# Module 3

## Data Transformation

```sql
BEGIN --get ready

    SELECT
        [Contents]
    FROM
        [Training]
    WHERE
        [Module] = '3';
```

- Data Flows
- Power Query Injection
- Spark Configuration
- Use Cases

**Workload:**

Process 100TB of Data

Scale
Out

**Workload:**

Process 100TB of Data

CPU CPU CPU CPU CPU CPU
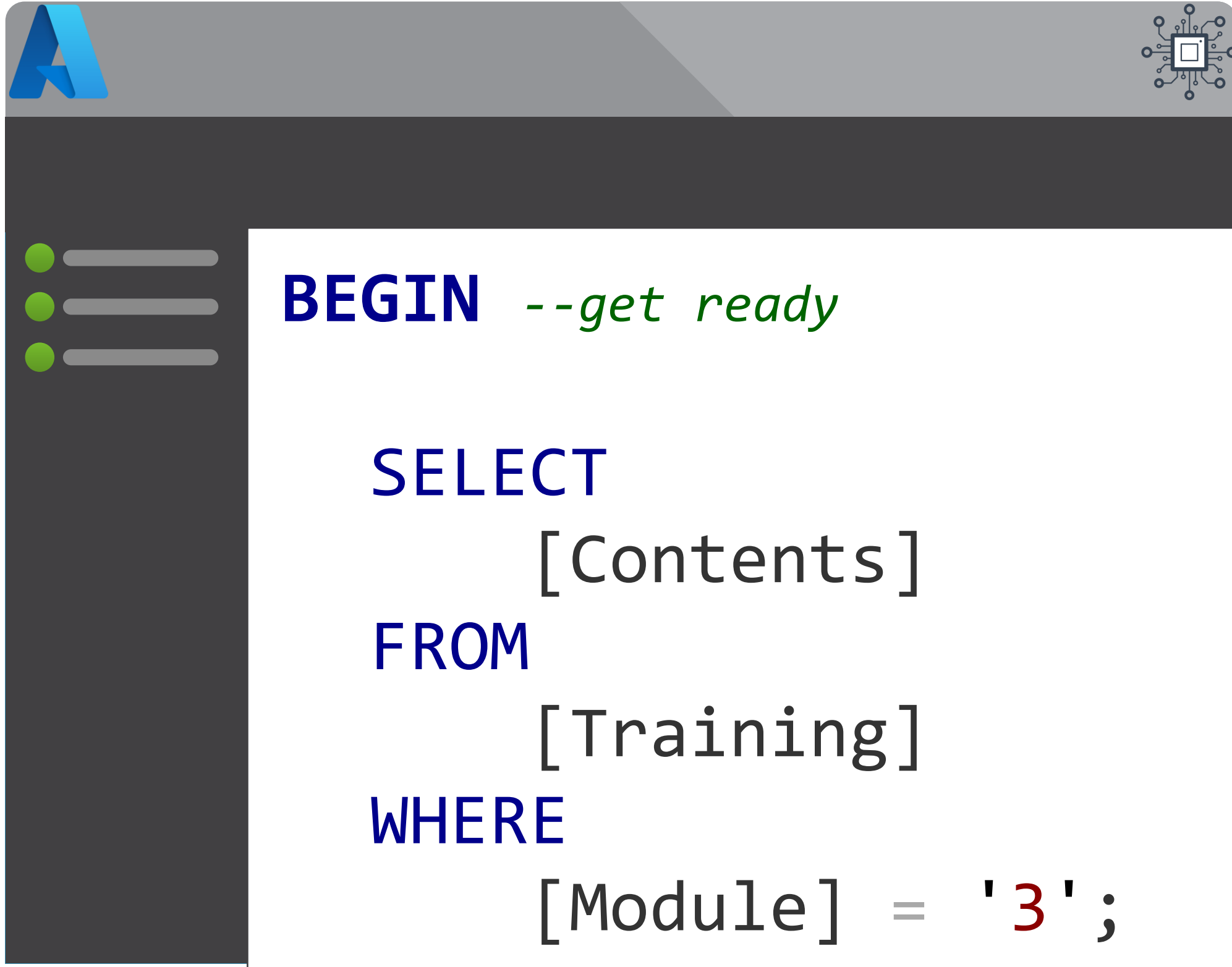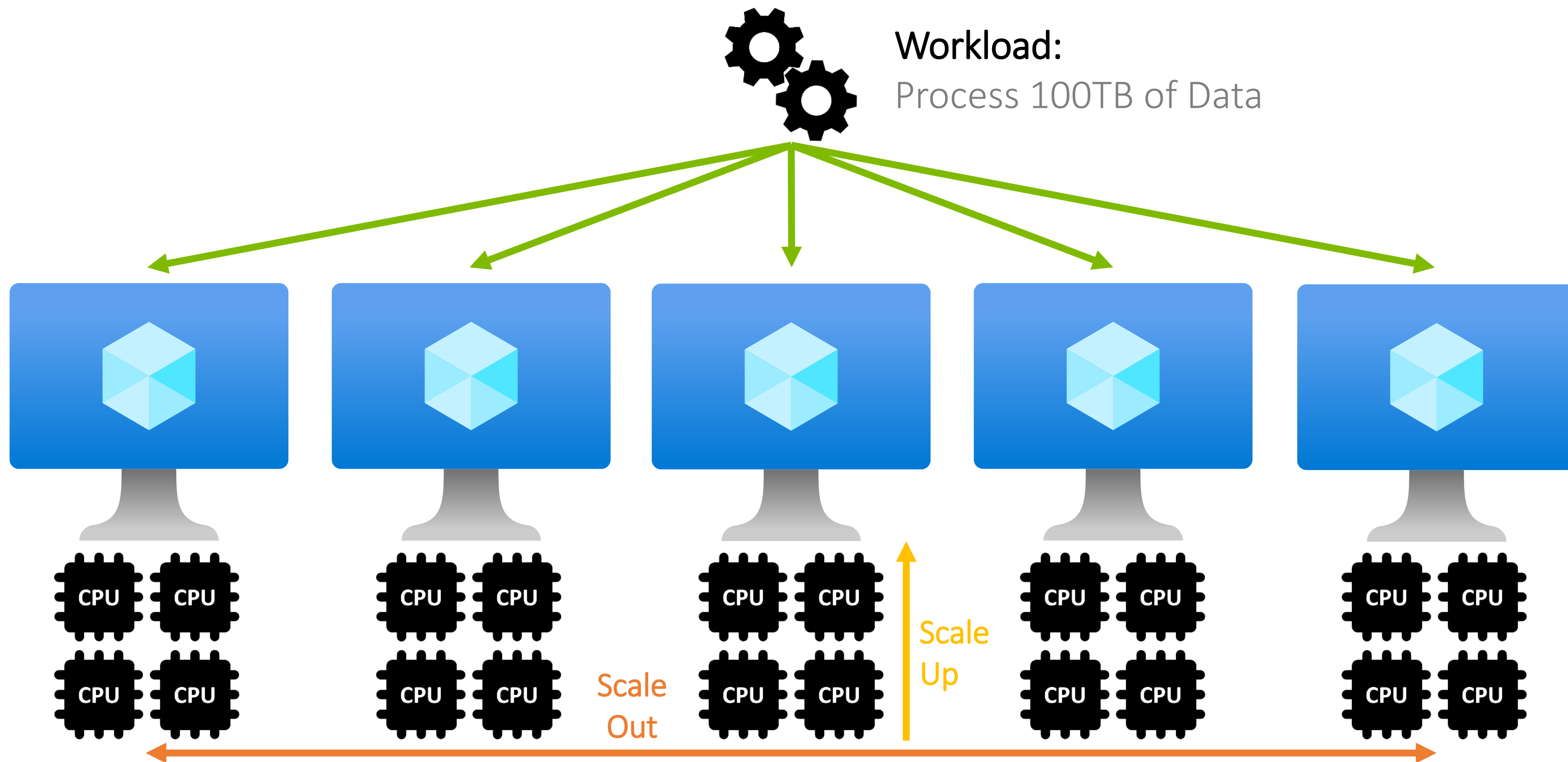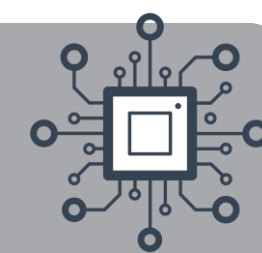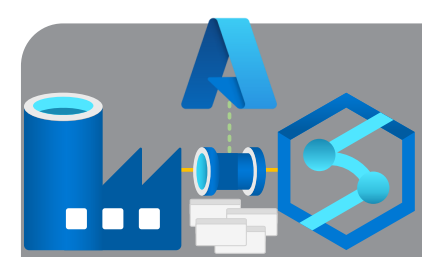
CPU CPU CPU

Scale Up

Scale Out

# Module 3

Data Transformation



- Data Flows
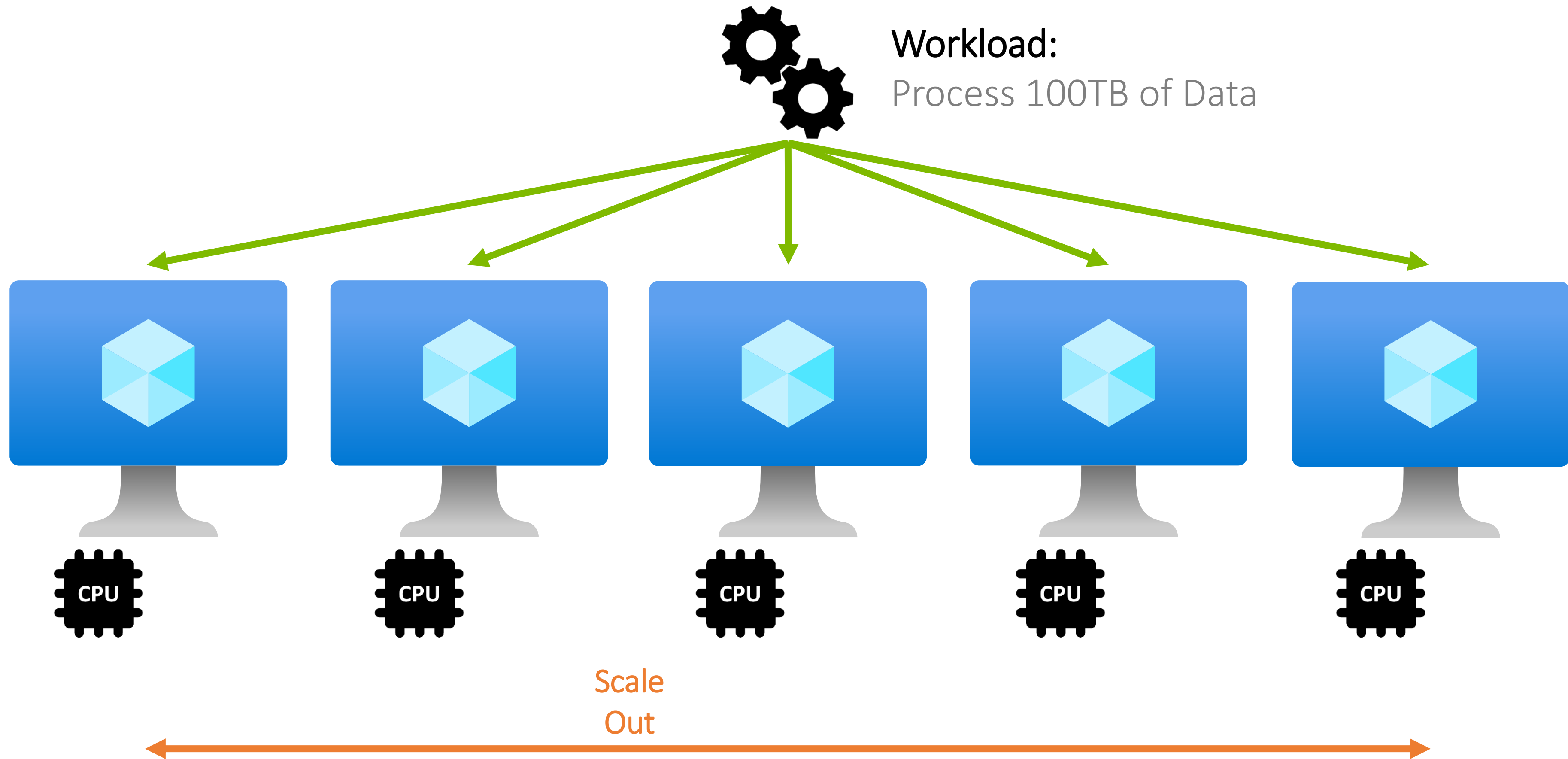- Power Query Injection
- Spark Configuration
- Use Cases
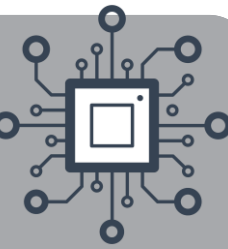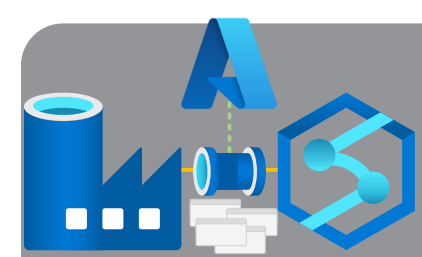
# Terminology Clarification

Data Flows

Mapping → <u>Data Flows</u>

Wrangling → Power Query

Integration Pipelines

# Control Flow Components



1. Linked Services
2. Datasets
3. Activities
4. Pipelines
5. Triggers

# ~~Control~~ Data Flow Components

CSV
XML
JSON
ZIP

JSON

SQL

JSON

1 Linked Services

2 Datasets

3 Activities

4 Pipelines

5 Triggers

Mapping → Data Flows

Wrangling → Power Query

# Other Data Transformation Services in Azure



| SSIS Packages | HD Insight | Data Lake Analytics | Synapse SQL Pools | SQL Database | Batch Service | Durable Functions | Synapse Spark | Databricks Spark | Analysis Services | Cosmos DB |

# Other Data Transformation Services in Azure

## When Should We Use These Integration Pipeline Transformation Activities?

| SSIS Packages | HD Insight | Data Lake Analytics | Synapse SQL Pools | SQL Database | Batch Service | Durable Functions | Synapse Spark | Databricks Spark | Analysis Services | Cosmos DB |
|---|---|---|---|---|---|---|---|---|---|---|

Data Flow

Power Query

# What is a ~~Mapping~~ Data Flow?

**Control Flow**

Data Flow
Transform Stuff
JSON

**Data Flow**



| OrderHeader | JoinHeaderToLineDetails | OrderLineCount | OrderSummary |
|---|---|---|---|
| Import data from LakeFileOrderHeaderParquet | Inner join on OrderHeader and OrderLineDetails | Aggregating data by 'SalesOrderNumber' producing columns 'RecordCount' | Export data to TableOrderSummary |

OrderLineDetails
Import data from LakeFileOrderDetailLinesParquet

Script

Spark

# Q: What is a ~~Mapping~~ Data Flow?

**Control Flow**

Data Flow — Transform Stuff — JSON

**Data Flow**

OrderHeader
Import data from LakeFileOrderHeaderParquet

JoinHeaderToLineDetails
Inner join on OrderHeader and OrderLineDetails

OrderLineCount
Aggregating data by 'SalesOrderNumber' producing columns 'RecordCount'

OrderSummary
Export data to TableOrderSummary

OrderLineDetails
Import data from LakeFileOrderDetailLinesParquet

Script — Spark

**A:** Graphic no low/low code data transformation tool that sits on top of Apache Spark.

# Data Flows – Inputs & Outputs

| | |
|---|---|
| Source & Sink | |
| Linked Services | |
| Source Types — Dataset | |
| Source Types — Inline | |

# Data Flows – Transformations

New Branch

Join

Conditional Split

Exists

Union

Lookup

Derived Column

Select

Aggregate

Surrogate Key

Pivot/Unpivot

Window

Rank

External Call

Cast

Flatten

Parse

Stringify

Filter

Sort

Alter Row

Assert

Flowlet

Key
Input & Output Modifiers
Schema Modifiers
Formatters
Row Modifiers

# Data Flows – Transformations

New Branch

Join

Conditional Split

Exists

Union

Lookup

Filter

Sort

Alter Row

## Components

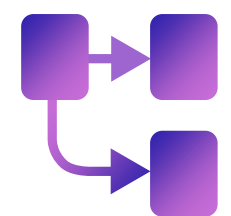| Operation / Activity | Description | SSIS equivalent | SQL Server equivalent |
|---|---|---|---|
| New branch | Create a new flow branch with the same data | Multicast (+icon) | `SELECT INTO` `SELECT OUTPUT` |
| Join | Join data from two streams based on a condition | Merge join | `INNER/LEFT/RIGHT JOIN,` `CROSS/FULL OUTER JOIN` |
| Conditional Split | Route data into different streams based on conditions | Conditional Split | `SELECT INTO WHERE condition1` `SELECT INTO WHERE condition2` `CASE ... WHEN` |
| Union | Collect data from multiple streams | Union All | `SELECT col1a UNION (ALL)` `SELECT col1b` |
| Lookup | Lookup additional data from another stream | Lookup | *Subselect, function,* `LEFT/RIGHT JOIN` |
| Derived Column | Compute new columns based on the existing once | Derived Column | `SELECT Column1 * 1.09 as NewColumn` |
| Aggregate | Calculate aggregation on the stream | Aggregate | `SELECT Year(DateOfBirth) as YearOnly,` `MIN(), MAX(), AVG()` `GROUP BY Year(DateOfBirth)` |
| Surrogate Key | Add a surrogate key column to output stream from a specific value | Script Component | `SELECT ROW_NUMBER()` `   OVER(ORDER BY name ASC) AS Row#,` `name` `FROM sys.databases` |

Key
Input & Output Modifiers
Schema Modifiers
Formatters
Row Modifiers

Staging

◢ PolyBase ⓘ

| | |
|---|---|
| Staging linked service | Select...  ∨  ⓘ  + New |
| Staging storage folder | Container  /  Directory  📁 Browse  \| ∨ |

# Data Flows – Expression Builder

Expression Builder

Transform Stuff

Count If

Sum If

Regex Extract

Regex Match

Regex Split

R-Like

# Data Flows – Data Distribution

Side Note

Broadcast Data

Optimize

Data Flow

Transform Stuff

Round Robin

Hashed

Conditional

Conditions

Dynamic Ranges

Key Based

Keys

Sample Data &
Row Limits

PARQUET

Data Preview

Data Flow

Transform Stuff

Spark

| SK | Field |
|----|-------|
| 1  | xxx   |
| 2  | yyy   |
| 3  | xxx   |
| 4  | xxx   |
| 5  | zzz   |

| SK | Field |
|----|-------|
| 1  | xxx   |
| 2  | yyy   |
| 3  | xxx   |
| 4  | xxx   |
| 5  | zzz   |

Enable Data Flow Debug Mode

Data Flow

Transform Stuff

# Data Flows – Monitoring

Transform Stuff

Data Flow

Spark

Mapping Order Aggregation
Data flow

Refresh    Auto refresh  On

Cluster startup time: **6s 878ms**  Number of transforms: **5**

| OrderHeader | JoinHeaderToLineD... | OrderLineCount | OrderSummary |
| Sink: ● | Sink: ● | Sink: ● | |

OrderLineDetails
Sink: ●

OrderHeader
Import data from
LakeFileOrderHeaderParquet

JoinHeaderToLineDetails
Inner join on OrderHeader and
OrderLineDetails

OrderLineCount
Aggregating data by
'SalesOrderNumber'
producing columns
'RecordCount'

OrderSummary
Export data to
TableOrderSummary

OrderLineDetails
Import data from
LakeFileOrderDetailLinesParque
t

## OrderLineCount
Aggregate

| | |
|---|---|
| Total columns | 2 |
| New columns | 1 |
| Updated columns | 1 |
| Dropped columns | 27 |
| Drifted columns | 0 |

### Stream information

| | |
|---|---|
| Rows calculated | 32 |
| Total partition | 190 |
| Stage time | 3s 405ms |
| Last update (GDT) | 25/08/2020, 14:44:44 |

### Partition chart

### OrderSummary                              ✓

**Processing time:** 7s 60ms

| TRANSFORM | ROWS | TIME |
|---|---|---|
| ● OrderHeader | - | - |
| ● JoinHeaderToLineDetails | 542 | |
| ● OrderLineDetails | 542 | 277ms |
| ● OrderSummary | 32 | |
| ● OrderLineCount | 32 | 3s 405ms |

| | |
|---|---|
| Skewness | 2.2425 |
| Kurtosis | 7.2667 |

Edit transformation

# Module 3

## Data Transformation

- Data Flows
- **Power Query Injection**
- Spark Configuration
- Use Cases

*(Public Preview)*

# What is a Data Flow?



Control Flow

Data Flow

Data Flow
Transform Stuff
JSON

OrderHeader
Import data from LakeFileOrderHeaderParquet

JoinHeaderToLineDetails
Inner join on OrderHeader and OrderLineDetails

OrderLineCount
Aggregating data by 'SalesOrderNumber' producing columns 'RecordCount'

OrderSummary
Export data to TableOrderSummary

OrderLineDetails
Import data from LakeFileOrderDetailLinesParquet

Script

Spark

# What is a Power Query Activity?

Control
Flow

Power Query

Transform Stuff

JSON

# What is a Power Query Activity?

# What can a Power Query Activity do?

# What can a Power Query Activity do?

Control Flow

Power Query — Transform Stuff

JSON

---

Power Query

## Power Query Editor Ribbon (ADF)

Home | Transform | Add column | View

Enter data | Options | Manage parameters | Properties | Advanced editor | Manage | Refresh | Choose columns | Remove columns | Keep rows | Remove rows | Data type: Whole number | Use first row as headers | Replace values | Split column | Group by | Merge queries | Append queries | Combine files

New query | Options | Parameters | Query | Manage columns | Reduce rows | Sort | Transform | Combine

Queries

▲ ADFResource [1]
  ▦ LakeFileOrderDetailL... 
  ▦ UserQuery

| | 1²₃ SalesOrderID | 1²₃ |
|---|---|---|
| 1 | 71774 | |
| 2 | 71774 | |
| 3 | 71776 | |
| 4 | 71780 | |
| 5 | 71780 | |
| 6 | 71780 | |
| 7 | 71780 | |
| 8 | 71780 | |
| 9 | 71780 | |
| 10 | 71780 | |
| 11 | 71780 | |
| 12 | 71780 | |
| 13 | 71780 | |
| 14 | 71780 | |
| 15 | 71780 | |
| 16 | 71780 | |
| 17 | 71780 | |

---

## Untitled - Power Query Editor

File | Home | Transform | Add Column | View | Tools | Help

Close & Apply | New Source | Recent Sources | Enter Data | Data source settings | Manage Parameters | Refresh Preview | Properties | Advanced Editor | Manage | Choose Columns | Remove Columns | Keep Rows | Remove Rows | Data Type: Whole Number | Use First Row as Headers | Replace Values | Split Column | Group By | Merge Queries | Append Queries | Combine Files | Text Analytics | Vision | Azure Machine Learning

Close | New Query | Data Sources | Parameters | Query | Manage Columns | Reduce Rows | Sort | Transform | Combine | AI Insights

Queries [1]

▦ OrderDetailLines

= Table.TransformColumnTypes(#"Promoted Headers",{{"SalesOrderID", Int64.Type}, {"SalesOrderDetailID", Int64.Type}

| | 1²₃ SalesOrderID | 1²₃ SalesOrderDetailID | 1²₃ OrderQty | 1²₃ ProductID | 1.2 UnitPrice | 1.2 UnitPriceⱼ |
|---|---|---|---|---|---|---|
| 1 | 71774 | 110562 | 1 | 836 | 356.898 | |
| 2 | 71774 | 110563 | 1 | 822 | 356.898 | |
| 3 | 71776 | 110567 | 1 | 907 | 63.9 | |
| 4 | 71780 | 110616 | 4 | 905 | 218.454 | |
| 5 | 71780 | 110617 | 2 | 983 | 461.694 | |
| 6 | 71780 | 110618 | 6 | 988 | 112.998 | |
| 7 | 71780 | 110619 | 2 | 748 | 818.7 | |
| 8 | 71780 | 110620 | 1 | 990 | 323.994 | |
| 9 | 71780 | 110621 | 1 | 926 | 149.874 | |
| 10 | 71780 | 110622 | 1 | 743 | 809.76 | |
| 11 | 71780 | 110623 | 4 | 782 | 1376.994 | |
| 12 | 71780 | 110624 | 2 | 918 | 158.43 | |
| 13 | 71780 | 110625 | 4 | 780 | 1391.994 | |
| 14 | 71780 | 110626 | 1 | 937 | 48.594 | |
| 15 | 71780 | 110627 | 6 | 867 | 41.994 | |
| 16 | 71780 | 110628 | 1 | 985 | 112.998 | |
| 17 | 71780 | 110629 | 2 | 989 | 323.994 | |

Query Settings

▲ PROPERTIES
Name
OrderDetailLines

All Properties

▲ APPLIED STEPS
Source
Promoted Headers
✕ Changed Type

# What can a Power Query Activity do?

## Transform



Control Flow

Power Query → Transform Stuff → JSON

Power Query

Home | Transform | Add column | View

Group by | Use first row as headers | Transpose | Reverse rows | Count rows | Replace values | Data type: Whole number | Detect data type | Mark as key | Rename | Pivot column | Unpivot columns | Fill | Move | Convert to list | Merge columns | ABC 123 Extract | abc Parse | Split column | Format | Statistics | Standard | Scientific | 10² | Trigonometry | Rounding | Information | Date | Time | Duration

Table

Queries
▲ ADFResource [1]
  LakeFileOrderDetailL...
UserQuery

| | 1²₃ SalesOrderID |
|---|---|
| 1 | 71774 |
| 2 | 71774 |
| 3 | 71776 |
| 4 | 71780 |
| 5 | 71780 |
| 6 | 71780 |
| 7 | 71780 |
| 8 | 71780 |
| 9 | 71780 |
| 10 | 71780 |
| 11 | 71780 |
| 12 | 71780 |
| 13 | 71780 |
| 14 | 71780 |
| 15 | 71780 |
| 16 | 71780 |
| 17 | 71780 |

Untitled - Power Query Editor

File | Home | Transform | Add Column | View | Tools | Help

Group By | Use First Row as Headers | Transpose | Reverse Rows | Count Rows | Data Type: Whole Number | Detect Data Type | Rename | Replace Values | Fill | Pivot Column | Unpivot Columns | Move | Convert to List | Merge Columns | ABC 123 Extract | Parse | Split Column | Format | Statistics | Standard | Scientific | 10² | Trigonometry | Rounding | Information | Date | Time | Duration | Structured Column | Run R script | Run Python script

Table | Any Column | Text Column | Number Column | Date & Time Column | Scripts

Queries [1]
OrderDetailLines

= Table.TransformColumnTypes(#"Promoted Headers",{{"SalesOrderID", Int64.Type}, {"SalesOrderDetailID", Int64.Type}

| | 1²₃ SalesOrderID | 1²₃ SalesOrderDetailID | 1²₃ OrderQty | 1²₃ ProductID | 1.2 UnitPrice | 1.2 UnitPriceD |
|---|---|---|---|---|---|---|
| 1 | 71774 | 110562 | 1 | 836 | 356.898 | |
| 2 | 71774 | 110563 | 1 | 822 | 356.898 | |
| 3 | 71776 | 110567 | 1 | 907 | 63.9 | |
| 4 | 71780 | 110616 | 4 | 905 | 218.454 | |
| 5 | 71780 | 110617 | 2 | 983 | 461.694 | |
| 6 | 71780 | 110618 | 6 | 988 | 112.998 | |
| 7 | 71780 | 110619 | 2 | 748 | 818.7 | |
| 8 | 71780 | 110620 | 1 | 990 | 323.994 | |
| 9 | 71780 | 110621 | 1 | 926 | 149.874 | |
| 10 | 71780 | 110622 | 1 | 743 | 809.76 | |
| 11 | 71780 | 110623 | 4 | 782 | 1376.994 | |
| 12 | 71780 | 110624 | 2 | 918 | 158.43 | |
| 13 | 71780 | 110625 | 4 | 780 | 1391.994 | |
| 14 | 71780 | 110626 | 1 | 937 | 48.594 | |
| 15 | 71780 | 110627 | 6 | 867 | 41.994 | |
| 16 | 71780 | 110628 | 1 | 985 | 112.998 | |
| 17 | 71780 | 110629 | 2 | 989 | 323.994 | |

Query Settings

▲ PROPERTIES
Name
OrderDetailLines
All Properties

▲ APPLIED STEPS
Source
Promoted Headers
✕ Changed Type

# What can a Power Query Activity do?

Control Flow

Power Query

Transform Stuff

JSON



Power Query

| Home | Transform | Add column | View |

Custom column · Conditional column · Index column · Duplicate column · General · Format · Merge columns · Extract · Parse · From text · Statistics · Standard · Scientific · Trigonometry · Rounding · Information · Date · Time · Duration

Queries
ADFResource [1]
LakeFileOrderDetailL...
UserQuery

SalesOrderID
71774
71774
71776
71780
71780
71780
71780
71780
71780
71780
71780
71780
71780
71780
71780
71780
71780

Untitled - Power Query Editor

| File | Home | Transform | Add Column | View | Tools | Help |

Column From Examples · Custom Column · Invoke Custom Function · Conditional Column · Index Column · Duplicate Column · General · Format · Merge Columns · Extract · Parse · From Text · Statistics · Standard · Scientific · Trigonometry · Rounding · Information · From Number · Date · Time · Duration · From Date & Time · Text Analytics · Vision · Azure Machine Learning · AI Insights

Queries [1]
OrderDetailLines

= Table.TransformColumnTypes(#"Promoted Headers",{{"SalesOrderID", Int64.Type}, {"SalesOrderDetailID", Int64.Type}

| | SalesOrderID | SalesOrderDetailID | OrderQty | ProductID | UnitPrice | UnitPrice[ |
|---|---|---|---|---|---|---|
| 1 | 71774 | 110562 | 1 | 836 | 356.898 | |
| 2 | 71774 | 110563 | 1 | 822 | 356.898 | |
| 3 | 71776 | 110567 | 1 | 907 | 63.9 | |
| 4 | 71780 | 110616 | 4 | 905 | 218.454 | |
| 5 | 71780 | 110617 | 2 | 983 | 461.694 | |
| 6 | 71780 | 110618 | 6 | 988 | 112.998 | |
| 7 | 71780 | 110619 | 2 | 748 | 818.7 | |
| 8 | 71780 | 110620 | 1 | 990 | 323.994 | |
| 9 | 71780 | 110621 | 1 | 926 | 149.874 | |
| 10 | 71780 | 110622 | 1 | 743 | 809.76 | |
| 11 | 71780 | 110623 | 4 | 782 | 1376.994 | |
| 12 | 71780 | 110624 | 2 | 918 | 158.43 | |
| 13 | 71780 | 110625 | 4 | 780 | 1391.994 | |
| 14 | 71780 | 110626 | 1 | 937 | 48.594 | |
| 15 | 71780 | 110627 | 6 | 867 | 41.994 | |
| 16 | 71780 | 110628 | 1 | 985 | 112.998 | |
| 17 | 71780 | 110629 | 2 | 989 | 323.994 | |

Query Settings

PROPERTIES
Name
OrderDetailLines

All Properties

APPLIED STEPS
Source
Promoted Headers
Changed Type

# What can a Power Query Activity do?

## Control Flow

Power Query

Transform Stuff

JSON

## Power Query

| Home | Transform | Add column | View |

Data view | Schema view | Go to column | Advanced editor

Preview | Columns | Advanced

Queries

▲ ADFResource [1]
  LakeFileOrderDetailL...
UserQuery

| 1²₃ SalesOrderID | | |
|---|---|---|
| 1 | 71774 | |
| 2 | 71774 | |
| 3 | 71776 | |
| 4 | 71780 | |
| 5 | 71780 | |
| 6 | 71780 | |
| 7 | 71780 | |
| 8 | 71780 | |
| 9 | 71780 | |
| 10 | 71780 | |
| 11 | 71780 | |
| 12 | 71780 | |
| 13 | 71780 | |
| 14 | 71780 | |
| 15 | 71780 | |
| 16 | 71780 | |
| 17 | 71780 | |

Untitled - Power Query Editor

| File | Home | Transform | Add Column | View | Tools | Help |

☑ Formula Bar
☐ Monospaced     ☐ Column distribution     ☐ Always allow
☑ Show whitespace   ☐ Column profile
☐ Column quality

Query Settings | Go to Column | Advanced Editor | Query Dependencies

Layout | Data Preview | Columns | Parameters | Advanced | Dependencies

Queries [1]

OrderDetailLines

= Table.TransformColumnTypes(#"Promoted Headers",{{"SalesOrderID", Int64.Type}, {"SalesOrderDetailID", Int64.Type}

| | 1²₃ SalesOrderID | 1²₃ SalesOrderDetailID | 1²₃ OrderQty | 1²₃ ProductID | 1.2 UnitPrice | 1.2 UnitPrice |
|---|---|---|---|---|---|---|
| 1 | 71774 | 110562 | 1 | 836 | 356.898 | |
| 2 | 71774 | 110563 | 1 | 822 | 356.898 | |
| 3 | 71776 | 110567 | 1 | 907 | 63.9 | |
| 4 | 71780 | 110616 | 4 | 905 | 218.454 | |
| 5 | 71780 | 110617 | 2 | 983 | 461.694 | |
| 6 | 71780 | 110618 | 6 | 988 | 112.998 | |
| 7 | 71780 | 110619 | 2 | 748 | 818.7 | |
| 8 | 71780 | 110620 | 1 | 990 | 323.994 | |
| 9 | 71780 | 110621 | 1 | 926 | 149.874 | |
| 10 | 71780 | 110622 | 1 | 743 | 809.76 | |
| 11 | 71780 | 110623 | 4 | 782 | 1376.994 | |
| 12 | 71780 | 110624 | 2 | 918 | 158.43 | |
| 13 | 71780 | 110625 | 4 | 780 | 1391.994 | |
| 14 | 71780 | 110626 | 1 | 937 | 48.594 | |
| 15 | 71780 | 110627 | 6 | 867 | 41.994 | |
| 16 | 71780 | 110628 | 1 | 985 | 112.998 | |
| 17 | 71780 | 110629 | 2 | 989 | 323.994 | |

Query Settings

▲ PROPERTIES
Name
OrderDetailLines

All Properties

▲ APPLIED STEPS
Source
Promoted Headers
✕ Changed Type

# What can a Power Query Activity do?

**Control Flow**

Power Query

Transform Stuff

JSON

Power Query

Home  Transform  Add column  View

Data view | Schema view | Go to column | Advanced editor

Preview | Columns | Advanced

Queries

ADFResource  [1]
  LakeFileOrderDetailL...
  UserQuery

## Advanced editor

```
1  let
2      AdfDoc = Web.Contents("https://traininglake01.dfs.core.windows.net/datawarehouse/Raw/OrderDetailLines.parquet"),
3      Parquet = Parquet.Document(AdfDoc),
4      #"Grouped rows" = Table.Group(Parquet, {"SalesOrderID"}, {{"Count", each Table.RowCount(_), Int64.Type}})
5  in
6      #"Grouped rows"
```
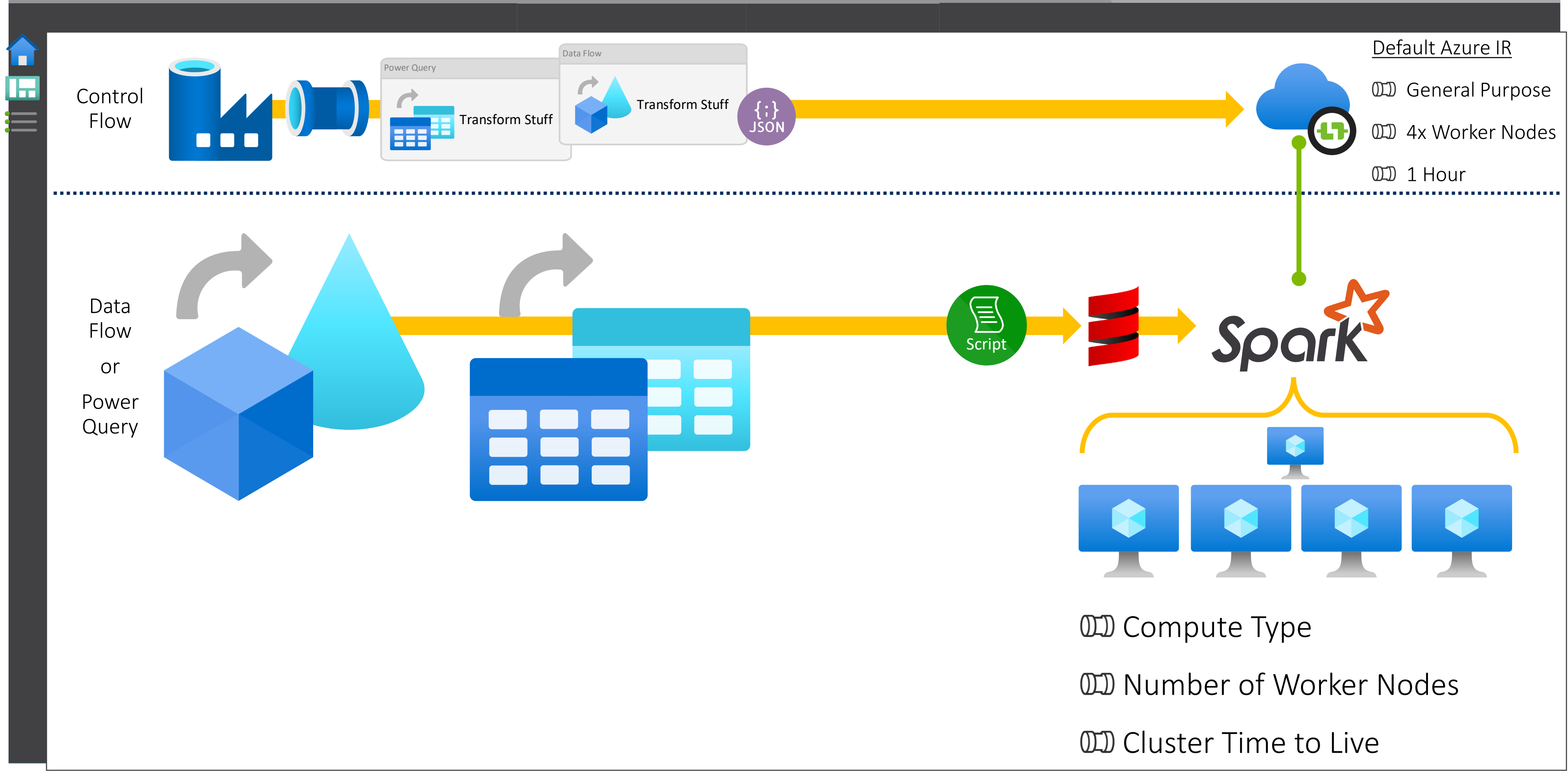
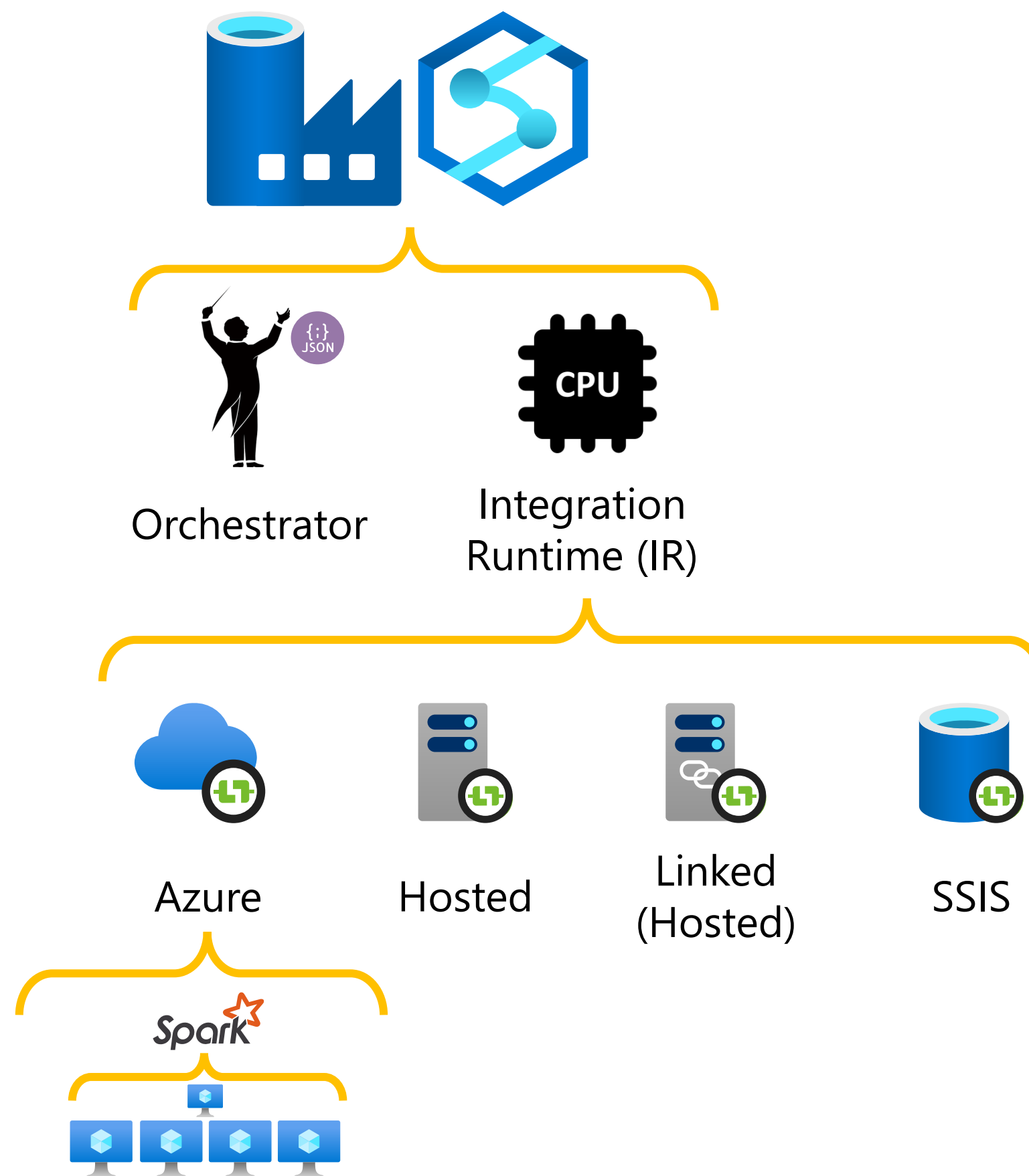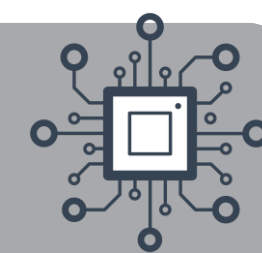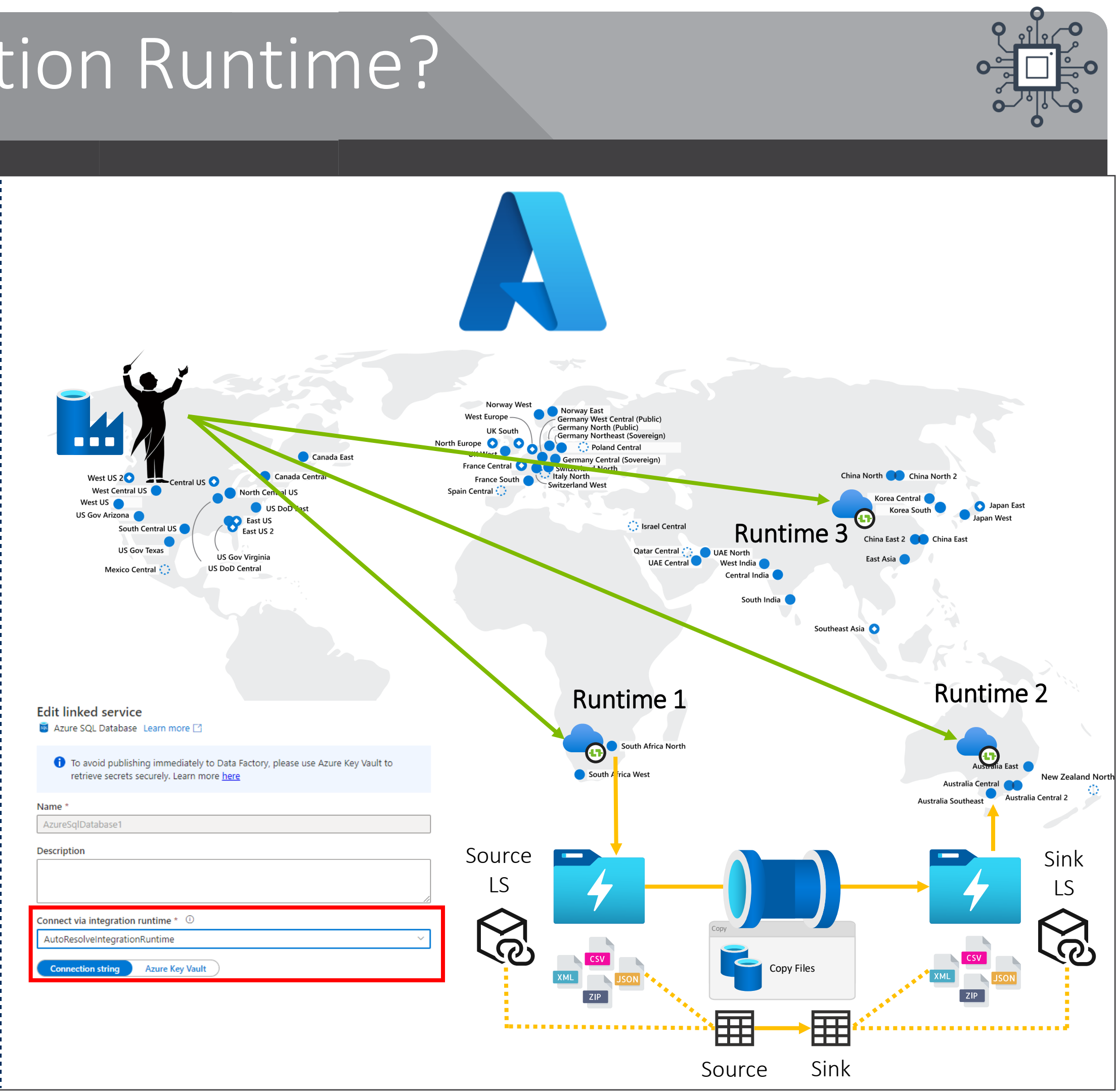OK    Cancel

# Module 3

## Data Transformation



- Data Flows
- Power Query Injection
- **Spark Configuration**
- Use Cases
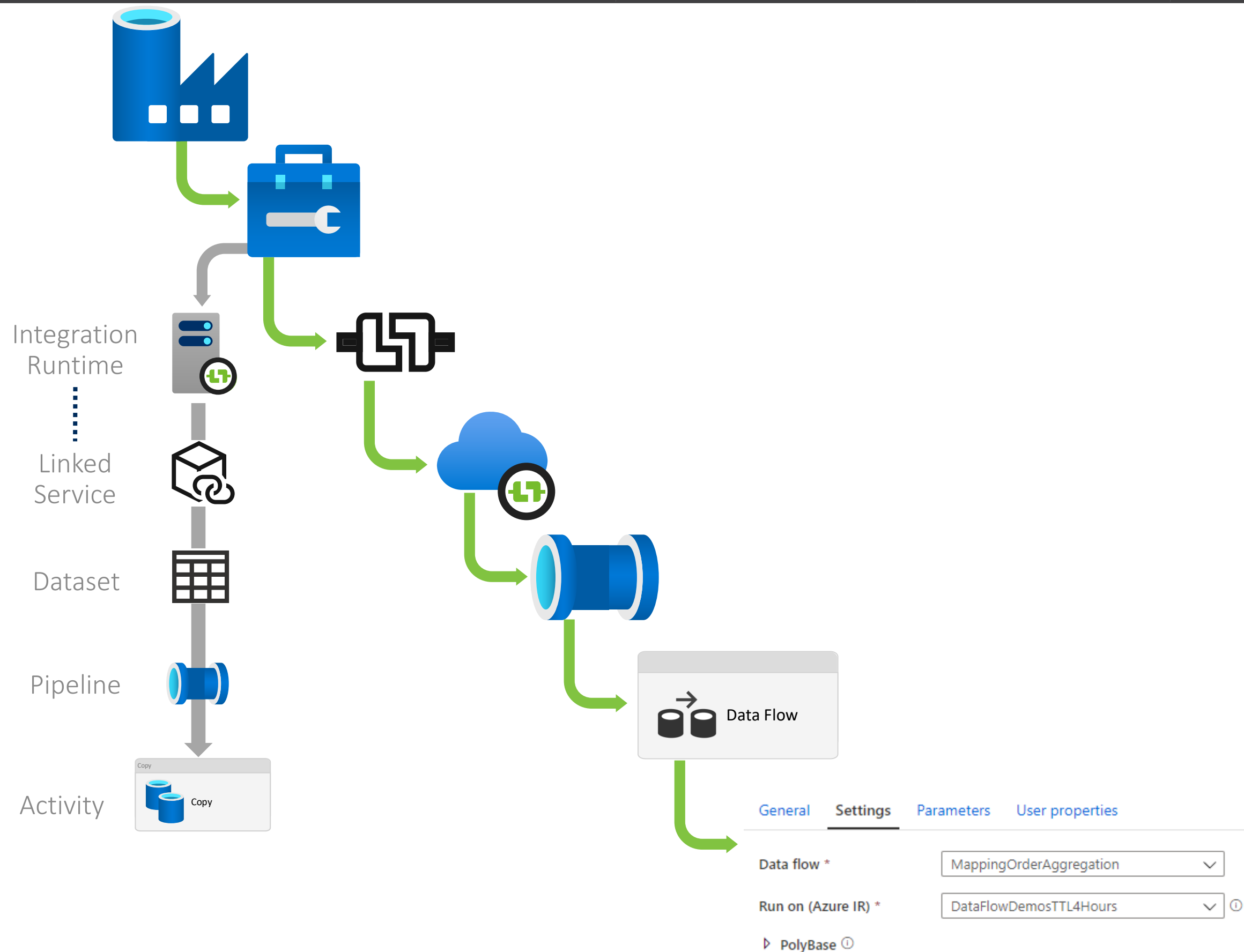
# Spark Configuration

Control Flow

Power Query — Transform Stuff

Data Flow — Transform Stuff

JSON

Default Azure IR

General Purpose

4x Worker Nodes

1 Hour

Data Flow or Power Query

Script

Spark

Compute Type

Number of Worker Nodes

Cluster Time to Live

Orchestrator

Integration Runtime (IR)

Azure

Hosted

Linked (Hosted)

SSIS

Spark

# What is an Integration Runtime?



Orchestrator

Runtime

Fixed Region

Flexible Location

*AutoResolveIntegrationRuntime*

Runtime 3

Runtime 1

Runtime 2

**Edit linked service**

Azure SQL Database  Learn more

To avoid publishing immediately to Data Factory, please use Azure Key Vault to retrieve secrets securely. Learn more here

**Name** *
AzureSqlDatabase1

**Description**

**Connect via integration runtime** *
AutoResolveIntegrationRuntime

Connection string    Azure Key Vault

Source LS

Copy Files

Sink LS
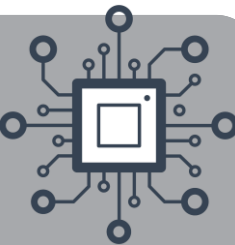
Source

Sink

# Setting the Data Flow Cluster (IR Configuration)



Data Factory

Manage
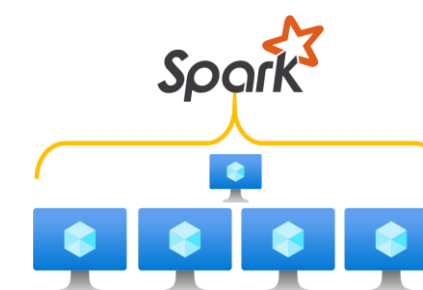
Integration Runtimes

Azure IR

Pipeline

Data Flow Activity
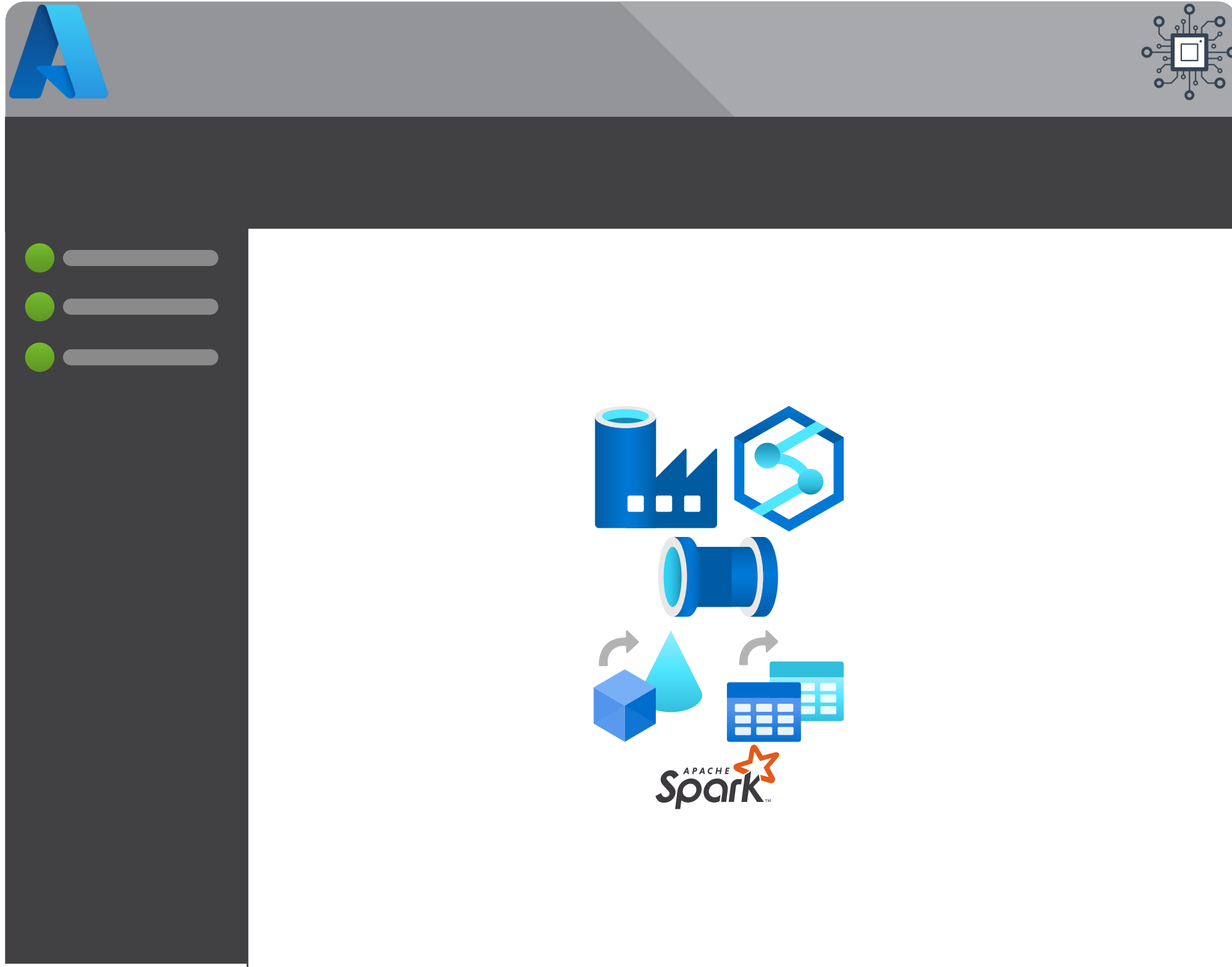
Settings

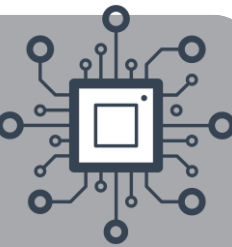Integration Runtime

Linked Service

Dataset

Pipeline

Activity

Copy

Data Flow

General    **Settings**    Parameters    User properties

Data flow *          MappingOrderAggregation

Run on (Azure IR) *    DataFlowDemosTTL4Hours
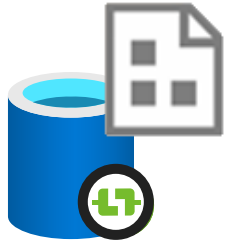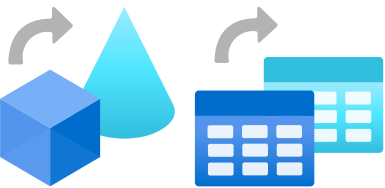
▷ PolyBase ⓘ

Spark

# Module 3

Data Transformation

- Data Flows
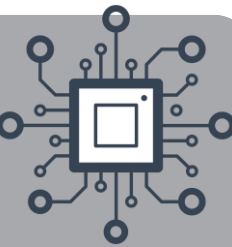- Power Query Injection
- Spark Configuration
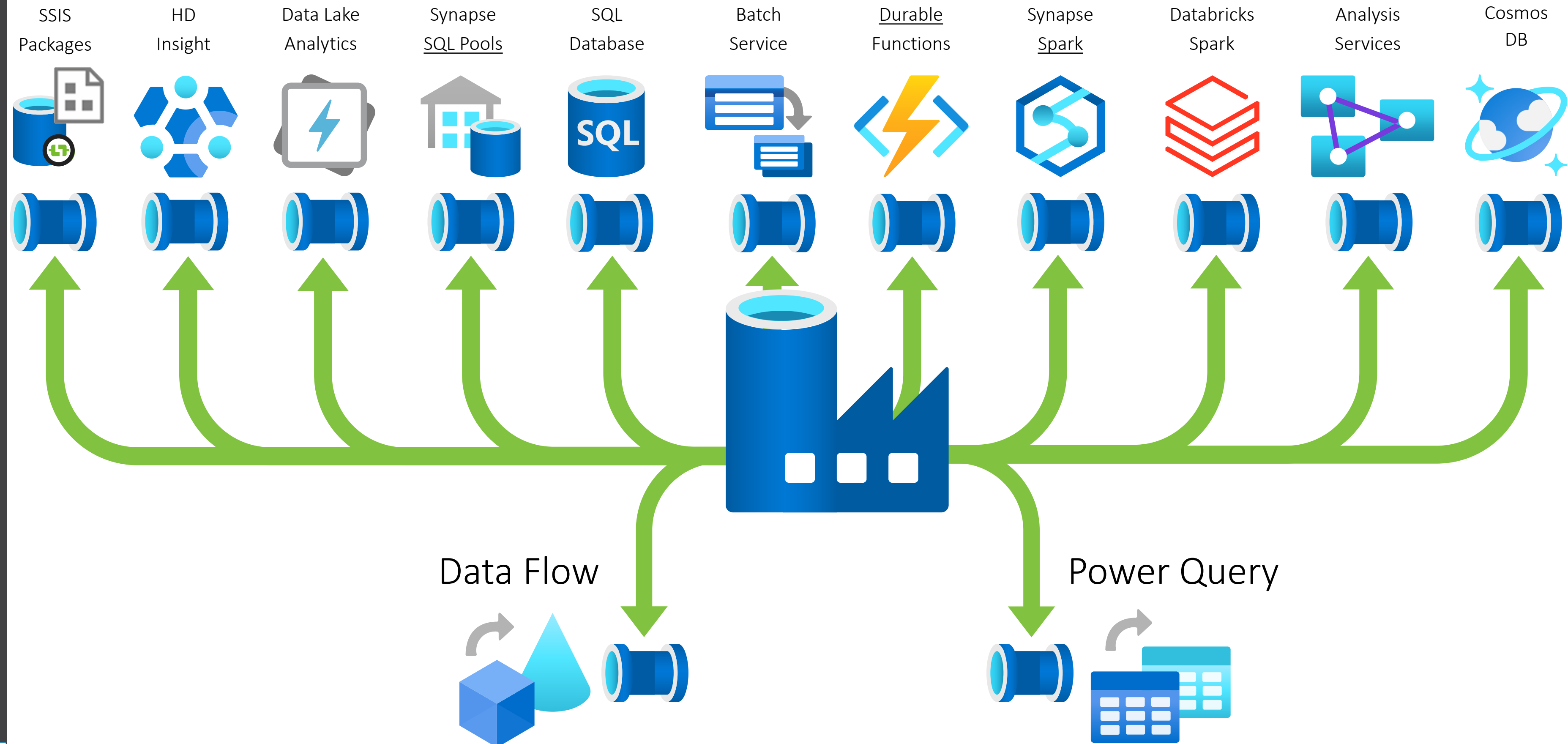- Use Cases

# Data Transformation Resources in Azure Comparison

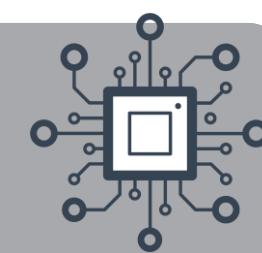| Transformation Tools | Graphical UI (Low/No Code) | Scales Out | Scales Up | Cloud Native Tech |
|---|:---:|:---:|:---:|:---:|
| T-SQL with SQLDB | ✘ | ✘ | ✔ | ✘ |
| SSIS Packages | ✔ | ✘ | ✔ | ✘ |
| Scala/Python/SQL with Databricks | ✘ | ✔ | ✔ | ✔ |
| Data Flows & Power Query | ✔ | ✔ | ✔ | ✔ |

# Other Data Transformation Services in Azure

## When Should We Use These Integration Pipeline Transformation Activities?

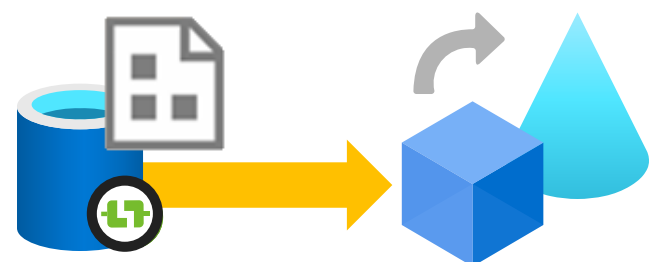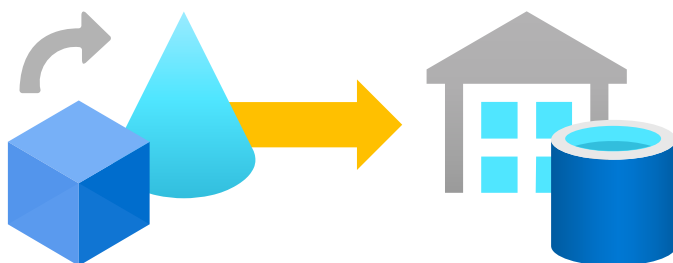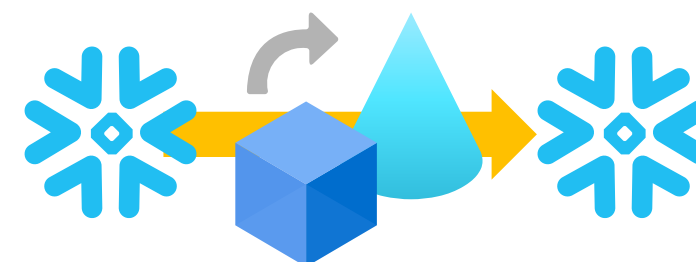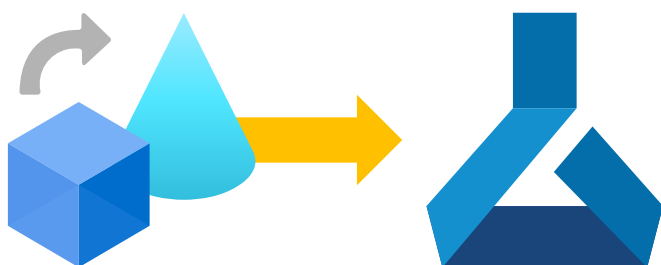| SSIS Packages | HD Insight | Data Lake Analytics | Synapse SQL Pools | SQL Database | Batch Service | Durable Functions | Synapse Spark | Databricks Spark | Analysis Services | Cosmos DB |
|---|---|---|---|---|---|---|---|---|---|---|

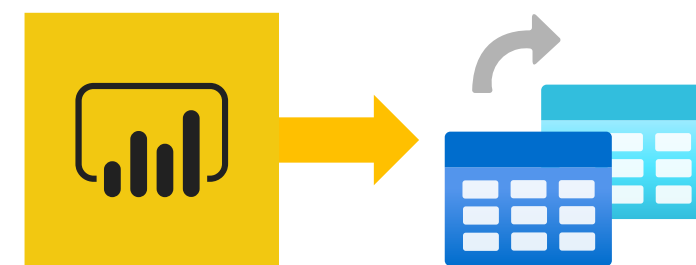Data Flow

Power Query

# Use Cases

**SSIS Package rebuild and skills migration.**

**Warehouse data distribution & loading.**

**Data model dataset preparation.**

**Inline dataset transformations.**

**Power Query industrialisation.**

# Module 3

## Data Transformation

```sql
SELECT
    [Contents]
FROM
    [Training]
WHERE
    [Module] = '3';

END; --module, fetch next
```

- Data Flows
- Power Query Injection
- Spark Configuration
- Use Cases