

# *Emerging DataOps Principles and Technologies for Analytics*

**MIDWEST ARCHITECTURE COMMUNITY COLLABORATION**

**Tim Garrod**



midwest  
architecture community  
collaboration

# Agenda

- What is DataOps?
- What is driving DataOps?
- Enabling DataOps...
- DataOps technologies and functional requirements

# What is DataOps?



# DevOps: Accelerating time to value

## WATERFALL



## AGILE



## DEVOPS



- Built upon the Agile development movement
- Collaboration between development and operations
- Continuous integration & delivery frameworks
- ...



**Agile, Collaborative & Continuous**

# WHAT IS DATAOPS? WE ARE IN THE EARLY DAYS

Figure 1. Hype Cycle for Data Management, 2018

Gartner



Plateau will be reached:

- less than 2 years
- 2 to 5 years
- 5 to 10 years
- ▲ more than 10 years
- ⊗ obsolete before plateau

< 1%  
ADOPTION

# LINKEDIN RESULTS

1,082  
DataOps

39,830  
Unicorns

488,972  
DevOps

# WHAT IS DATAOPS

Emerging discipline to build and manage efficient, effective data pipelines

PEOPLE

PROCESS

TECHNOLOGY



**CODE**



**TOOLS**



**INFRASTRUCTURE**

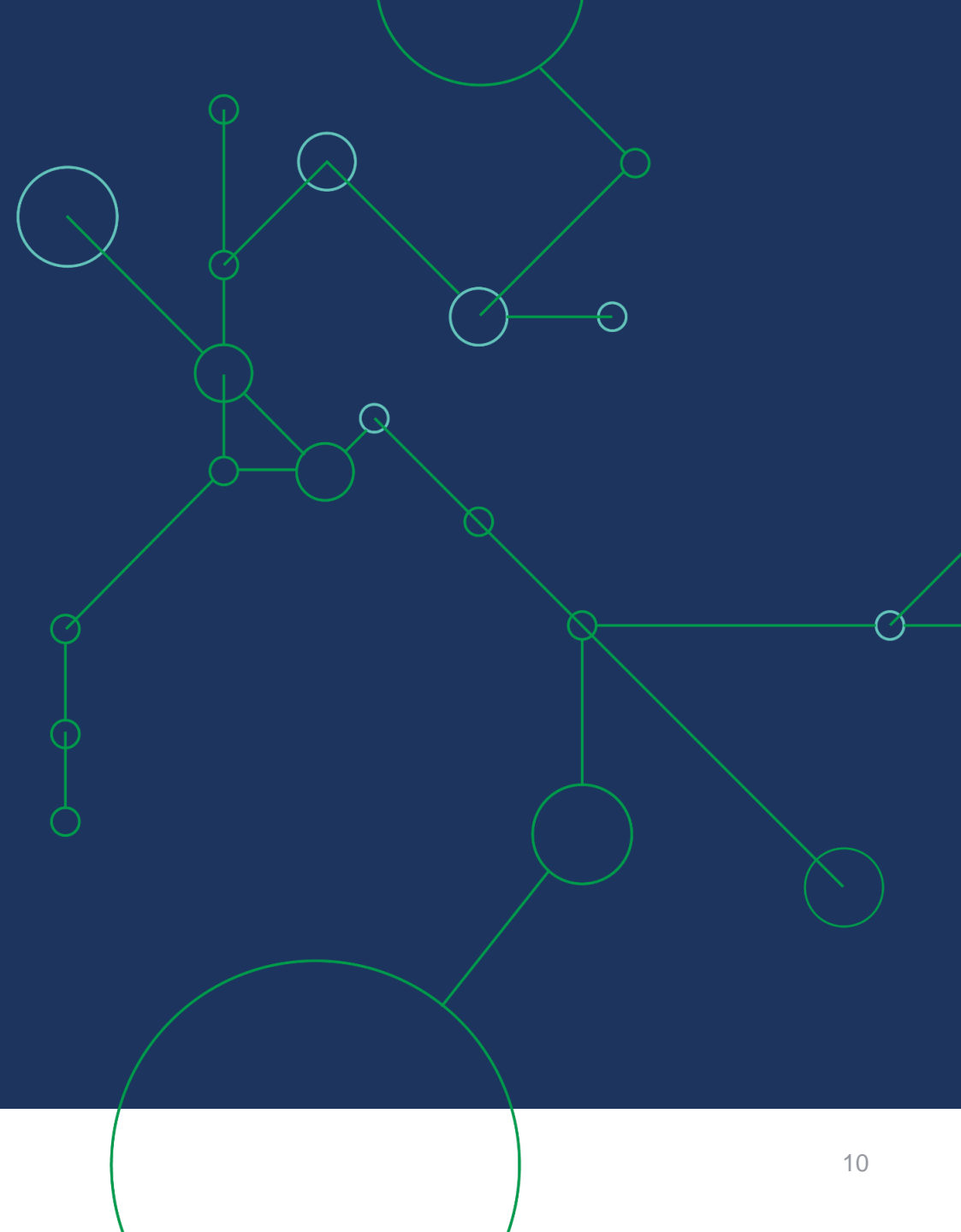


**DATA**

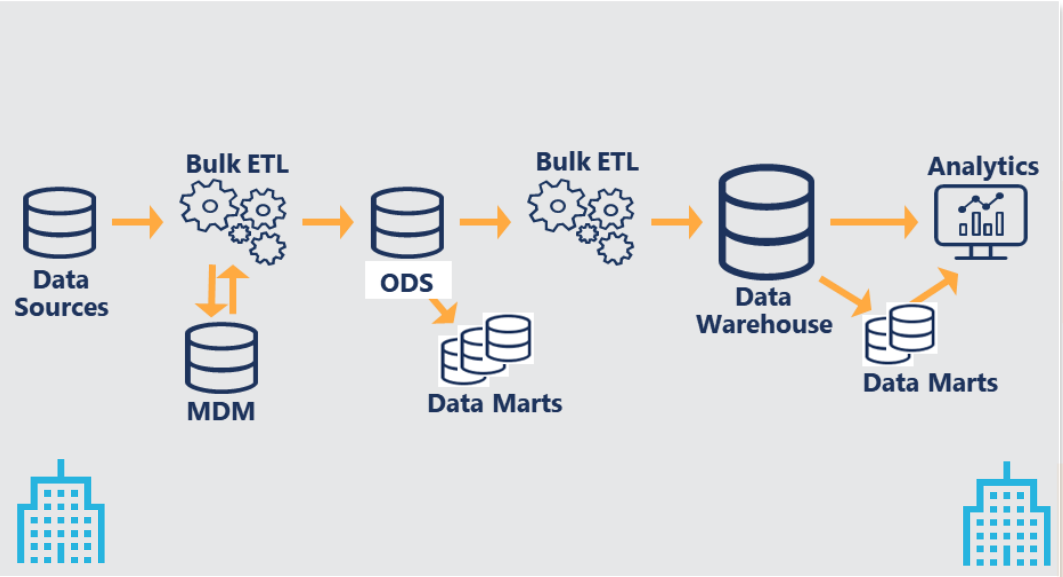
- **Applies DevOps principles of agile development and continuous integration**
- **Seeks to improve collaboration between data managers and consumers**
- **Continuous data delivery in the data pipeline**



# What is driving DataOps?



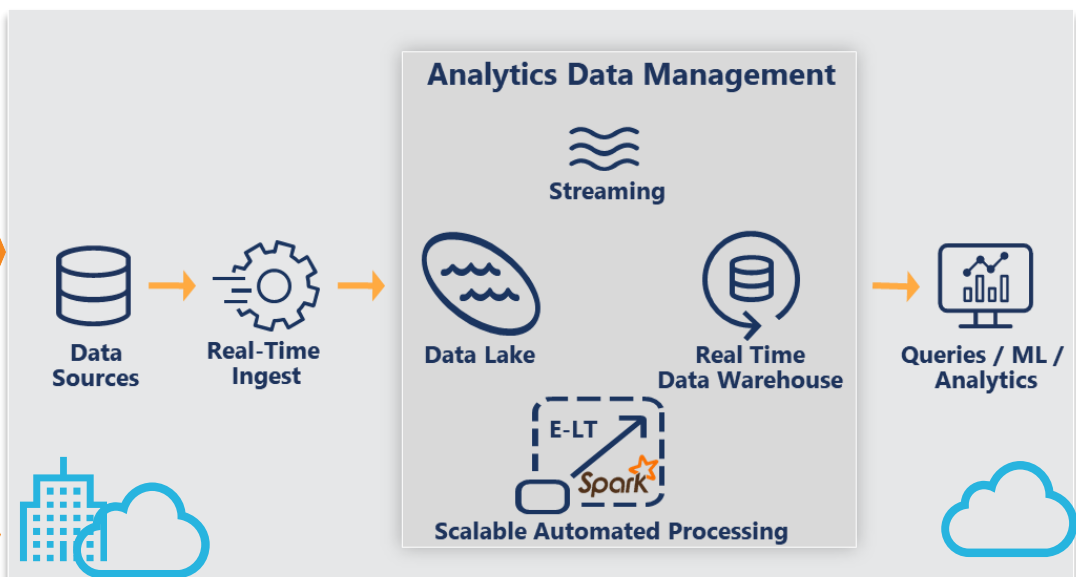
# DATA ARCHITECTURE EVOLUTION



Architecture

Platform

Processes



## Legacy Data Warehouse Architecture

- ✗ Bulk data movement
- ✗ Brittle **hand-coded** ETL
- ✗ **Monolithic** / appliance driven architectures (one size fits all)
- ✗ **Slow time to market / react to change**

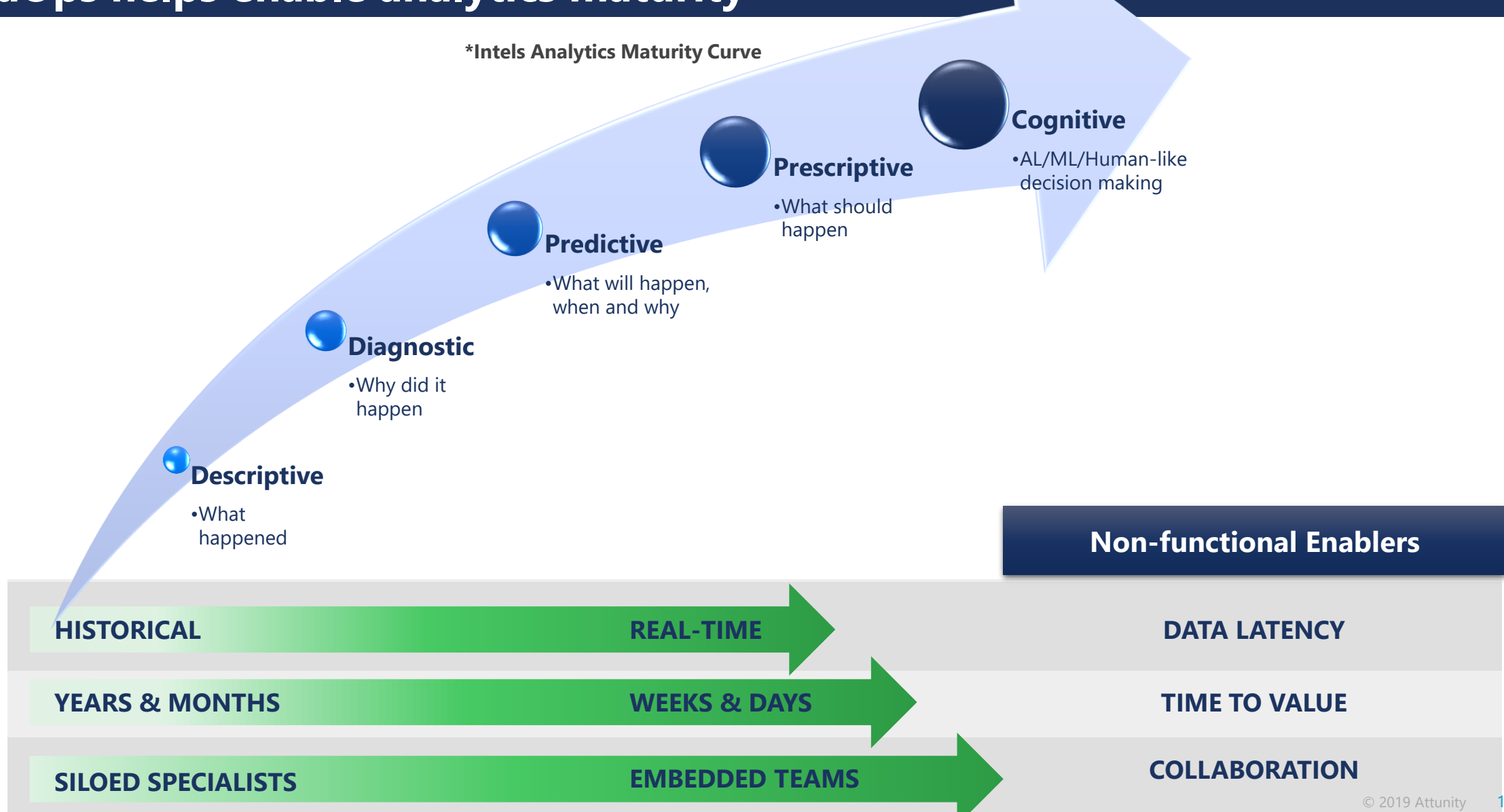
## Modern Analytics Data Management

- ✓ **Real-time** data movement
- ✓ **Automated** design and code generation
- ✓ Use-case driven **scalable** technologies
- ✓ **Faster Time to Market / Value**

# ANALYTICS MATURITY CURVE

## DataOps helps enable analytics maturity

\*Intels Analytics Maturity Curve



# HOW WE USED TO DO THINGS!



```
{code:"etl", type: "manual", speed:"batch"}
```



We need data!!

That's not exactly what I wanted

We don't need that anymore

Why is this data old?

# WHY DATAOPS?

- Increasing analytics requirements create complexity and data flow bottlenecks
- Data consumers drive demands that IT cannot meet with existing processes and technologies
- Projects are failing due to this friction

*"In every pipeline, data must be identified, captured, formatted, tagged, validated, profiled, cleaned, transformed, combined, aggregated, secured, cataloged, governed, moved, queried, visualized, analyzed, and acted upon. Phew!"*

Wayne Eckerson PRESIDENT, ECKERSON GROUP



## RISING CHALLENGES

**DATA VOLUME,  
VARIETY,  
VELOCITY**

**NEW  
PLATFORMS**

**SPEED OF  
CHANGE**

**CODING  
COMPLEXITY**

# Enabling DataOps



# YOU CAN'T BUY DATAOPS

PEOPLE

PROCESS

TECHNOLOGY



**CODE**



**TOOLS**



**INFRASTRUCTURE**



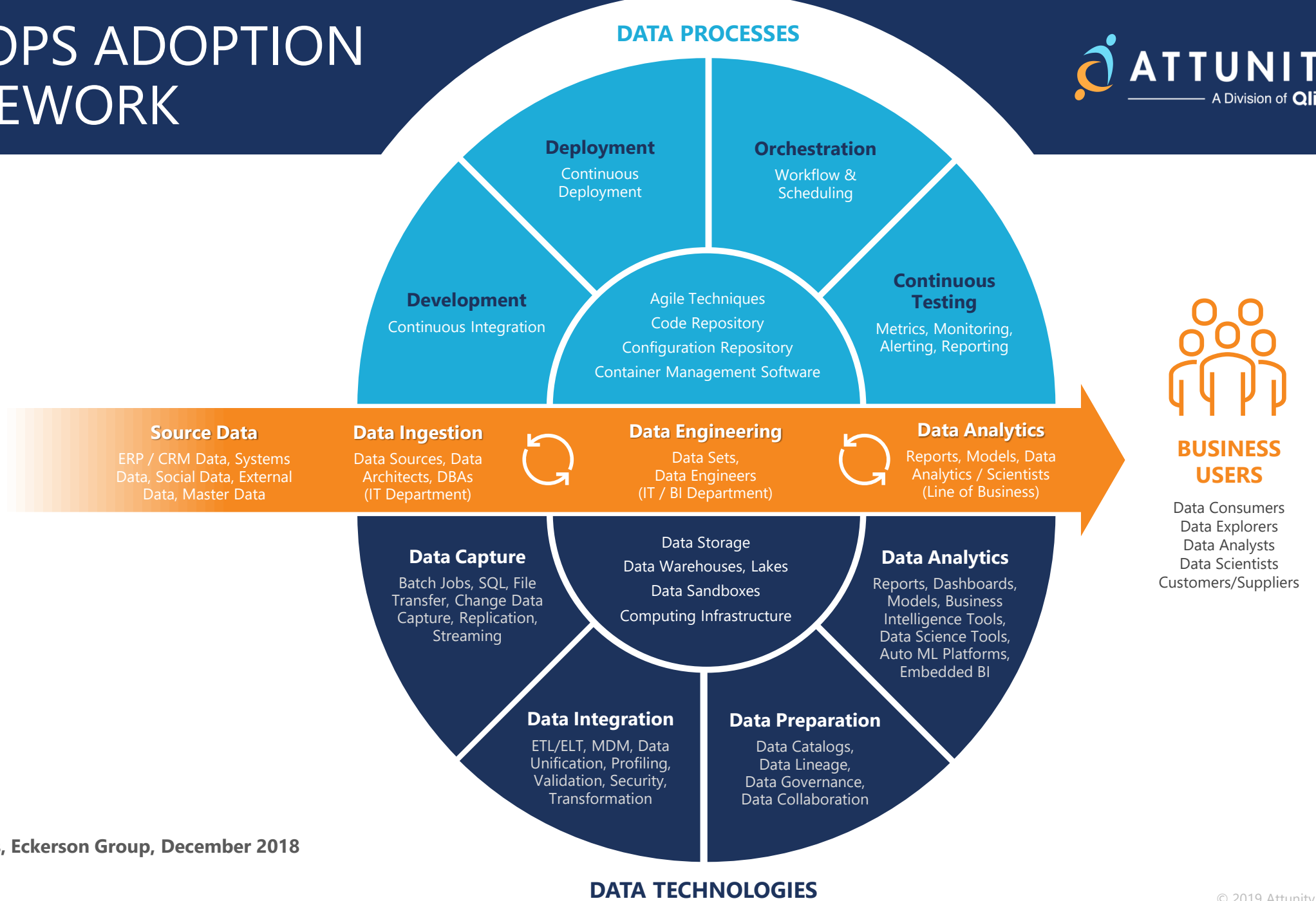
**DATA**

**You should invest in modern technology which supports key DataOps principles and significantly accelerates your time to insight**

# DATAOPS ADOPTION FRAMEWORK



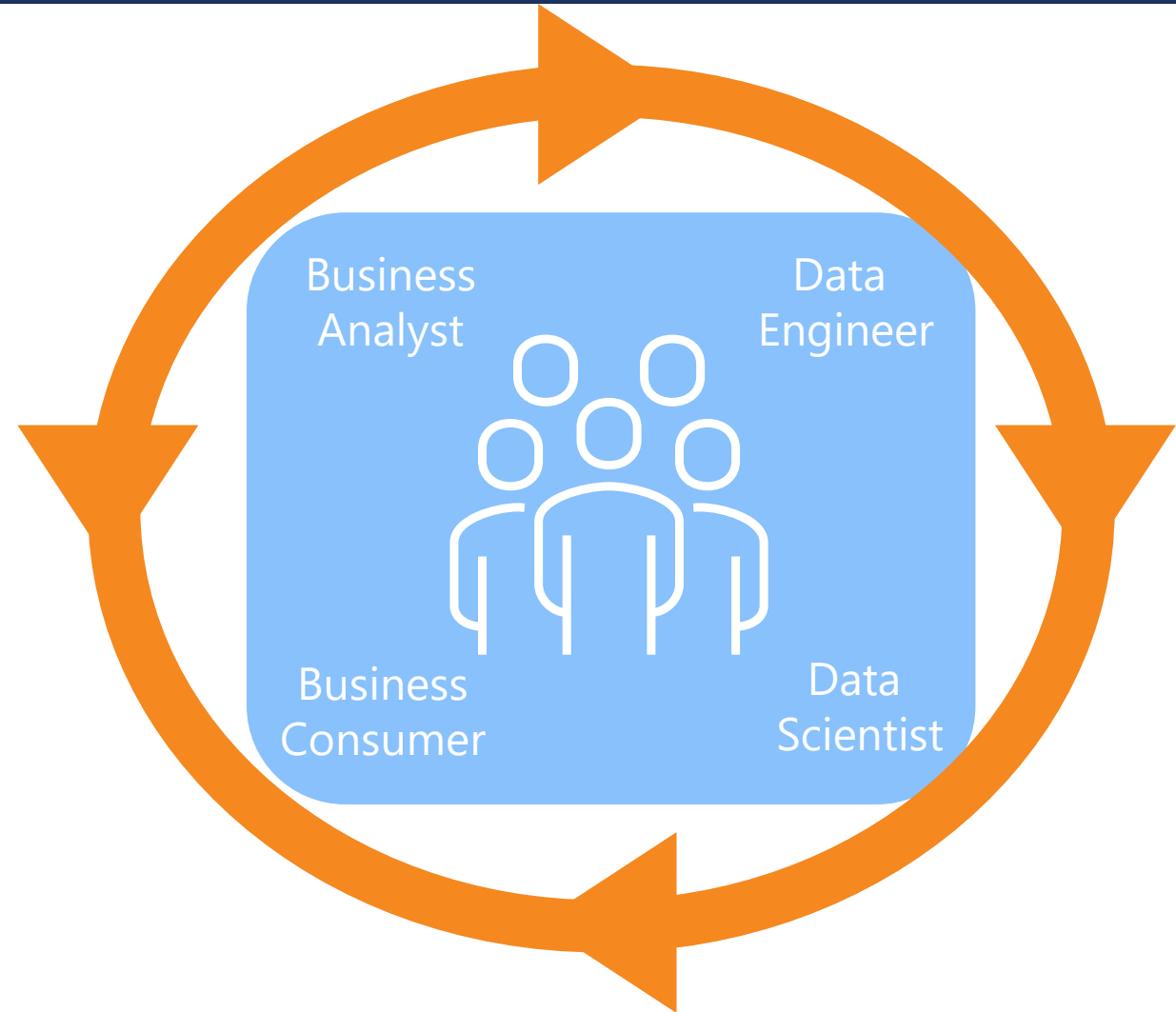
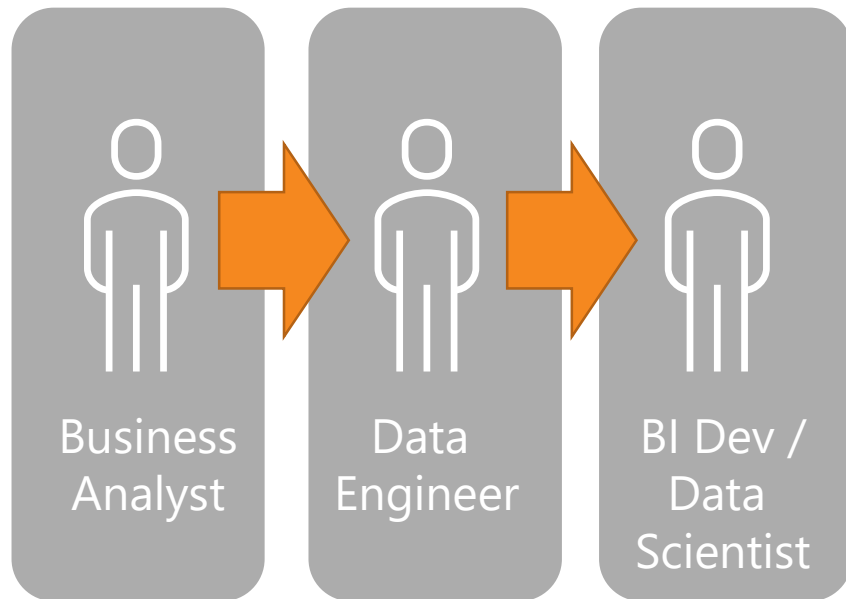
## DATA PIPELINES



Source: Diving into DataOps, Eckerson Group, December 2018



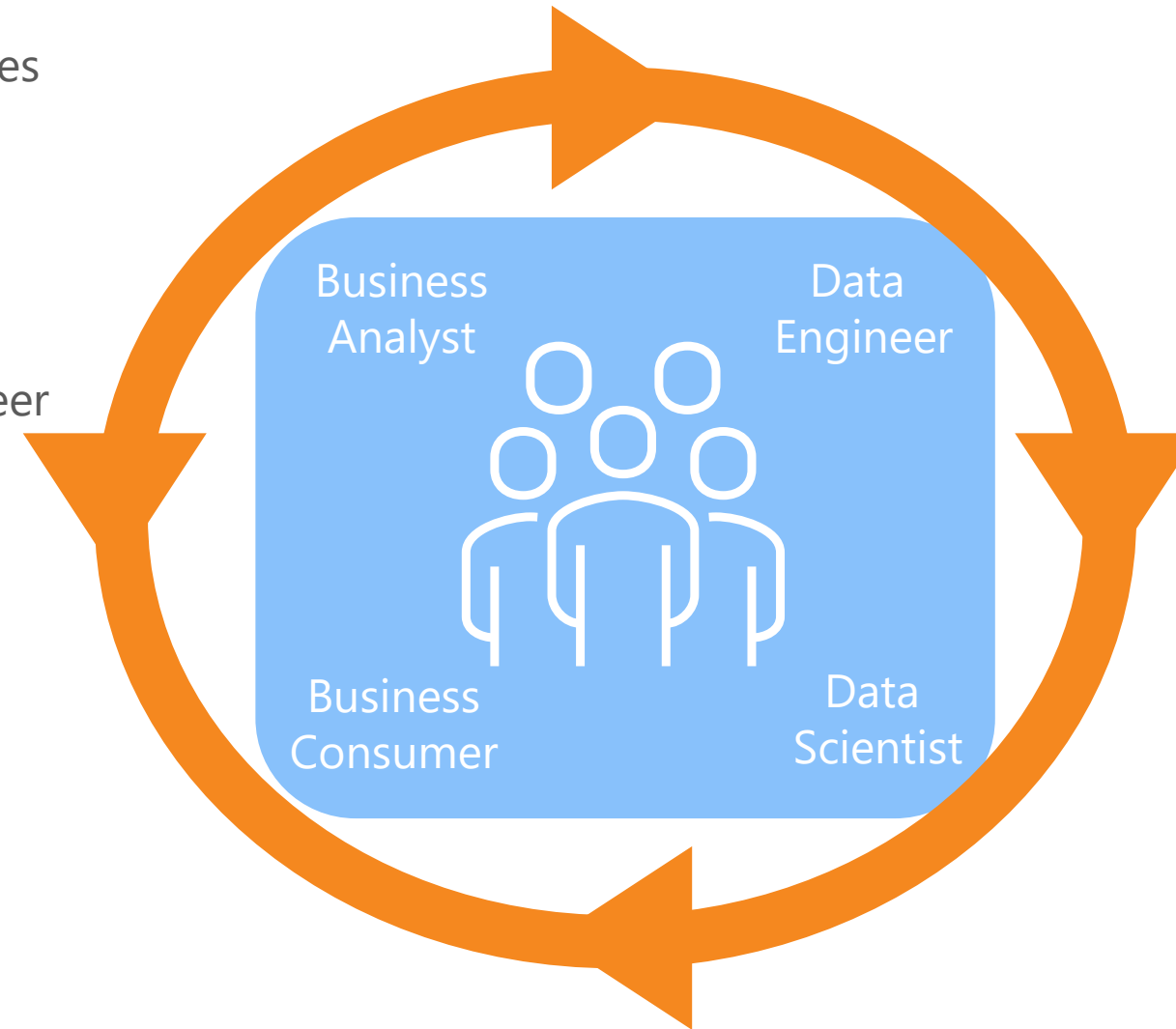
# DATAOPS team structure



From siloed to centralized – DataOps teams need to be enabled to collaborate

# DATAOPS team structure

- **Data Engineer**
  - Responsible for building Data Lakes, Data Warehouses
  - DBA, ETL Engineer
  - ETL, Spark programming
- **Data Analyst**
  - Responsible for visualizations, charts, dashboards
  - BI Engineer, Report analyst, Data visualization engineer
  - Data storytelling
- **Data Scientist**
  - Responsible for algorithms, models
  - ML Engineer, AI programmers w/ domain expertise
  - Spark, R, Python, Notebooks
- **Business Consumer\***
  - Virtual team member
  - Responsible for business feedback loop



# DATAOPS BENEFITS FROM A MODERN DATA INFRASTRUCTURE

## MODERN ANALYTICS

AI/ML

IoT

Predictive

Real-Time

## MODERN PLATFORMS

Big Data

Cloud

Data Lakes

Streaming

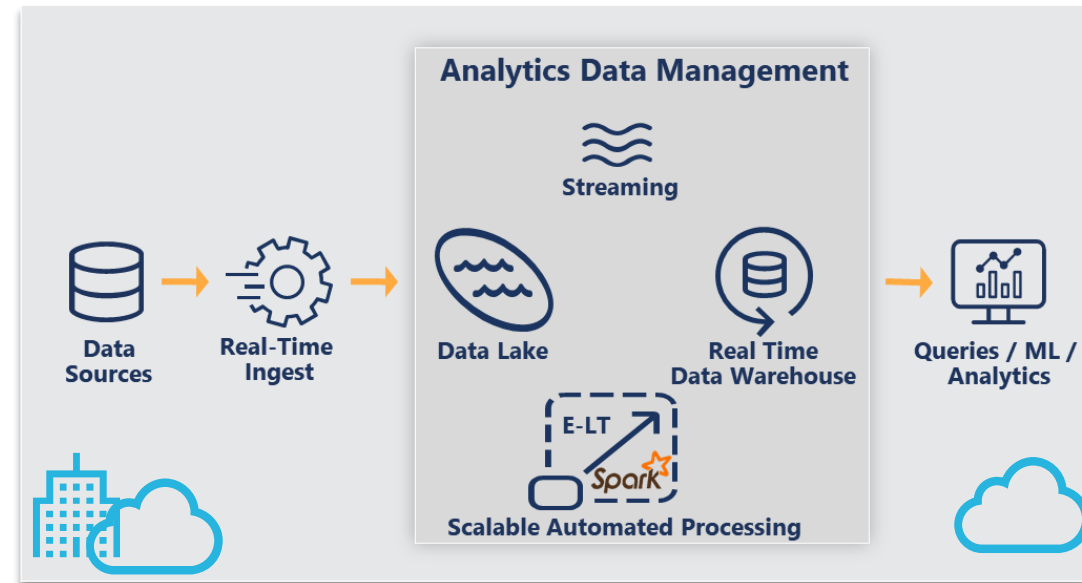
## MODERN INTEGRATION

# DATAOPS BENEFITS FROM A MODERN DATA INFRASTRUCTURE

## MODERN **INTEGRATION**

# A MODERN DATA INFRASTRUCTURE

- Cloud platforms – AWS, Azure, Google, Snowflake
- Data democratization



## Modern Analytics Data Management

- ✓ **Real-time** data movement
- ✓ **Automated** design and code generation
- ✓ Use-case driven **scalable** technologies
- ✓ **Faster Time to Market / Value**

## MODERN INTEGRATION

**1** Continuous

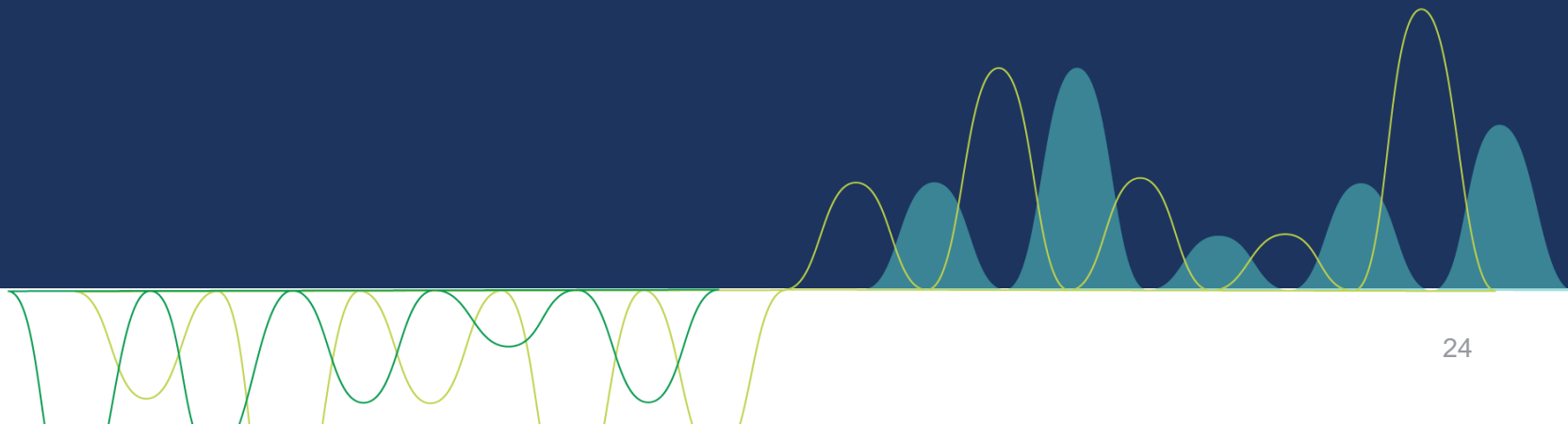
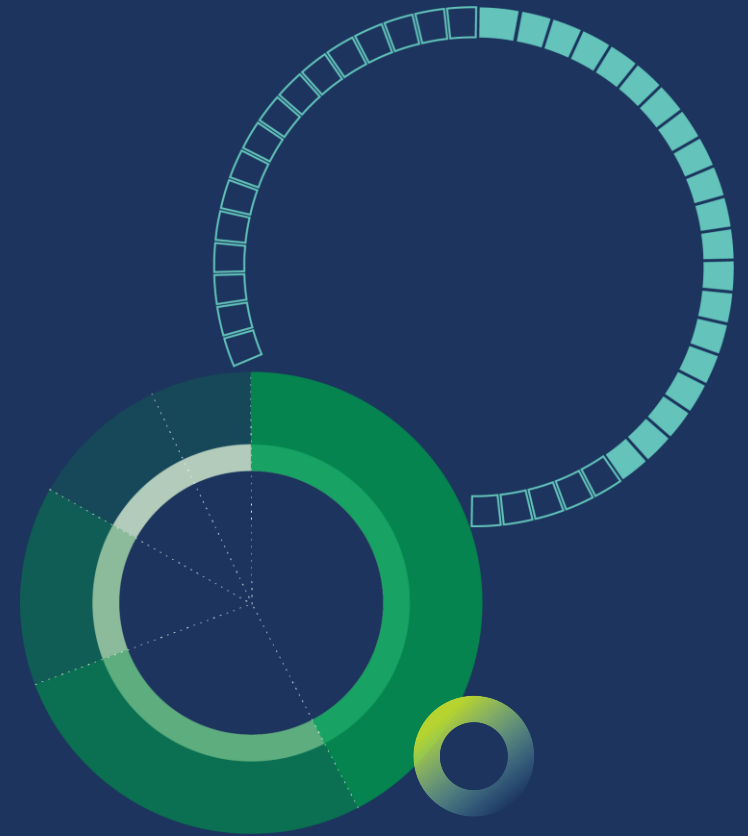
**2** Universal

**3** Automation

**4** Agility

**5** Trust

# DataOps technologies



# DATAOPS REQUIREMENT #1

## CONTINUOUS INTEGRATION

### MODERN INTEGRATION

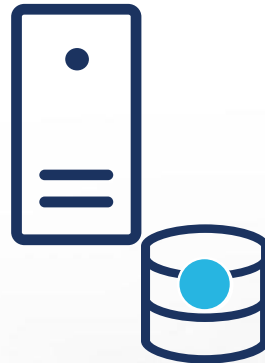
**1** Continuous

**2** Universal

**3** Automation

**4** Agility

**5** Trust



AGENTLESS CDC



REAL-TIME STREAMS



STREAM DATA  
AND METADATA  
CHANGES



MINIMAL  
IMPACT USING  
TRANSACTION LOGS  
WITH NO AGENTS



OPTIMIZED FOR  
EVERY SOURCE  
AND TARGET



# DATAOPS REQUIREMENT #2

## UNIVERSAL

### MODERN INTEGRATION

1 Continuous





2 Universal

3 Automation

4 Agility

5 Trust

### 30 SOURCES

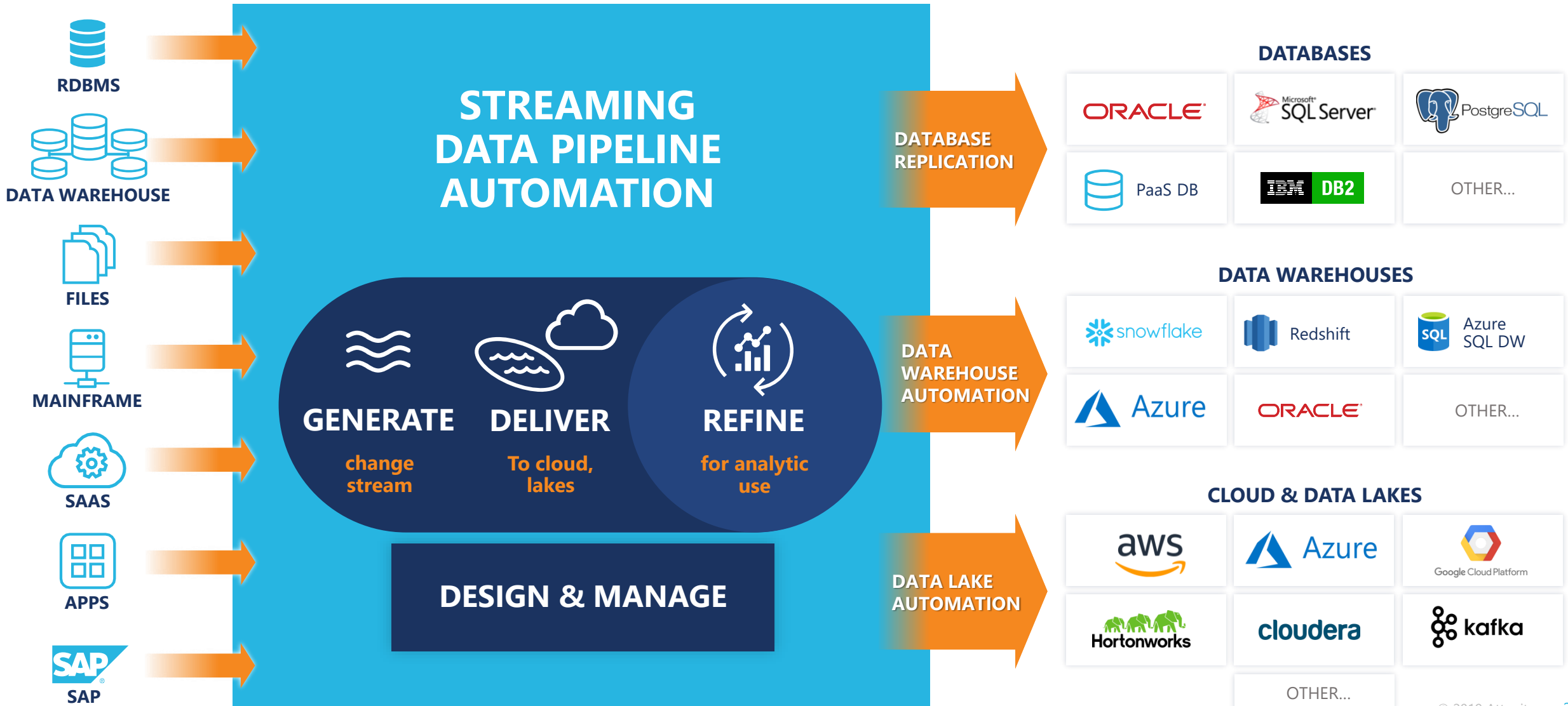
 <p><b>DATABASE</b></p> <p>Oracle SQL Server DB2 iSeries DB2 z/OS DB2 LUW MySQL PostgreSQL Sybase ASE Informix ODBC</p>	 <p><b>EDW</b></p> <p>Exadata Teradata Netezza Vertica Pivotal</p>	 <p><b>HADOOP</b></p> <p>Hortonworks Cloudera MapR</p>
 <p><b>SAP</b></p> <p>ECC on Oracle ECC on SQL ECC on DB2 ECC on HANA S4 HANA</p>	 <p><b>CLOUD</b></p> <p>AWS RDS Amazon Aurora Amazon Redshift Salesforce</p>	 <p><b>OTHER LEGACY</b></p> <p>SQL/MP Enscribe RMS</p>
 <p><b>MAINFRAME</b></p> <p>DB2 for z/OS IMS/DB VSAM</p>	 <p><b>FLAT FILES</b></p> <p>Delimited (e.g., CSV, TSV)</p>	

### 40 TARGETS

 <p><b>DATABASE</b></p> <p>Oracle SQL Server DB2 LUW MySQL PostgreSQL Sybase ASE Informix MemSQL</p>	 <p><b>EDW</b></p> <p>Microsoft PDW Exadata Teradata Netezza Vertica Sybase IQ Amazon Redshift Actian Vector SAP HANA</p>	 <p><b>CLOUD</b></p> <p>Amazon RDS Amazon Redshift Amazon S3 Amazon Aurora Google Cloud SQL Azure SQL DW Azure SQL DB Azure MySQL Azure PostgreSQL Snowflake</p>
 <p><b>STREAMING</b></p> <p>Kafka Azure Event Hubs MapR-ES AWS Kinesis</p>	 <p><b>SAP</b></p> <p>HANA</p>	 <p><b>FLAT FILES</b></p> <p>Delimited (e.g., CSV, TSV)</p>
 <p><b>HADOOP</b></p> <p>Hortonworks Cloudera MapR Amazon EMR HDInsight</p>		

# DATAOPS REQUIREMENT #2

## UNIVERSAL



### MODERN INTEGRATION

1 Continuous

2 Universal

3 **Automation**

4 Agility

5 Trust

## AUTOMATE THE DATA PIPELINE

GENERATE



DELIVER



REFINE



MANAGE

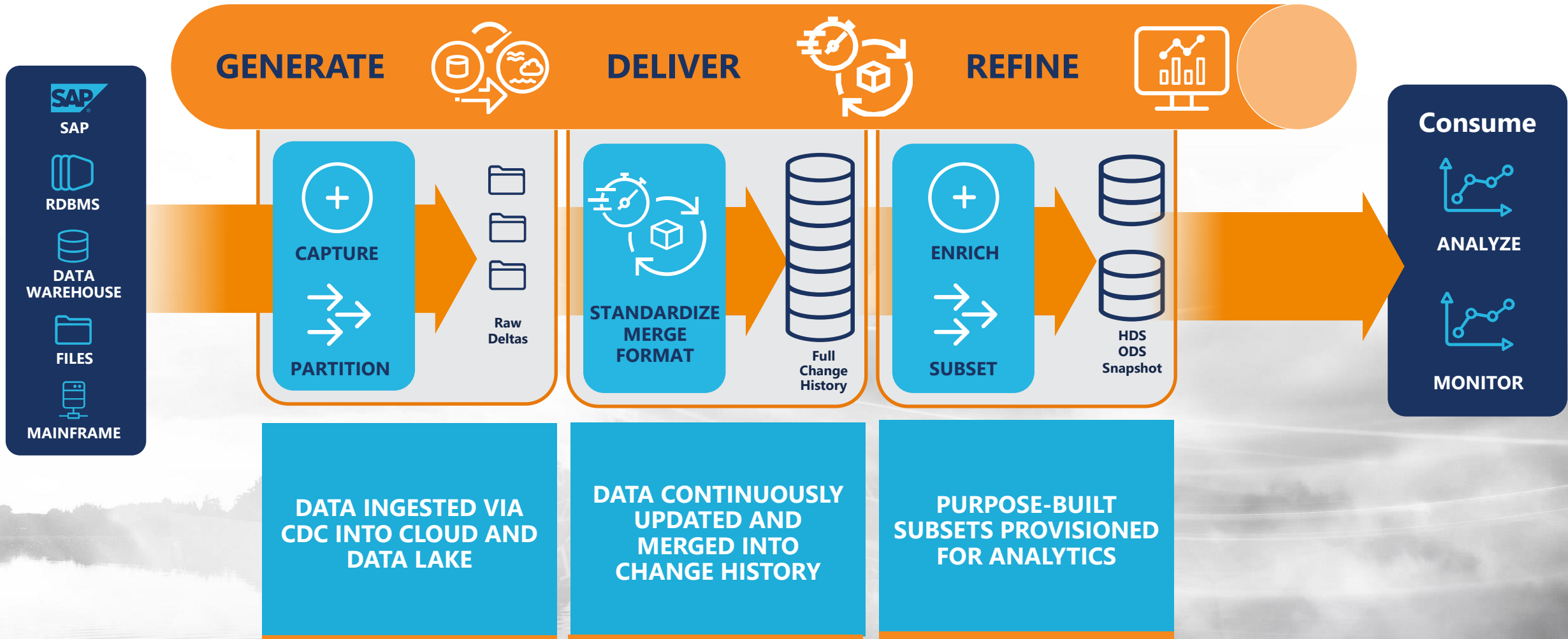
**AUTOMATE MULTI-STAGE PROCESSING WITHOUT CODING**

**ADAPTS TO SOURCE AND TARGET CHANGES**

**HETEROGENEOUS & DISTRIBUTED WORKLOADS**

# DATAOPS REQUIREMENT #3

## AUTOMATION



### MODERN INTEGRATION

1 Continuous

2 Universal

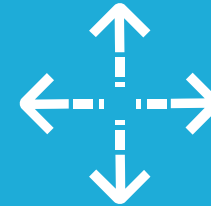
3 Automation

4 **Agility**

5 Trust



**RUNS IN THE  
CLOUD(S), ON-  
PREM, OR BOTH**



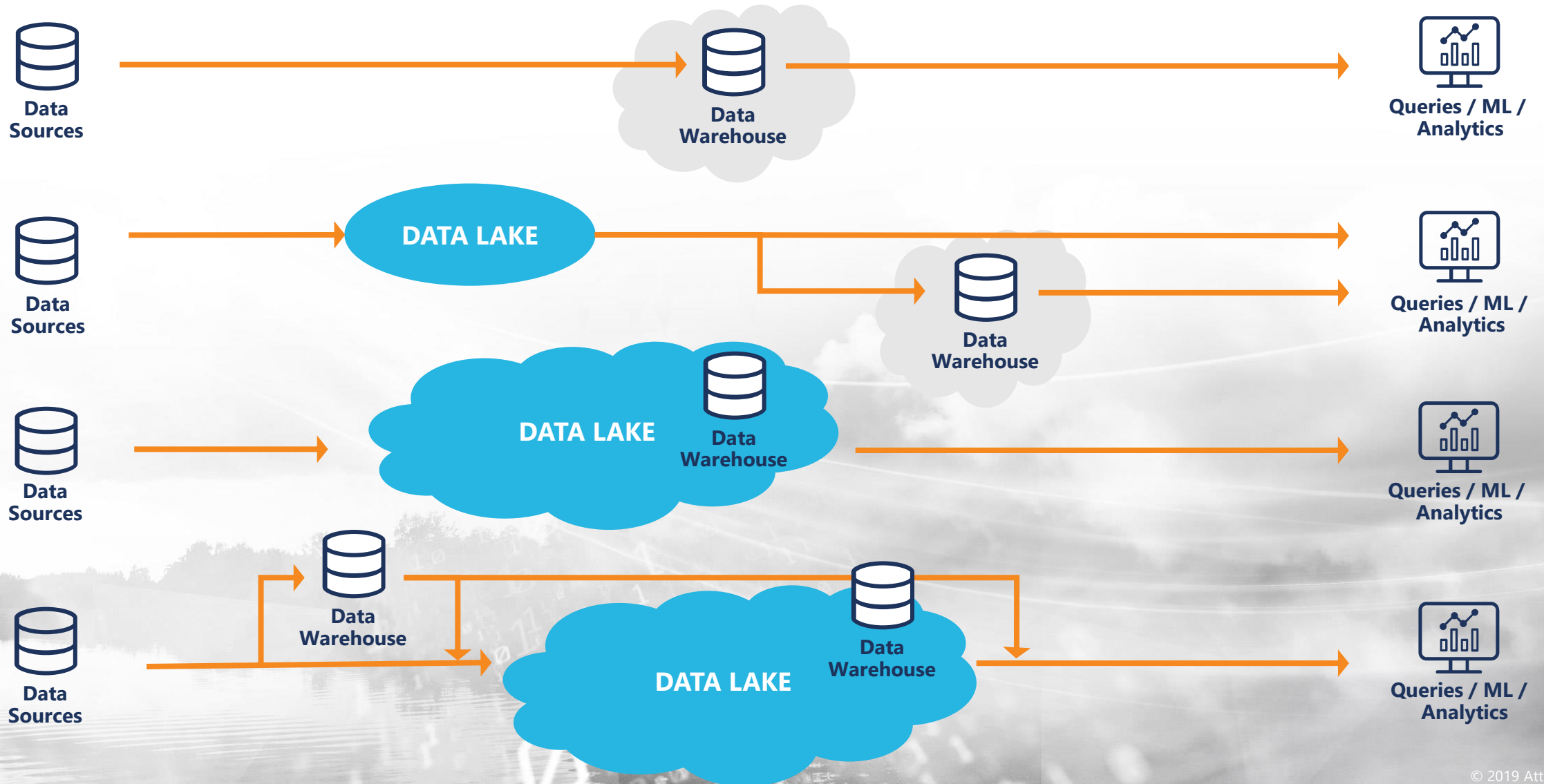
**DATA LAKE AND  
WAREHOUSE  
AUTOMATION**



**FUTURE-PROOF  
FOR  
'ARCHITECTURES  
IN MOTION'**

# DATAOPS REQUIREMENT #3

## AUTOMATION



MODERN  
**INTEGRATION**

**1** Continuous

**2** Universal

**3** Automation

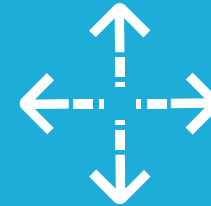
**4** Agility

**5** Trust

## Metadata creates greater trust and confidence from data consumers



**DATA CATALOG  
TO ACCESS IT-  
MANAGED AND  
USER-CURATED  
DATA SETS**



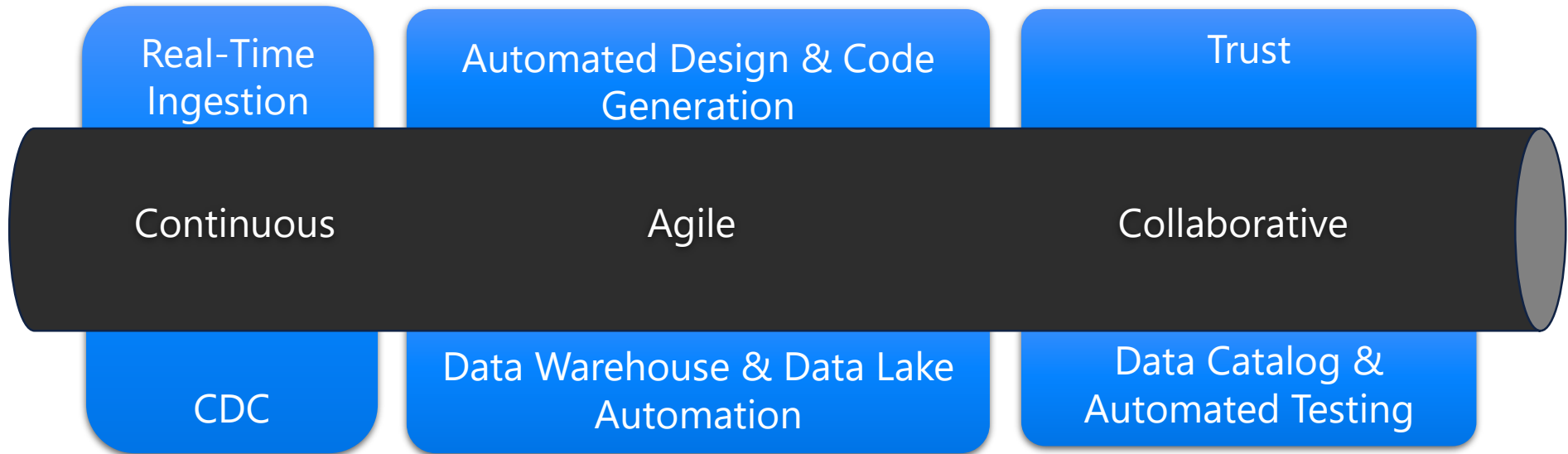
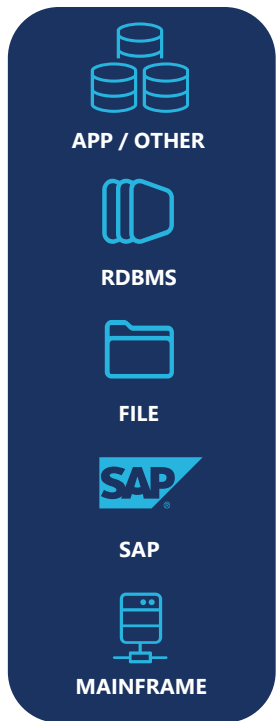
**DATA LINEAGE  
TO UNDERSTAND  
WHERE THE DATA  
CAME FROM AND  
HOW IT WAS  
TRANSFORMED**



**DATA VALIDATION  
TO ENSURE THAT  
ALL OF THE SOURCE  
DATA WAS  
REPLICATED**

# Technologies to support DataOps Data Pipelines

## ENTERPRISE DATA SOURCE





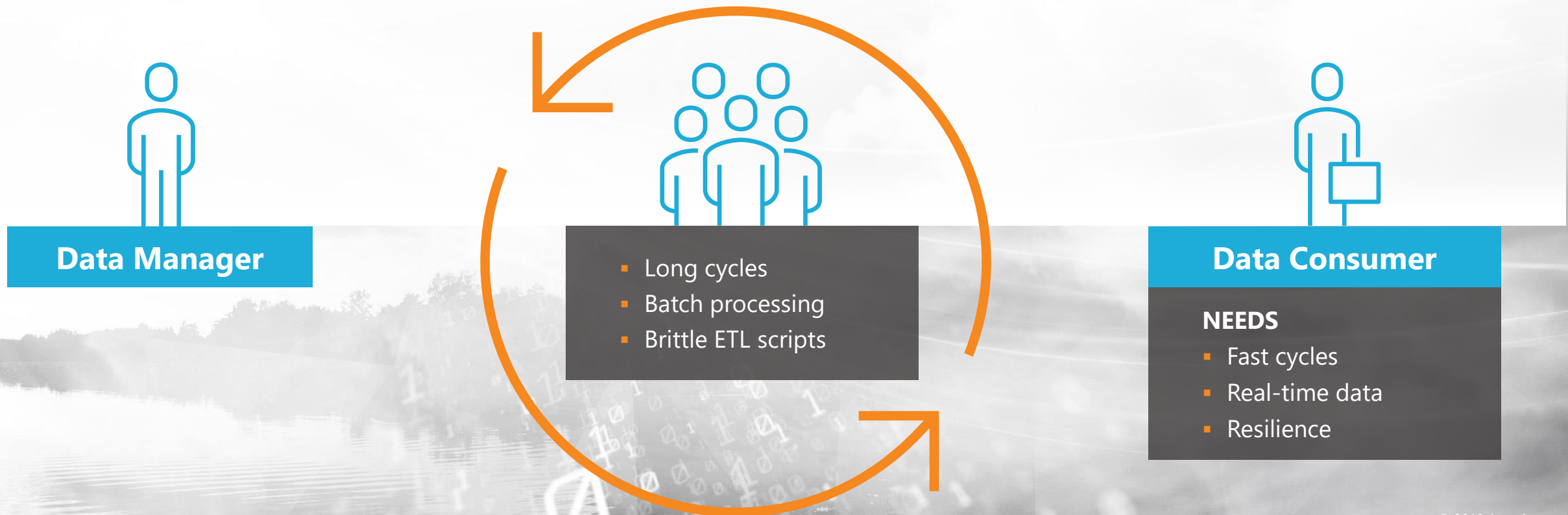
# Demonstration

DataOps pipeline imagined



# DATAOPS: ACCELERATING TIME TO INSIGHT

## TRADITIONAL



# DATAOPS: ACCELERATING TIME TO INSIGHT



**Data Manager**



- ✓ Code Automation
- ✓ Real-time delivery
- ✓ Run-time evolution



**Data Consumer**

**NEEDS**

- Fast cycles
- Real-time data
- Resilience

## DATAOPS



# Questions?



**Thank you...**

