 0.05 seconds -- by koen.verbeeck@outlook.com at 10/31/2019, 10:39:25 AM on my-spark-cluster

Cmd 6

```
1  df.take(10)
```

▸ (1) Spark Jobs

Out[3]: [Row(Movies='1. The Shawshank Redemption'),
 Row(Movies='2. The Godfather'),
 Row(Movies='3. The Godfather: Part II'),

# DEMO TIME

azure logic apps

# azure logic apps

"business workflows without coding"

lots of built-in connectors and actions

Azure version of ~~Microsoft Flow~~

Power Automate

# azure logic apps

extensibility through customer connectors

e.g. an API and its swagger file

each API method can be translated into an action

# DEMO TIME

azure data factory

# data factory v2

data integration service

ELT/ETL without writing code

scalable data pipelines

git integration

## Repository settings

Enter Git repository information to be associated with your Data Factory:

Repository type *

Select... ▼

- Azure DevOps Git
- GitHub

# linked services

# datasets

What type of data? Where is it located?
What is the schema?



dbo.dimEmployee



/clean/2019/12/10/logs.csv

# pipelines

# mapping data flow

# wrangling data flow

# runtime



Azure Integration Runtime

Self-hosted Integration Runtime

Azure SSIS Integration Runtime

architects
for business
& ict

# Integration Runtime Setup ✕

Name *ⓘ

mssqltips-ir

Description ⓘ

SSIS IR for mssqltips

Type

Azure-SSIS

Location *ⓘ

West Europe ▾

Node Size *ⓘ

Standard_D1_v2 (1 Core(s), 3584 MB) ▾

Node Number *ⓘ

2

Edition/License *ⓘ

Standard ▾

Save Money

Save with a license you already own. Already have a SQL Server license?

Yes **No**

CONCLUSION

architects
for business
& ict

# conclusion

Azure Functions

      if you like coding

      event-based

      more suitable for smaller, specific tasks

Azure Databricks

      big data solution

      for data engineering / data science

ae

architects
for business
& ict

# conclusion

Logic Apps

      easy to build workflows

      lots of connectors

      event-based

Azure Data Factory

      most comprehensive and flexible

      visual tool

      big & small data