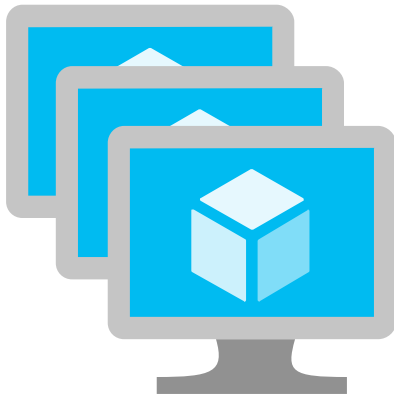# GitHub

https://github.com/mrpaulandrew

**CommunityEvents**

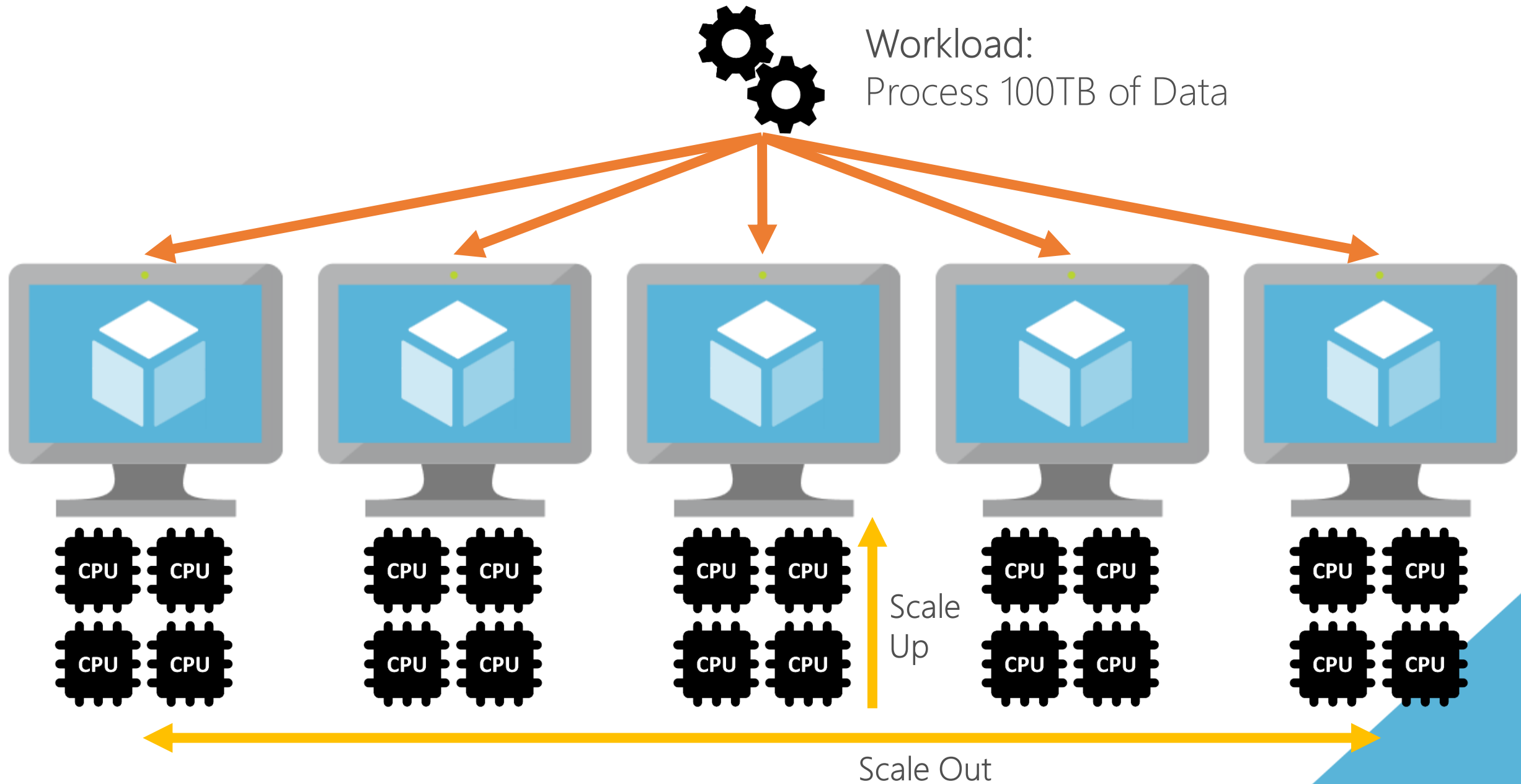Demo code, content and slides from various community events.

● C++

{Event/Location}-{Month}-{Year}

Scale Up vs Scale Out
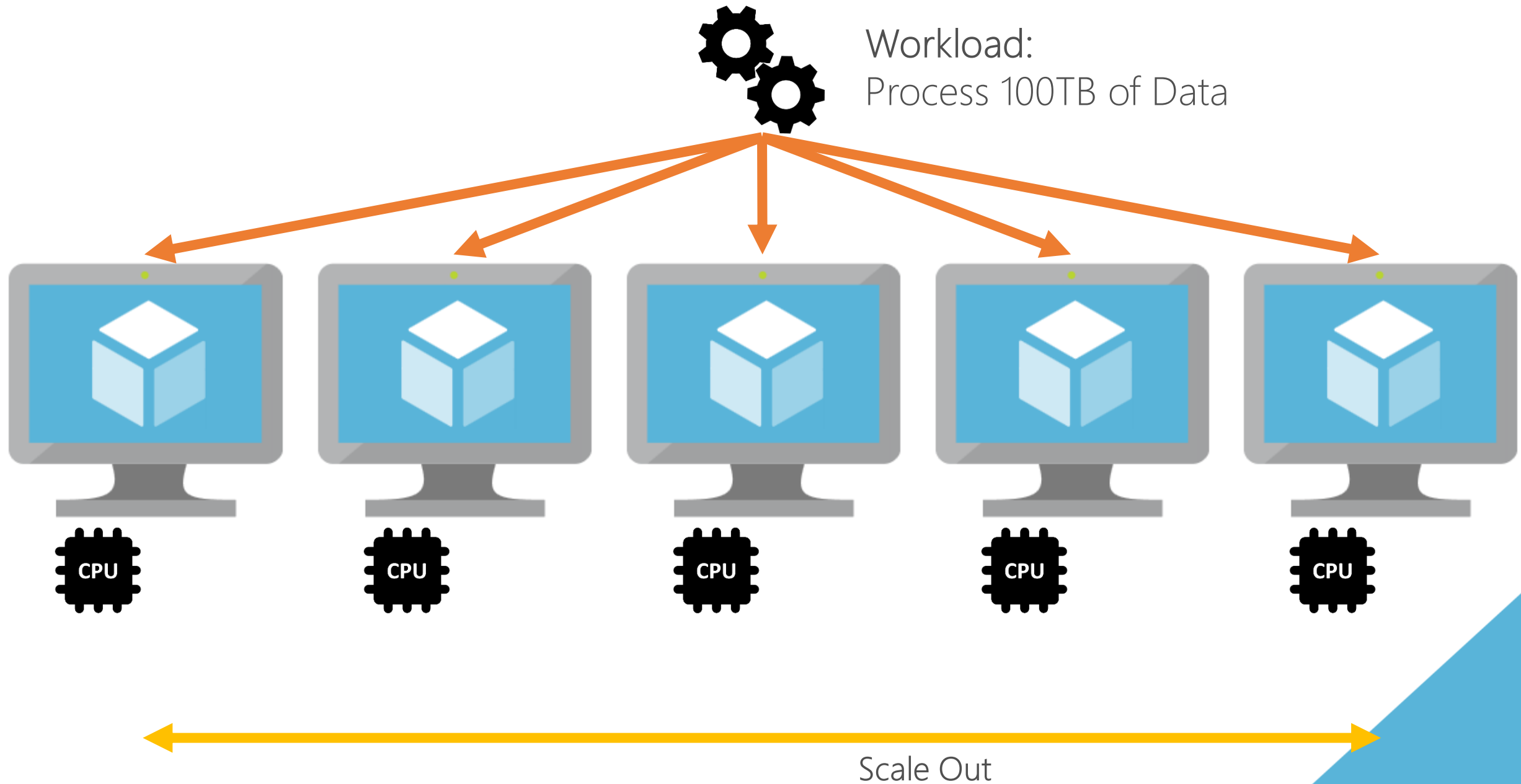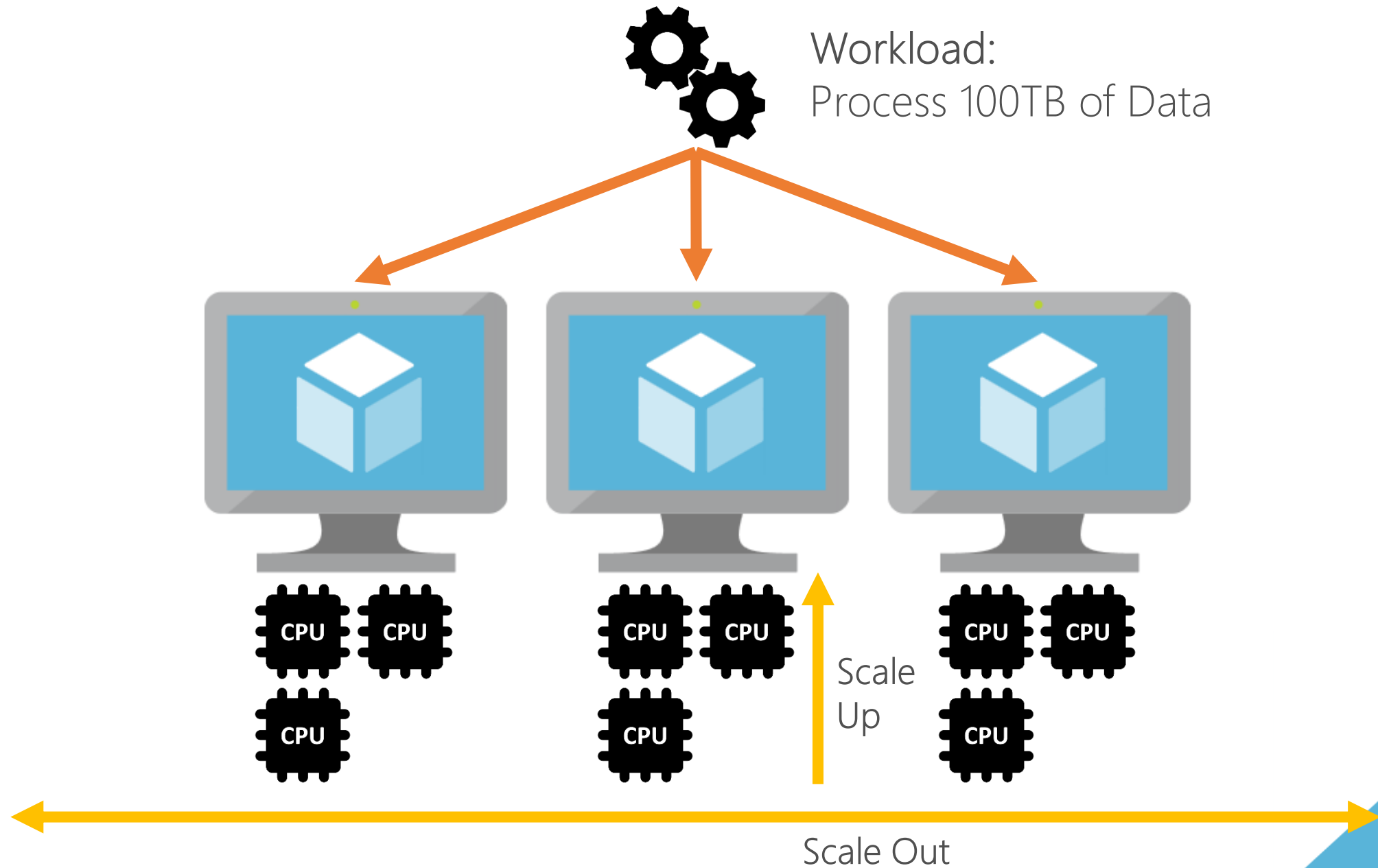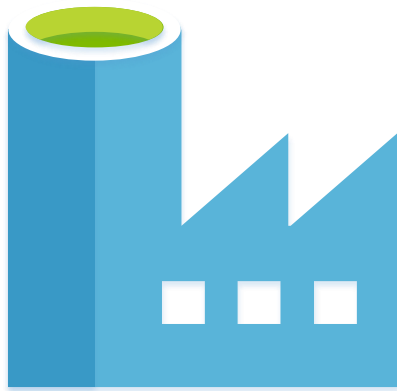
# Scale Up and Scale Out



Workload:
Process 100TB of Data

Scale Up

Scale Out

# Scale Up and Scale Out



Workload:
Process 100TB of Data

CPU     CPU     CPU     CPU     CPU

Scale Out

# Scale Up and Scale Out

# Azure Data Factory

# What is Azure Data Factory?

# What is Azure Data Factory?



Copy Data

Transform

# What is Azure Data Factory?



Copy Data                                        Transform

# Data Factory Components

Copy Data

Transform

**1** **Linked Services** – How and what to connect to. Like the SSIS connection manager.

SQLDBLinkedService

ConnectionString: *Server=MyServer;Database=myDataBase*
*UserName: "MrPaulAndrew"*
*Password: ***************

# Data Factory Components



**1** **Linked Services**

# Data Factory Components

Copy Data

Transform

**SQL**

**SQL**

**①** **Linked Services**

# Data Factory Components



Copy Data

Transform

**1** **Linked Services**

**2** **Data Sets** – Where is my data? What format? What file path/table do I need?

dbo.DimCustomer

/RAW/Orders/2018/01/01/Orders.csv

# Data Factory Components

Copy Data          Transform



**1** **Linked Services**

**2** **Data Sets**

**3** **Activities** – What do we want to happen?
With what conditions?

**Databricks Notebook Activity**

notebookPath: */Playground/Playing*
baseParameters: *Testing*
libraries[jar]: dbfs:/lib1.jar
linkedServiceName: *BricksOfData01*

# Data Factory Components

Copy Data

Transform

**1** **Linked Services**

**2** **Data Sets**

**3** **Activities**

**4** **Pipelines** – What groups of work do I want to do?

Sequence Container

Execute Package Task

Execute Pipeline Activity

# Data Factory Components

Copy Data

Transform

**1** **Linked Services**

**2** **Data Sets**

**3** **Activities**

**4** **Pipelines** – What groups of work do I want to do?

# Data Factory Components

Copy Data

Transform

**SQL**

**SQL**

| 1 | **Linked Services** |
| 2 | **Data Sets** |
| 3 | **Activities** |
| 4 | **Pipelines** |
| 5 | **Triggers** – How are we going to tell our pipeline(s) to execute? |

- Manual via UI
- Tumbling Windows
- Scheduled
- Blob File Events
- Logic App Calls

# Data Factory Components



Copy Data

Transform

**SQL**

SQL

1 **Linked Services**

2 **Data Sets**

3 **Activities**

4 **Pipelines**

5 **Triggers**

- **Manual**
- Tumbling Windows
- Scheduled
- Blob File Events
- Logic App Calls

```
Invoke-AzureRmDataFactoryV2Pipeline
-DataFactoryName $dataFactoryName
-ResourceGroupName $resourceGroupName
-PipelineName $pipelineName
```

PASS

# Data Factory Components

Copy Data

Transform

**SQL**

**SQL**

**1** **Linked Services**

**2** **Data Sets**

**3** **Activities**

**4** **Pipelines**

**5** **Triggers**

- Manual via UI
- **Tumbling Windows** - AKA Time Slices
- Scheduled
- Blob File Events
- Logic App Calls

Loading

2019

2020

# Data Factory Components



Copy Data

Transform

**SQL**

1 **Linked Services**

2 **Data Sets**

3 **Activities**

4 **Pipelines**

5 **Triggers**

- Manual via UI
- Tumbling Windows
- **Scheduled**
- Blob File Events
- Logic App Calls

- Every 1 minute.
- UTC

# Data Factory Components



Copy Data

Transform

1  **Linked Services**

2  **Data Sets**

3  **Activities**

4  **Pipelines**

5  **Triggers**

- Manual via UI
- Tumbling Windows
- Scheduled
- **Blob File Events**
- Logic App Calls

{Path} Created
{Path} Deleted

# Data Factory Components



Copy Data

Transform

**SQL**

SQL

1. **Linked Services**

2. **Data Sets**

3. **Activities**

4. **Pipelines**

5. **Triggers**

- Manual via UI
- Tumbling Windows
- Scheduled
- Blob File Events
- **Logic App Calls**

Azure Data Factory

← Search all actions

Triggers (0)    Actions (3)

Azure Data Factory - Cancel a pipeline run Preview

Azure Data Factory - Create a pipeline run Preview

Azure Data Factory - Get a pipeline run Preview

# Data Factory Components



Copy Data          Transform

1. **Linked Services**
2. **Data Sets**
3. **Activities**
4. **Pipelines**
5. **Triggers**

# Data Factory Control Flow Components



Copy Data

Transform

SQL

SQL

1 Linked Services

2 Data Sets

3 Activities

4 Pipelines

5 Triggers

# Integration Runtimes

**1** **Azure** Integration Runtime

Movement Hours

Activity Orchestration

Flexible Region

**2** **SSIS** Integration Runtime

SSIS Package Execution

Specified Region

**3** **Self Hosted** Integration Runtime

Gateway Access

Activity Orchestration

Virtual Machine

Why use Azure Data Factory?

Security

Privacy

Predict

Deploy

Land

Raw

TXT
TXT

CSV

CSV
CSV

Ingest

XLS

ZIP

JSON
JSON

Clean

PARQUET PARQUET PARQUET
PARQUET PARQUET PARQUET

PARQUET

Align

PARQUET PARQUET PARQUET
PARQUET PARQUET PARQUET

PARQUET

Conform

CSV
CSV

PARQUET
PARQUET

Serve

SQL

SQL

Transform

Monitor

Azure PaaS
Data Warehouse &
Analytics Platform

Why use Azure Data Factory?

Security | Privacy | Predict | Deploy

Land | Raw | Clean | Align | Conform | Serve

Ingest

Transform

Monitor

Orchestrate

Azure PaaS
Data Warehouse &
Analytics Platform

PASS

# Data Factory What & Why - Recap

**1** Linked Services

**2** Data Sets

**3** Activities

**4** Pipelines

**5** Triggers

**1** **Azure**
Integration Runtime

**2** **SSIS**
Integration Runtime

**3** **Self Hosted**
Integration Runtime

# Data Factory Control Flow Components

# Data Transformation in Azure

Data Transformation in Azure with SSIS

# Data Transformation in Azure with SSIS



SSIS Package

8x Packages per Node

SSIS Integration Runtime

10x Nodes

Express Route

VNet

Visual Studio – SQL Server Data Tools (SSDT)

- Sales Order Header
- Sales Order Details
- Merge Join
- Aggregate
- Order Line Count Table

Data Factory

Logical or Managed Instance

SSISDB

# Data Transformation in Azure with SSIS



SSIS Package

SSIS Integration Runtime

**D64_v3**
64 Cores
256GB RAM
£1,700 a month

1x Package per Node

1x Node

Sales Order Header

Sales Order Details

Merge Join

Aggregate

Order Line Count Table

Visual Studio – SQL Server Data Tools (SSDT)

Data Factory

Logical or Managed Instance

SSISDB

PASS

# Data Transformation in Azure

SSIS Integration
Runtime

# Data Transformation in Azure

# Data Transformation in Azure

Transform

# Data Factory Data Flows

# Data Factory Control Flow Components

# Data Factory Data Flows



Wrangling
Data Flow

Mapping
Data Flow

1 Linked Services
2 Data Sets
3 **Activities**
4 Pipelines
5 Triggers

PASS

# What is a Mapping Data Flow?



Data Factory Azure Portal Interface

Mapping Data Flow

Azure Databricks

Data Factory

# Mapping Data Flows

# Mapping Data Flows – Settings & Concepts



New Control Flow Activity

Azure Databricks

Azure SQLDW

PolyBase

# Integration Runtimes



**Azure** Integration Runtime — 1
- Movement Hours
- Activity Orchestration
- Flexible Region

**SSIS** Integration Runtime — 2
- SSIS Package Execution
- Specified Region

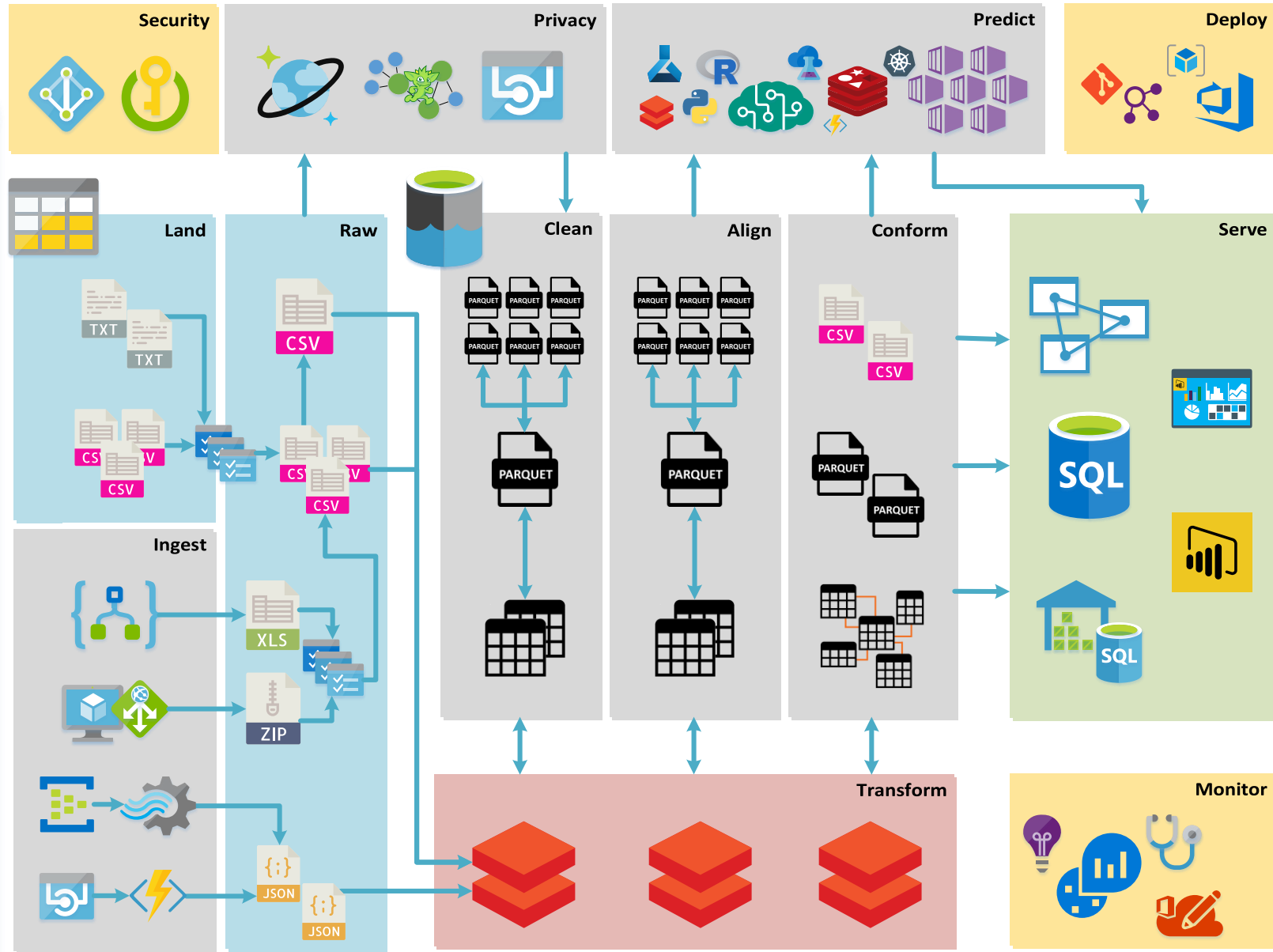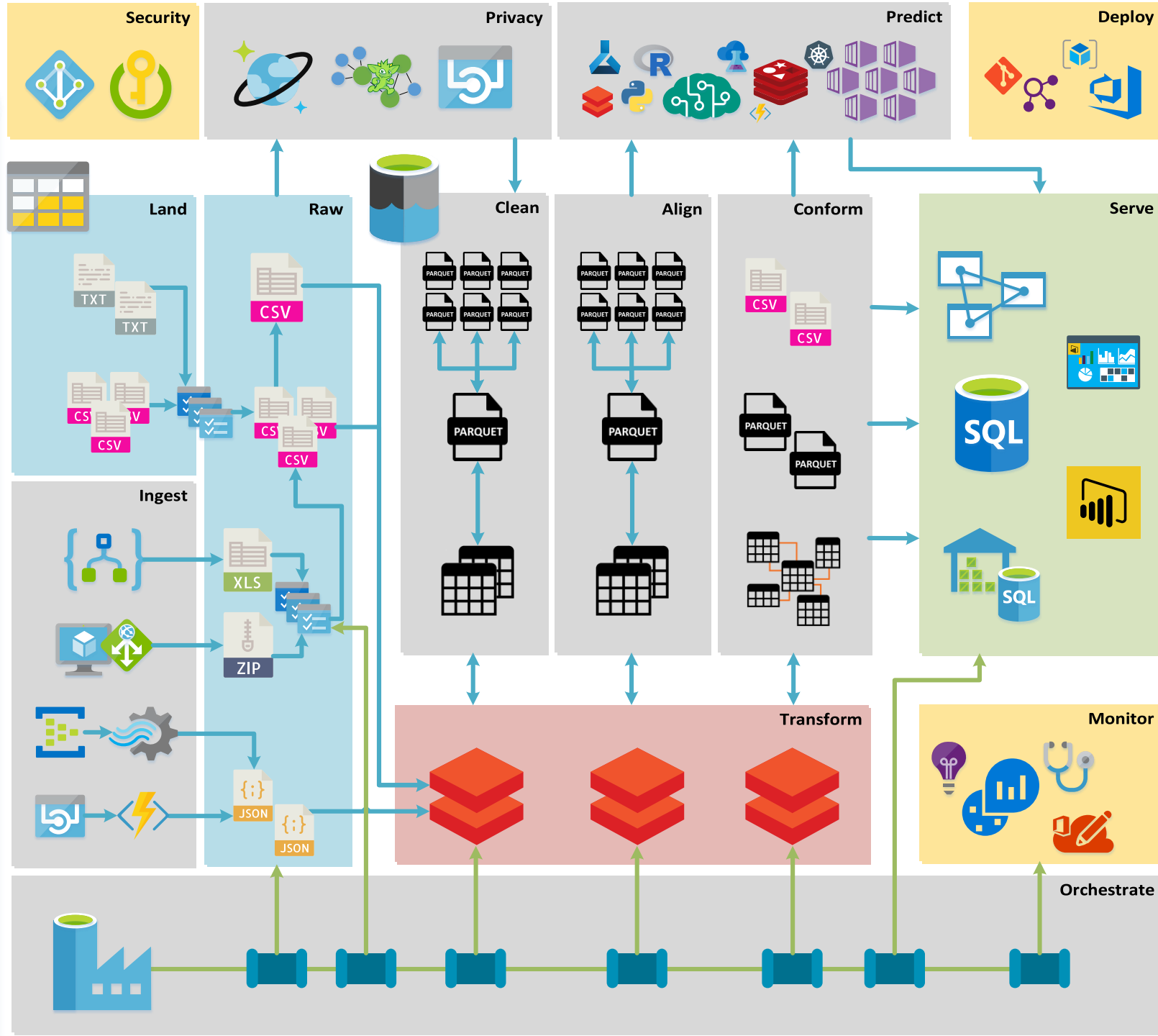**Self Hosted** Integration Runtime — 3
- Gateway Access
- Activity Orchestration
- Virtual Machine

# Integration Runtimes – Mapping Data Flow Cluster

# Mapping Data Flows – Settings & Concepts



New Control Flow Activity

Azure
Databricks

# Mapping Data Flows – Settings & Concepts

source1

Add source dataset

+

sink1

Add sink dataset

Currently Available:

Azure Blob Storage

Azure Data Lake Storage Gen1

Azure Data Lake Storage Gen2

Azure SQL Data Warehouse

Azure SQL Database

# Mapping Data Flows – Settings & Concepts

# Mapping Data Flows – Settings & Concepts

# Mapping Data Flows – Transformations

**Multiple inputs/outputs**

- ⬚ New Branch —————— Multicast
- ⬚ Join —————— Merge Join
- ⬚ Conditional Split
- ⬚ Union
- ⬚ Lookup

**Schema modifier**

- ⬚ Derived Column
- Σ Aggregate
- ⬚ Surrogate Key
- ⬚ Pivot
- ⬚ Unpivot
- ⬚ Window

https://docs.microsoft.com/en-gb/azure/data-factory/data-flow-surrogate-key

**Row modifier**

- ⬚ Exists
- ⬚ Select
- ⬚ Filter
- ⬚ Sort

**Custom**

- ⬚ Extend —————— Script Component

PASS

# Mapping Data Flows – Expression Builder

# Mapping Data Flows – Debug Mode



Only gives you a General Purpose cluster

# Mapping Data Flows – Monitoring

# Mapping Data Flows

**1** **Activity**
https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-overview

**2** **Source & Sink**
https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-schema-drift

**3** **Transformations**
https://docs.microsoft.com/en-gb/azure/data-factory/data-flow-aggregate

**4** **Expression Builder**
https://docs.microsoft.com/en-gb/azure/data-factory/data-flow-expression-functions

**5** **Debug Mode**
https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-debug-mode

**6** **Monitoring**
https://docs.microsoft.com/en-gb/azure/data-factory/concepts-data-flow-monitoring

**Mark Kromer**
https://github.com/kromerm/adfdataflowdocs

PASS

# What is a Wrangling Data Flow?



Data Factory Azure Portal Interface

Wrangling Data Flow

Azure Databricks

Data Factory

# What is a Wrangling Data Flow?



Power BI Desktop

Data Factory

# What is a Wrangling Data Flow?

Note to self - start clusters!

# Demo Summary

| Transformation Method | Graphical UI | Scales Out | Scales Up | Cloud Native Tech |
|---|---|---|---|---|
| T-SQL (SQLDB) | ✖ | ✖ | ✔ | ✖ |
| SSIS | ✔ | ✖ | ✔ | ✖ |
| Scala (Databricks) | ✖ | ✔ | ✔ | ✔ |
| Mapping Data Flow | ✔ | ✔ | ✔ | ✔ |

# Design Patterns

# Mapping Data Flow Future Design Patterns ???



Copy

Mapping Data Flow

Databricks

Source Systems

Data Lake Storage

Data Factory

**On Premises** **Azure**

# Mapping Data Flow Future Design Patterns ???

# Mapping Data Flow Future Design Patterns ???



Mapping Data Flow

**OrderHeader** — Import data from ADWSalesOrderHeader

**Join1** — Inner join on OrderHeader and OrderDetails

**Aggregate1** — Aggregating data by 'SalesOrderNumber' producing columns 'DetailLineCount'

**sink1** — Export data to ADWOrderLineCountTable

**OrderDetails** — Import data from ADWSalesOrderDetail

```
"fileName": {
    "value": "@dataset().FileName",
    "type": "Expression"
},
"folderPath": {
    "value": "@dataset().SourceDIR",
    "type": "Expression"
}
```

```
"transformations": [
    {
        "name": "Join1",
        "script": "OrderHeader, OrderDetail join(OrderHeader@SalesOrderID == OrderDetail@SalesOrderID,\n\tjoinType:'inner',\n\tbroadcast: 'none')~> Join1"
    },
    {
        "name": "Aggregate1",
        "script": "Join1 aggregate(groupBy(SalesOrderNumber),\n\tDetailLineCount = count(SalesOrderDetailID)) ~> Aggregate1"
    }
]
```

# What else?

PASS

**Wrangling Data Flow**

**Mapping Data Flow**

- Exploring
- Experimenting
- Analysing

- Productionising
- Engineering
- Warehousing

Future Design Patterns ???

# Conclusions

Why use Azure Data Factory?

Security
Privacy
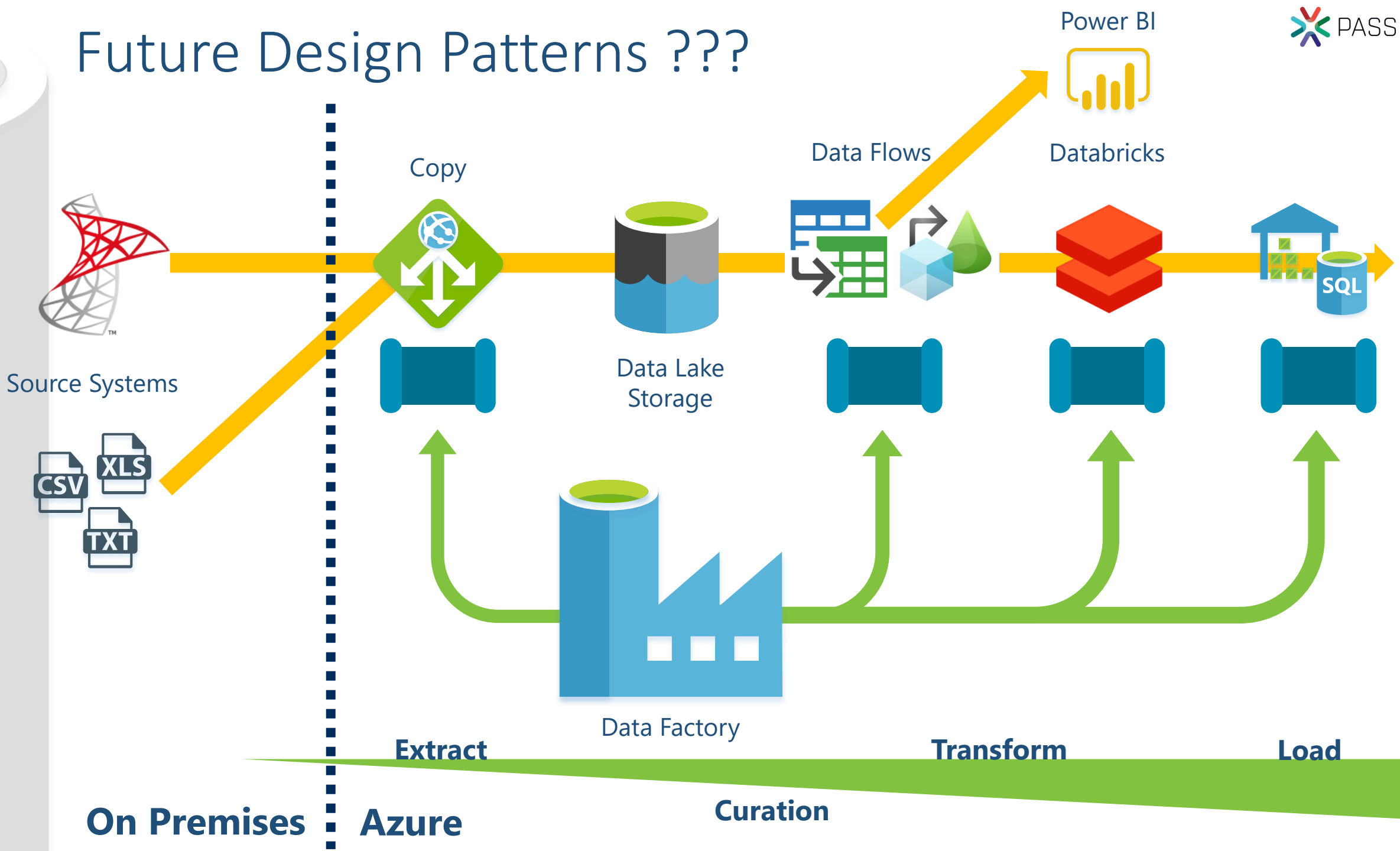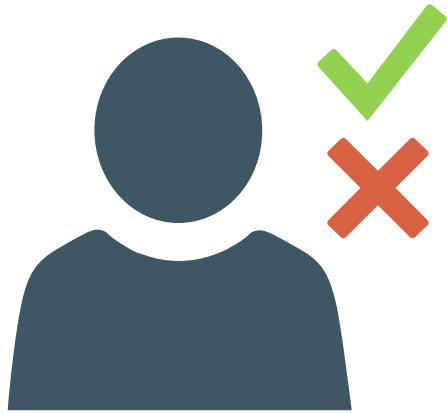Predict
Deploy

Land
Raw
Clean
Align
Conform
Serve

Ingest

Transform

Monitor

Orchestrate

Azure PaaS
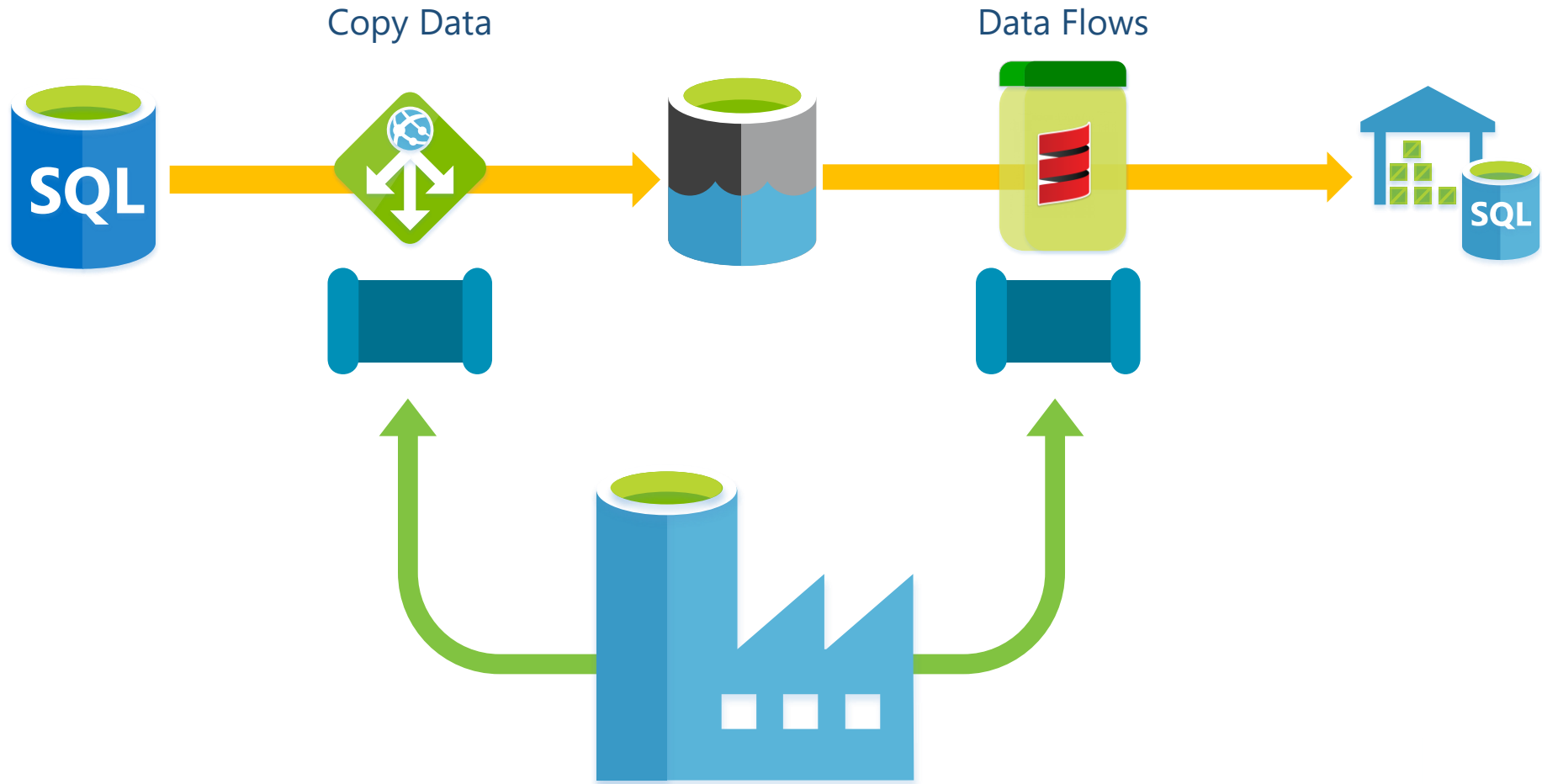Data Warehouse &
Analytics Platform

PASS

# Data Factory for the Data Scientist

# What is Azure Data Factory?

Copy Data

Data Flows

Orchestrator of our solution <u>Control Flow</u> operations.
Orchestrator of our solution <u>Data Flow</u> transformations.

… using cloud native technology in Azure and now with a user interface for both.

# Thanks for Listening
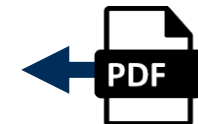
## Paul Andrew

🐦 @MrPaulAndrew

**Microsoft MVP Most Valuable Professional**     altius

**Email:**     paul@mrpaulandrew.com
**Blog:**      mrpaulandrew.com
**GitHub:**   github.com/mrpaulandrew ← PDF Slides