# Dave Ruijter

**Azure Solution Architect**

Macaw Netherlands

@DaveRuijter

ModernData.ai

## Battle Of Modern Data Architectures

Have you had that moment, where you are in doubt which Azure service to use for your 'Modern Data Warehousing' solution? So many good options..
Like the Mapping/Wrangling Data Flows capabilities in Azure Data Factory, or the Delta feature in Databricks!

In this session we will look at the different services, compare them using real use-cases, and learn how to choose the best fit for each scenario.
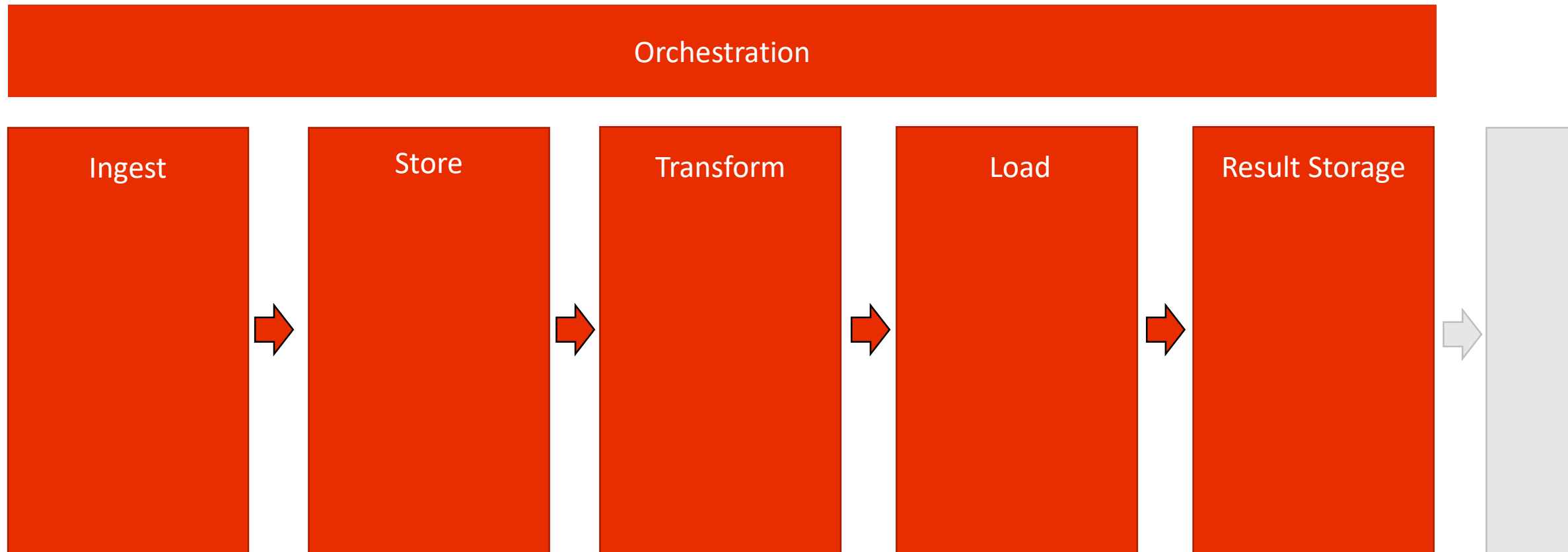
# Our Partners

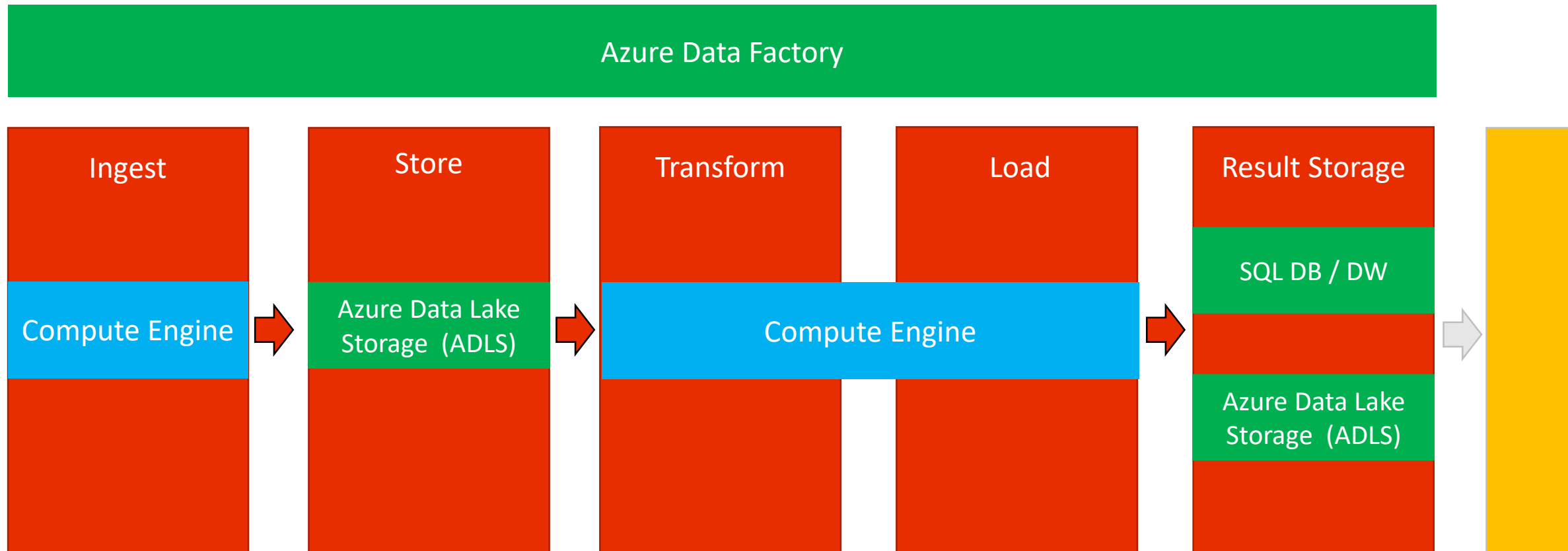# Datawarehousing in the cloud…

# Datawarehousing in the cloud…

# Datawarehousing in the cloud…

# Why is it so complicated?

- Various sources/formats
- Schema drift
- Scalability / "big data"
- Monitoring & Auditing
- Keep up with the business needs
- Version control
- ALM / DevOps

# Datawarehousing in the cloud…

# Datawarehousing in the cloud…

- Azure SQL /w PolyBase
- HDInsight
- Azure Data Lake Analytics
- Azure SSIS Integration Runtime
- Azure Data Factory (Mapping Data Flows)
- Azure Databricks
- Power BI dataflows

# Datawarehousing in the cloud…

- <mark>Azure SQL /w PolyBase</mark>
- HDInsight
- Azure Data Lake Analytics
- Azure SSIS Integration Runtime
- Azure Data Factory (Mapping Data Flows)
- Azure Databricks
- Power BI dataflows

dataMinds Connect 2019

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~
- <mark>HDInsight</mark>
- Azure Data Lake Analytics
- Azure SSIS Integration Runtime
- Azure Data Factory (Mapping Data Flows)
- Azure Databricks
- Power BI dataflows

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~
- ~~HDInsight~~
- <mark>Azure Data Lake Analytics</mark>
- Azure SSIS Integration Runtime
- Azure Data Factory (Mapping Data Flows)
- Azure Databricks
- Power BI dataflows

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~
- ~~HDInsight~~
- ~~Azure Data Lake Analytics~~
- <mark>Azure SSIS Integration Runtime</mark>
- Azure Data Factory (Mapping Data Flows)
- Azure Databricks
- Power BI dataflows

dataMinds Connect 2019

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~
- ~~HDInsight~~
- ~~Azure Data Lake Analytics~~
- <mark>Azure SSIS Integration Runtime</mark>
- <mark>Azure Data Factory (Mapping Data Flows)</mark>
- Azure Databricks
- Power BI dataflows

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~

- ~~HDInsight~~

- ~~Azure Data Lake Analytics~~

- Azure SSIS Integration Runtime

- Azure Data Factory (Mapping Data Flows)

- Azure Databricks

- Power BI dataflows

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~
- ~~HDInsight~~
- ~~Azure Data Lake Analytics~~
- Azure SSIS Integration Runtime
- Azure Data Factory (Mapping Data Flows)
- Azure Databricks (Delta)
- Power BI dataflows

# Datawarehousing in the cloud…

- ~~Azure SQL /w PolyBase~~

- ~~HDInsight~~

- ~~Azure Data Lake Analytics~~

- <mark>Azure SSIS Integration Runtime</mark>

- <mark>Azure Data Factory (Mapping Data Flows)</mark>

- <mark>Azure Databricks (Delta)</mark>

- ~~Power BI dataflows~~

# The battle!

- round #1: capabilities
- round #2: developer experience
- round #3: operator experience
- round #4: security
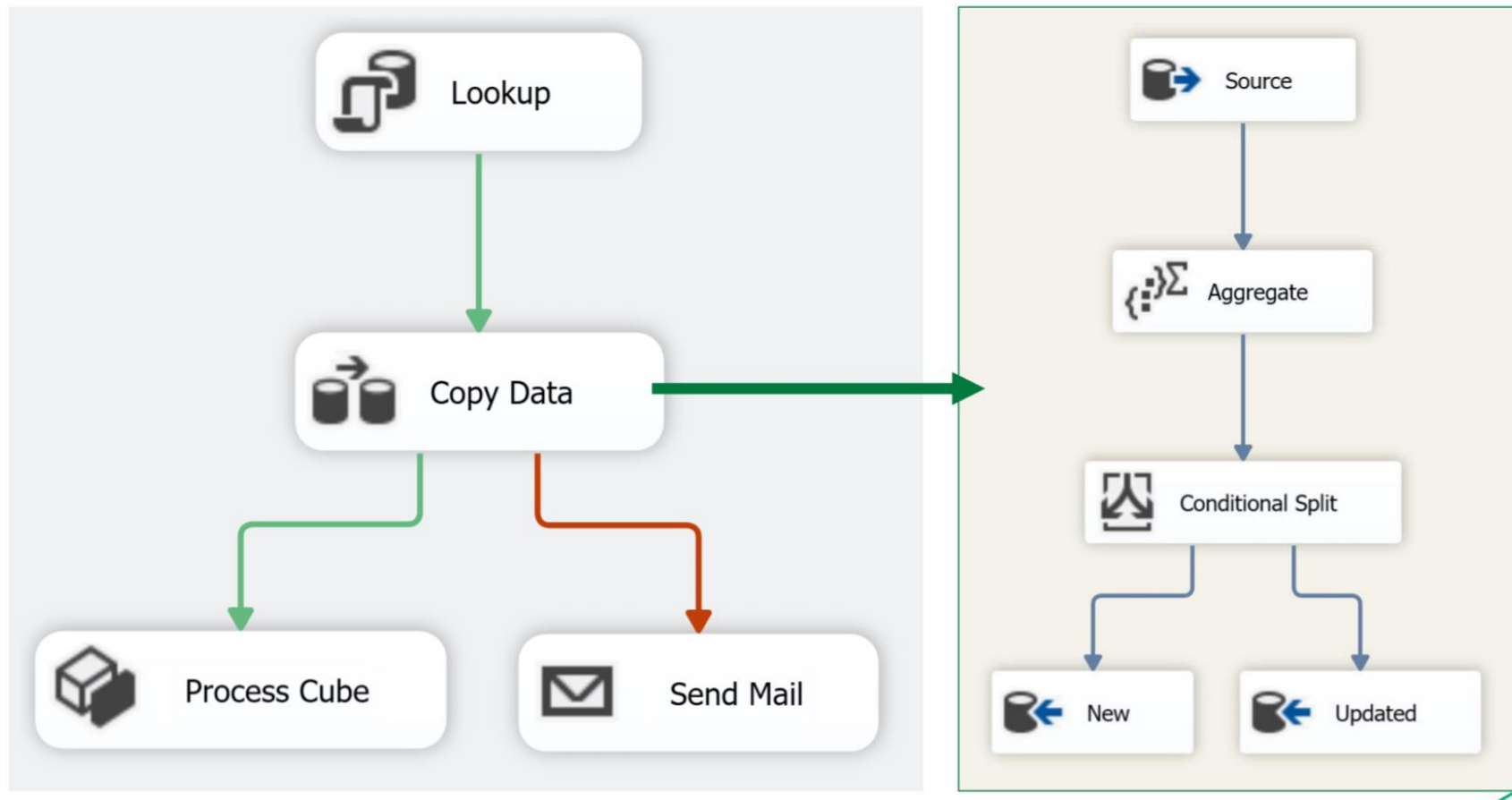- round #5: roadmap / future readiness
- round #6: coolness

# capabilities

round.set(1)

01

# Capabilities - SSIS in Azure

# Capabilities - SSIS in Azure

- Pros
  - Combination of simple & advanced options
  - Continue with existing solution / less initial investment
  - Mature tooling for meta-driven generation
  - Git integration in VS
- Cons:
  - Not serverless
  - No further development (continuity?)
  - Integration with Azure PaaS (data Source connectivity)
  - No advanced analytics / data science
  - No streaming data support

# Capabilities - ADF (MDF!)

# Capabilities - ADF (MDF!)

# Capabilities - ADF (MDF!)

# Capabilities - ADF (MDF!)

- Pros
  - Simplicity / easy to understand / no programming required
  - Low time to value
  - DevOps ready (bit clumsy design)
- Cons:
  - New / immature
  - Poor generation using meta-data driven frameworks
  - Limited set of (advanced) transformations
  - Not suited for complex tasks, fallback to functions/notebooks
  - Limited v-net support. MDF / Web Activities can't access Key Vault.
  - Publishing is a manual action in the browser
  - No advanced analytics / data science
  - No streaming data support

# Capabilities - Azure Databricks

- Azure Databricks is a **first party** service on Azure.
  - Unlike with other clouds, it is not an Azure Marketplace or a 3rd 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 Billing: You get a single bill from Azure
  - Azure Power BI: For rich data visualization

- Eliminates need to create a separate account with Databricks.

# Capabilities – Databricks Delta

Essentially, it's an optimized Spark table with SQL-like features:

- **ACID transactions**
- **DELETES / UPDATES / UPSERTS**
- **Statistics, data skipping and ZORDER clustering**

# Capabilities - Azure Databricks

- Pros
  - Extremely versatile and scalable
  - Easily add streaming data
  - Not only applicable for data engineering "unified analytics"
  - Interactive notebook experience
  - Cloud agnostic / open source
- Cons:
  - Steep learning curve
  - Not serverless (don't underestimate cluster management)
  - Poor Git integration
  - Longer time to value
  - Poor Service Principal support

# developer experience

round.set(2)

Dataminds test presentation

02

# Developer experience - SSIS

- Tool: Visual Studio (crashes / manual updates)

- Infrastructure-as-a-Service: Virtual Machine

- Good options to generate code

- Poor collaboration

- Testing via dataviewers

- Disconnected from ADF

- Source code XML

- Schema drift

# Developer experience - ADF MDF

- Tool: Browser
- Platform-as-a-Service: Azure Portal
- Simple to start
- Can feel limited
- Breaking changes
- Poor collaboration
- Seamless integration ADF
- Source code is JSON
- Testing via debug-mode
- Schema drift

# Developer expierience – Databricks

- Tool: Browser
- Platform-as-a-Service: Azure Portal
- Just code
- Good collaboration
- Almost anything is possible
- Multilingual
- Future VS Code support
- Schema drift

# operator expierience

round.set(3)

**03**

# Operator experience - SSIS in Azure

- SQL Server Management Studio (SSMS)

# Operator experience - SSIS in Azure

- Deployment can be complicated
- Debugging / troubleshooting can be intimidating
- Limited integration of monitoring (ADF / Databricks)

# Operator experience - MDF

- Internet Browser (Azure Portal)
- (Azure) PowerShell

# Operator experience - Databricks

- Deployment complicated (clusters/notebooks)
- Internet Browser (Azure Portal)
- (Azure) PowerShell
- Debugging / troubleshooting can be intimidating
- Limited integration of monitoring (ADF / Databricks)

# pricing

round.set(4)

04

# Pricing - SSIS in Azure

- Integration Runtime costs ("it depends")
- Benefit from existing SQL Server licensing

# Pricing - SSIS in Azure

## SQL Server Integration Services Enterprise E-series V3 VM

| INSTANCE | CORES | RAM | TEMPORARY STORAGE | LICENSE INCLUDED PRICE PER NODE | PRICE WITH AZURE HYBRID BENEFIT PER NODE (% SAVINGS) |
|----------|-------|-----|-------------------|--------------------------------|------------------------------------------------------|
| E2 V3 | 2 | 16.00 GiB | 50 GiB | €1.568/hour | €0.335/hour (~79%) |
| E4 V3 | 4 | 32.00 GiB | 100 GiB | €1.903/hour | €0.670/hour (~65%) |
| E8 V3 | 8 | 64.00 GiB | 200 GiB | €3.806/hour | €1.340/hour (~65%) |
| E16 V3 | 16 | 128.00 GiB | 400 GiB | €7.612/hour | €2.679/hour (~65%) |
| E32 V3 | 32 | 256.00 GiB | 800 GiB | €15.230/hour | €5.356/hour (~65%) |
| E64 V3 | 64 | 432.00 GiB | 1,600 GiB | €29.830/hour | €10.096/hour (~66%) |

https://azure.microsoft.com/en-us/pricing/details/data-factory/ssis/

# Pricing - SSIS in Azure

- License Visual Studio

(don't forget about the DevOps server)

## For organizations

An unlimited number of users within an organization can use Visual Studio Community for the following scenarios: in a classroom learning environment, for academic research, or for contributing to open source projects.

For all other usage scenarios:
In non-enterprise organizations, up to five users can use Visual Studio Community. In enterprise organizations (meaning those with >250 PCs or >$1 Million US Dollars in annual revenue), no use is permitted beyond the open source, academic research, and classroom learning environment scenarios described above.

# Pricing - SSIS in Azure

• Development workload (3 days a week)

# Pricing - SSIS in Azure

- Production workload (2h a day)



TIER:

Enterprise ⌄

INSTANCE:

E4V3: 4 Cores(s), 32 GB RAM, 100GB Disks, €1.9025/hour ⌄

Save up to 30% with SQL Server licenses you already own with
Azure Hybrid Benefit for SQL Server

| 2 | ✕ | 60 | = | €228.30 |

Virtual machines

Hours ⌄

# ADF MDF - Pricing

- Debug Mode:
  - "Preview Pricing"
  - 8 cores default
  - $0.112 / hour
  - 60 minutes default Time To Live (TTL)
  - Example dev-day: 10h.  x 8 (cores) x $0.112 = **$8.96**
- Transform data in Blob Store (scheduled):
  - "Preview Pricing"
  - 8 cores default
  - $0.112 / hour
  - 10 minutes default Time To Live (TTL)
  - Example: 10m. compute + 10m. TTL = 0,33h. x 8 (cores) x $0.112 = **$0.299**

# Azure Databricks - Pricing

## Standard tier features

| FEATURE | DATA ANALYTICS | DATA ENGINEERING | DATA ENGINEERING LIGHT |
|---------|----------------|------------------|------------------------|
|         | Interactive workloads to analyze data collaboratively with notebooks | Automated workloads to run fast and robust jobs via API or UI | Automated workloads to run robust jobs via API or UI |

## Premium tier features

| FEATURE | DATA ANALYTICS | DATA ENGINEERING | DATA ENGINEERING LIGHT |
|---------|----------------|------------------|------------------------|
|         | Interactive workloads to analyze data collaboratively with notebooks | Automated workloads to run fast and robust jobs via API or UI | Automated workloads to run robust jobs via API or UI |
|         | Includes standard features | Includes standard features | Includes standard features |
| Role-based access control for notebooks, clusters, jobs, and tables | ✔ | ✔ | ✔ |
| JDBC/ODBC Endpoint Authentication | ✔ | ✔ | ✔ |
| Audit logs (In preview) | ✔ | ✔ | ✔ |

# Azure Databricks - Pricing

- Data Analytics:
  - Interactive Clusters only here
  - Power BI connection to data in cluster
  - Notebook collaboration experience

- 'Data Engineering Light'
  - Delta not available
  - Notebooks not available (also no scheduling of notebooks)

- Premium:
  - Role-based access control for notebooks, clusters, jobs, and tables
  - Audit Logs (preview)
  - JDBC/ODBC Endpoint Authentication

# Azure Databricks - Pricing

- Development: Premium Tier - Data Analytics
- Production:    Premium Tier – Data Engineering

## Pay as you go

Azure Databricks bills* you for virtual machines (VMs) provisioned in clusters and Databricks Units (DBUs) based on the VM instance selected. A DBU is a unit of processing capability, billed on a per-second usage. The DBU consumption depends on the size and type of instance running Azure Databricks.

| WORKLOAD | DBU PRICES—STANDARD TIER | DBU PRICES—PREMIUM TIER |
| --- | --- | --- |
| Data Analytics | €0.34/DBU-hour | €0.464/DBU-hour |
| Data Engineering | €0.13/DBU-hour | €0.253/DBU-hour |
| Data Engineering Light | €0.06/DBU-hour | €0.186/DBU-hour |

*In addition to virtual machines, Azure Databricks will also bill for managed, disk, blob storage, Public IP Address.

# Azure Databricks - Pricing

- Development: Premium Tier - Data Analytics

- Workload: 5 days a week

INSTANCE:

| F4: 4 Core(s), 8 GB RAM, 0.5 Databricks Unit(s), €0.191/hour | ⌄ |
|---|---|

## Billing Option

Save up to 72% on pay-as-you-go prices with 1-year or 3-year Reserved Virtual Machine Instances. Reserved Instances are great for applications with steady-state usage and applications that require reserved capacity. Learn more about Reserved VM Instances pricing.

- ● Pay as you go
- ○ 1 year reserved (~27% savings)
- ○ 3 year reserved (~51% savings)

| 2 | **✕** | 200 | **=** | €76.57 |
|---|---|---|---|---|
| Virtual machines | | Hours ⌄ | | Per month |

## DBU (Databricks Unit) ⓘ

| 1.00 | **✕** | €0.464 | **✕** | 200 | **=** | €92.76 |
|---|---|---|---|---|---|---|
| DBU | | Per DBU per hour | | Hours ⌄ | | |

Sub-total     €169.33

# Azure Databricks - Pricing

- Production :
  Premium Tier –
  Data Engineering

- Workload:
  2 hours a day

INSTANCE:

F16: 16 Core(s), 32 GB RAM, 2 Databricks Unit(s), €0.767/hour

## Billing Option

Save up to 72% on pay-as-you-go prices with 1-year or 3-year Reserved Virtual Machine Instances. Reserved Instances are great for applications with steady-state usage and applications that require reserved capacity. Learn more about Reserved VM Instances pricing.

- ● Pay as you go
- ○ 1 year reserved (~27% savings)
- ○ 3 year reserved (~51% savings)

| 2 | × | 62 | | = | €95.05 |
|---|---|----|----|---|--------|
| Virtual machines | | Hours | | | Per month |

## DBU (Databricks Unit) ⓘ

| 4.00 | × | €0.253 | × | 62 | | = | €62.74 |
|------|---|--------|---|----|----|---|--------|
| DBU | | Per DBU per hour | | Hours | | | |

Sub-total        €157.79

# roadmap

round.set(5)

05

# Roadmap - SSIS in Azure

- ?

# Roadmap – ADF MDF

- Active Monitoring (watch progress live)

# Roadmap – Databricks

- C# as notebook language
- Integration with Visual Studio Code

# coolness

round.set(6)

06

# Recruiting

- People are increasingly looking for new tooling in job offers:
  - Azure
  - Azure Databricks
  - Azure Data Factory
  - Data Lake
  - Datawarehouse
  - DevOps

# Job offers - examples

**Wat je doet als Technisch Specialist Data & Analytics**

Je adviseert klanten over de technische (on)mogelijkheden en randvoorwaarden voor data & analytics. Samen met je team ontwerp, realiseer en implementeer je de data science en analytics toepassingen, slimme big-data platformen en self service BI oplossingen die zij nodig hebben. Daarbij hou jij je onder andere bezig met het ontwerpen van architecturen, deployment van software en Azure services en het configureren van ETL en Azure Data Factory pipelines en data engineering. Je werkt voor diverse klanten van Macaw, volgt markttrends en innovaties en werkt behalve aan het succes van onze klanten, ook aan je eigen ontwikkeling en die van je collega's. Je maakt onderdeel uit van een groot, deskundig en gedreven team van Functionele- en Technische Specialisten, Consultants en Architecten, waarbinnen je zelf vorm geeft aan de rol, het specialisme en de richting die het beste bij jou en je ambities past.

# Job offers - examples

Innovatie en data gedreven werken raakt de kern van onze klanten. Data goed kunnen organiseren en er business waarde uit kunnen halen is niet meer een bijzaak, maar een hoofdzaak geworden. Door de dynamische ontwikkelingen in de Big Data & Analytics markt, is het zeer uitdagend geworden om te weten wanneer je wat moet gebruiken. Als Big Data Engineer werk jij samen met onze opdrachtgevers en partners, zoals: AWS, Azure, Databricks, Cloudera & Hortonworks, om orde in deze chaos te scheppen - je bent bepalend voor het success van onze Data Science, IoT, Realtime en Data Lake projecten.

- Een afgeronde hbo of wo opleiding in de richting van computerwetenschappen, software engineering of andere technische studie.
- Minimaal 2 jaar werkervaring binnen IT.
- Ervaring met een of meerdere big data technologieën (bv. Databricks Spark, Cloudera / Hortonworks Hadoop)
- Ervaring met cloud computing omgeving (bv. AWS, Azure)
- En uitstekende Nederlandse én Engelse communicatie skills.

# Job offers - examples

**Wat vragen wij van je?**

- Je bent enthousiast, leergierig en oplossingsgericht
- Je hebt HBO/WO werk- en denkniveau
- Je hebt kennis van en ervaring met het Microsoft (Azure) Data platform. In ieder geval bestaande uit:
    - (Azure) SQL Server Database
    - (Azure) SQL Server Analysis Services
    - (Azure) SQL Server Integration Services
    - Azure Data Factory
    - Azure Data Lake (Store & Analytics)
    - Azure Data Warehouse
- Je hebt gewerkt met modelleringstechnieken van Kimball en Dan Linstedt (Data Vault)
- Je hebt kennis van (Microsoft) AI technieken zoals:
    - Cognitive Services
    - Machine Learning Services
    - R
- Je hebt kennis van (Microsoft) reporting tools zoals:
    - Microsoft SQL Server Reporting Server
    - Microsoft Power BI
- Je hebt ervaring als consultant en in een (pre-/ technical)sales rol
- Je werkt graag samen in een team

# The battle!

- round #1: capabilities
- round #2: developer experience
- round #3: operator experience
- round #4: security
- round #5: roadmap / future readiness
- round #6: coolness

# Thank You

# What do you think?

1. Open the form
2. Provide constructive feedback
3. Be eligible for an ==amazing prize==!

http://bit.ly/dataMindsConnectSession
bit.ly is CASE SENSITIVE!

99

# Q&A

# Our Partners