

Data processing dans Azure

Venez découvrir les différentes options qui s'offrent à vous pour la manipulation des données dans le cloud Azure.

Nous présenterons les différents services vous permettant de déplacer, transformer, appliquer des calculs sur vos données, en prenant en compte les notions bien connues de variétés, vélocités, volumes.

Agenda (2h00)

- Revue des options de processing des données dans le cloud Azure
- Azure Synapse (Pipeline, Data Flow)
- Azure Batch
- processing Spark et serverless
- Compute Azure ML
- Orchestration



GPS Data/AI Strategy FY23

Delivered by CSA Team



Franck Gaillard
Cloud Solution Architect
Data AI
frgail@microsoft.com



Narjes Majdoub
Cloud Solution Architect
Data AI
nmajdoub@microsoft.com



Ali Bouhaddou
Cloud Solution Architect
Data Analytics
albouhad@microsoft.com



Frederic Gisbert
Cloud Solution Architect
Data Analytics
frgisber@microsoft.com

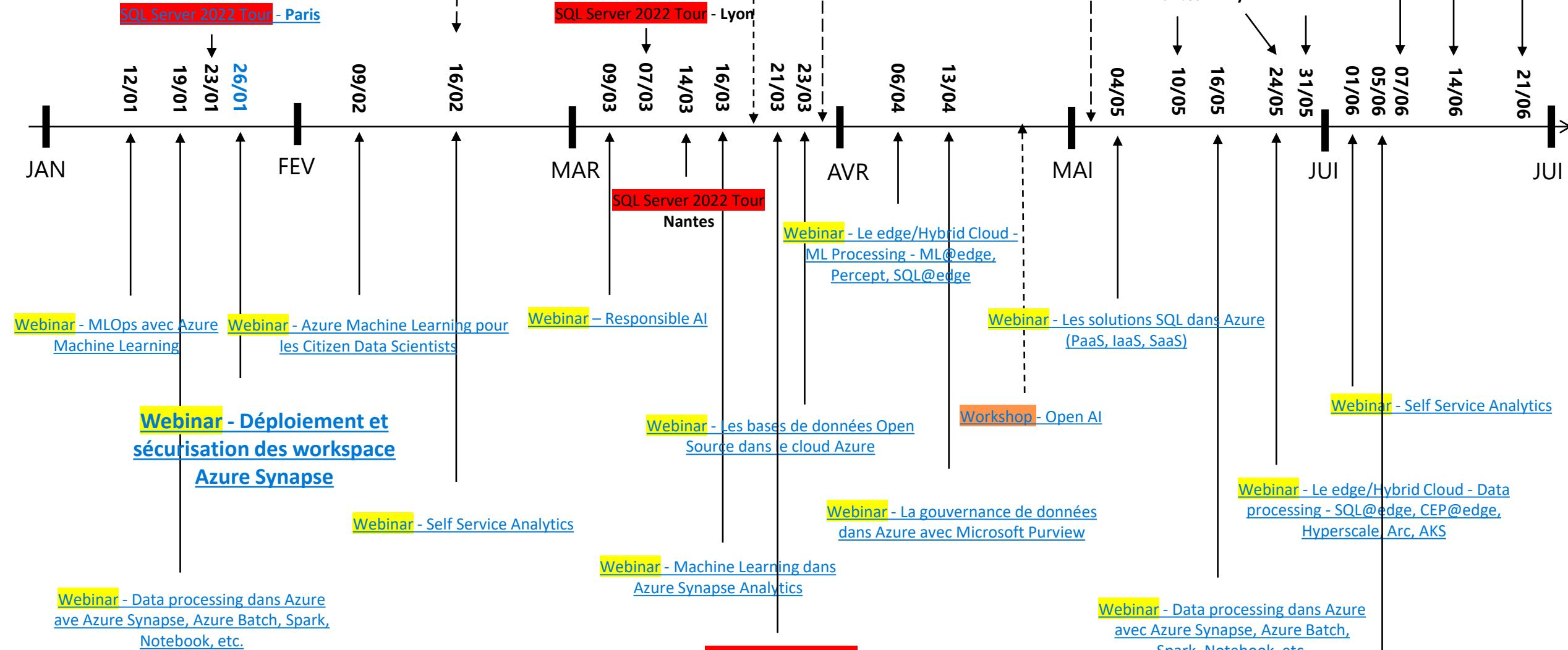


Azure Data & AI technical intensity plan

- From June 2022 to June 2023
- Focus on "Azure Data & AI" tech intensity
- Many content, from L100 Beginner to L400 Expert level:
 - Academy L100
 - Webinar L200/L300
 - Workshop L300/L400
 - Certification kickstart L300/L400
 - Openhack / Microhack L400

Plan GPS Data/AI global FY23 (H2)

cademy DP900/AI900



Webinar/Academy- 200/300

Workshop/ OpenHack/ Certifications - 300/400

SQl Server 2022 Tour - 300/400

Mons - Belgique

Workshop - Actualités des plateformes Azure Synapse Spark & Machine Learning services

+ Le low code avec Azure Machine Learning designer, pipelines et Datamart Power BI

Liste des évènements de type Webinar 2H

Event Webinar (Les jeudis de la Data & AI) - L200/300	Date	Duration (min)	Link
Azure Machine Learning pour les Data Scientists	15/09/2022	120	https://msevents.microsoft.com/event?id=2454281594
Azure Synapse	22/09/2022	120	https://msevents.microsoft.com/event?id=857781749
Les solutions SQL dans Azure (PaaS, IaaS, SaaS)	29/09/2022	120	https://msevents.microsoft.com/event?id=502366997
Déploiement et sécurisation des workspaces Azure Machine learning	06/10/2022	120	https://msevents.microsoft.com/event?id=1505714138
Azure Scale Analytics - Architectures Data Mesh dans Azure avec Azure Synapse, Microsoft Purview et Azure Data Share	13/10/2022	120	https://msevents.microsoft.com/event?id=139685175
MLOps avec Azure Machine Learning	20/10/2022	120	https://msevents.microsoft.com/event?id=1245885767
SQL Server 2022 et hybridation native avec Azure SQL Managed Instance	10/11/2022	120	https://msevents.microsoft.com/event?id=145826476
Machine Learning dans Azure Synapse Analytics	17/11/2022	120	https://msevents.microsoft.com/event?id=3637723312
Azure Cosmos DB et IA	24/11/2022	120	https://msevents.microsoft.com/event?id=2646013445
Azure et les Services Cognitifs	08/12/2022	120	https://msevents.microsoft.com/event?id=3772037220
La gouvernance de données dans Azure avec Microsoft Purview	15/12/2022	120	https://msevents.microsoft.com/event?id=1499560981
MLOps avec Azure Machine Learning	12/01/2023	120	https://msevents.microsoft.com/event?id=4115194515
Data processing dans Azure ave Azure Synapse, Azure Batch, Spark, Notebook, etc.	19/01/2023	120	https://msevents.microsoft.com/event?id=1537241181
Déploiement et sécurisation des workspace Azure Synapse	26/01/2023	120	https://msevents.microsoft.com/event?id=1806467748
Azure Machine Learning pour les Citizen Data Scientists	09/02/2023	120	https://msevents.microsoft.com/event?id=2367521229
PowerBI - Self Service Analytics	16/02/2023	120	https://msevents.microsoft.com/event?id=1401519679
L'IA responsable avec Azure machine learning	09/03/2023	120	https://msevents.microsoft.com/event?id=2072953112
Machine Learning dans Azure Synapse Analytics	16/03/2023	120	https://msevents.microsoft.com/event?id=3413014857
Les bases de données Open Source dans le cloud Azure	23/03/2023	120	https://msevents.microsoft.com/event?id=2727487131
Hybridation des services de Machine Learning Azure	06/04/2023	120	https://msevents.microsoft.com/event?id=1624914222
La gouvernance de données dans Azure avec Microsoft Purview	13/04/2023	120	https://msevents.microsoft.com/event?id=3909342839
Les solutions SQL dans Azure (PaaS, IaaS, SaaS)	04/05/2023	120	https://msevents.microsoft.com/event?id=1162207895
Data processing dans Azure ave Azure Synapse, Azure Batch, Spark, Notebook, etc.	16/05/2023	120	https://msevents.microsoft.com/event?id=3517068442
Hybridation des services de données Azure	24/05/2023	120	https://msevents.microsoft.com/event?id=2996507398
Self Service Analytics	01/06/2023	120	https://msevents.microsoft.com/event?id=2387937988

Total

25

Liste des évènements de type Workshop/Prepa Cert/Academy

Event Workshop L300/400	Date	Duration (min)	Link
Synapse et Lakehouse Database	18/10/2022	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdURE1RMVgwTDNISTE1TDFYSDVLR0cyS1kwWS4u
Le NoSQL dans Azure - Migration CosmosDB / OSS	08/11/2022	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdURE1RMVgwTDNISTE1TDFYSDVLR0cyS1kwWS4u
Lab Lyon - Data & AI	22/11/2022	240	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUMIZZOURET0RSWjcyTERYRkJGTIFFUJaUi4u
Lab Nantes - Data & AI	29/11/2022	240	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUMIZZOURET0RSWjcyTERYRkJGTIFFUJaUi4u
L'IA dans les services de données Azure	12/12/2022	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdURE1RMVgwTDNISTE1TDFYSDVLR0cyS1kwWS4u
Open AI	H2	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdURE1RMVgwTDNISTE1TDFYSDVLR0cyS1kwWS4u

Event Academy, kickstart certifications, workshop certifications	Date	Duration (min)	Link
MCPP - DP-420	10/10/2022	420	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUMkJSIRKSU1RRFA00VgzSFdTSTY0RE9WQy4u
Micro Hack CosmosDB	20/10/2022	420	H1 - Inscriptions PTA
Academy DP900	17-21/10/2022	300	https://msevents.microsoft.com/event?id=3250818161
Academy AI900	17-21/10/2022	300	https://msevents.microsoft.com/event?id=2717528090
Kickstart DP-500	17/10/2022	60	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUNEk3WFQ1TEdNNTQ2Uk85V0cxQzM3TE9ZRS4u
Dry Run DP-500	14/11/2022	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUNEk3WFQ1TEdNNTQ2Uk85V0cxQzM3TE9ZRS4u
Q&A DP-500	05/12/2022	90	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUNEk3WFQ1TEdNNTQ2Uk85V0cxQzM3TE9ZRS4u
Kickstart DP-100	17/10/2022	60	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUNDAxV0hSN0FHM1YzUzI3OUNMFYxSkRIMi4u
Dry Run DP-100	14/11/2022	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUNDAxV0hSN0FHM1YzUzI3OUNMFYxSkRIMi4u
Q&A DP-100	05/12/2022	90	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUNDAxV0hSN0FHM1YzUzI3OUNMFYxSkRIMi4u
Kickstart DP-203	17/10/2022	60	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUOVFWOUVCNFcyQk5SVjFBUFczNktCUFpLMi4u
Dry Run DP-203	14/11/2022	120	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUOVFWOUVCNFcyQk5SVjFBUFczNktCUFpLMi4u

Q&A DP-203	05/12/2022	90	https://forms.office.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHB3zwJTO3s11AuapNnBbrwdUOVFWOUVCNFcyQk5SVjFBUFczNktCUFpLMi4u
------------	------------	----	---

Liste des évènements de type Workshop/Prepa Cert/Academy

Event Type	Date	Duration (min)	Link
SQL Server 2022 Tour L300/400			
Session Paris	23/01/2023	480	https://msevents.microsoft.com/event?id=1969809187
Session Lyon	07/03/2023	480	
Session Nantes	14/03/2023	480	
Session Toulouse	21/03/2023	480	
Tour de France - Data/AI services L300/400			
Nantes - Session Azure Synapse SQL-Spark-ML / Azure Machine Learning	10/05/2023	480	https://forms.office.com/r/P6fUq2VJmv
Nantes - Session low code - Azure Machine Learning, Pipelines, Datamarts Power BI	07/06/2023	480	https://forms.office.com/r/P6fUq2VJmv
Toulouse - Session Azure Synapse SQL-Spark-ML / Azure Machine Learning	31/05/2023	480	https://forms.office.com/r/P6fUq2VJmv
Toulouse - Session low code - Azure Machine Learning, Pipelines, Datamarts Power BI	21/06/2023	480	https://forms.office.com/r/P6fUq2VJmv
Lyon - Session Azure Synapse SQL-Spark-ML / Azure Machine Learning	24/05/2023	480	https://forms.office.com/r/P6fUq2VJmv
Lyon - Session low code - Azure Machine Learning, Pipelines, Datamarts Power BI	14/06/2023	480	https://forms.office.com/r/P6fUq2VJmv

Data as a strategic asset



Data-driven transformations yield significant benefits

54%

increase in
revenue performance

44%

faster time
to market

62%

improvement in
customer satisfaction

54%

increased
profit results

Today's data realities

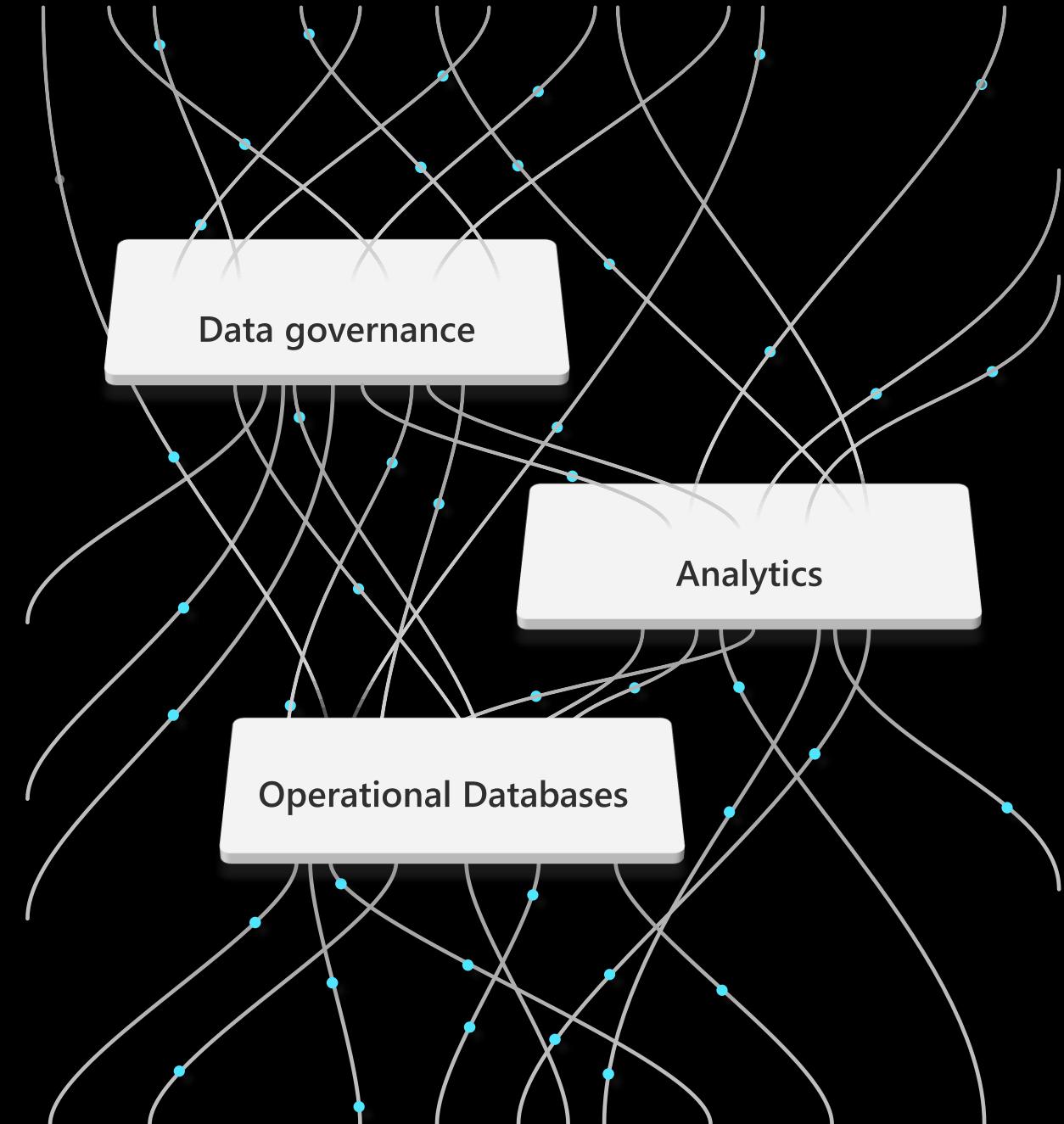
What **data** do I have?

Is it **trustworthy**?

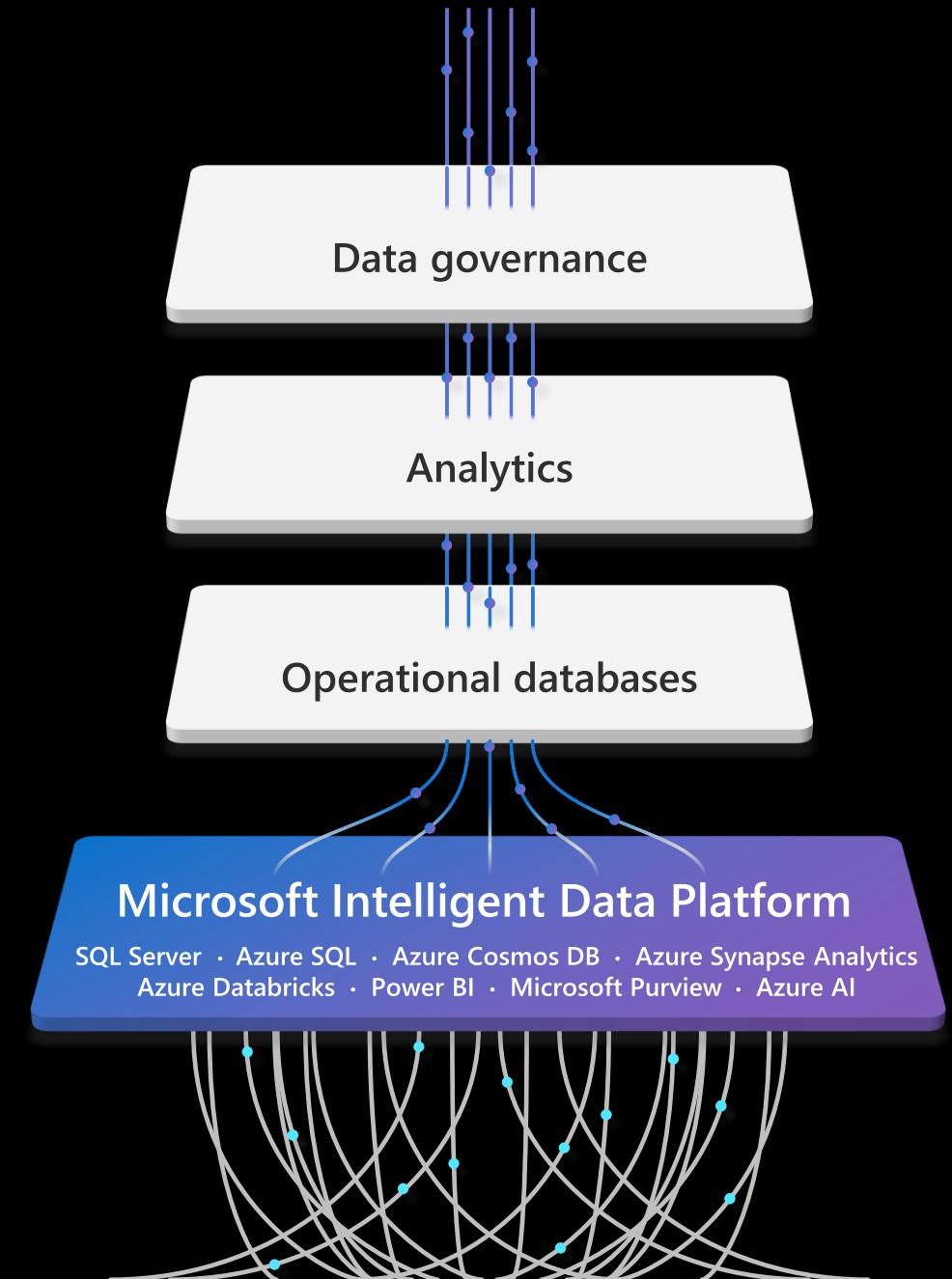
Can people access the **data** needed
to make the right decisions?

How can I **enable faster**
business insights?

What's my **compliance exposure**?



Introducing Microsoft Intelligent Data Platform



Transform with the Microsoft Cloud

Microsoft Cloud

Dynamics 365 · Microsoft 365

Power Platform · Microsoft Azure

Microsoft Intelligent Data Platform

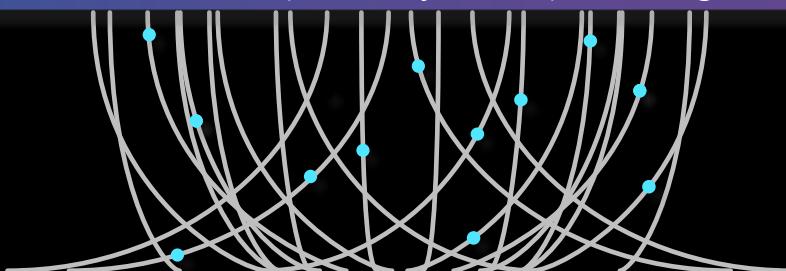
SQL Server · Azure SQL · Azure Cosmos DB · Azure Synapse Analytics

Azure Databricks · Power BI · Microsoft Purview · Azure AI

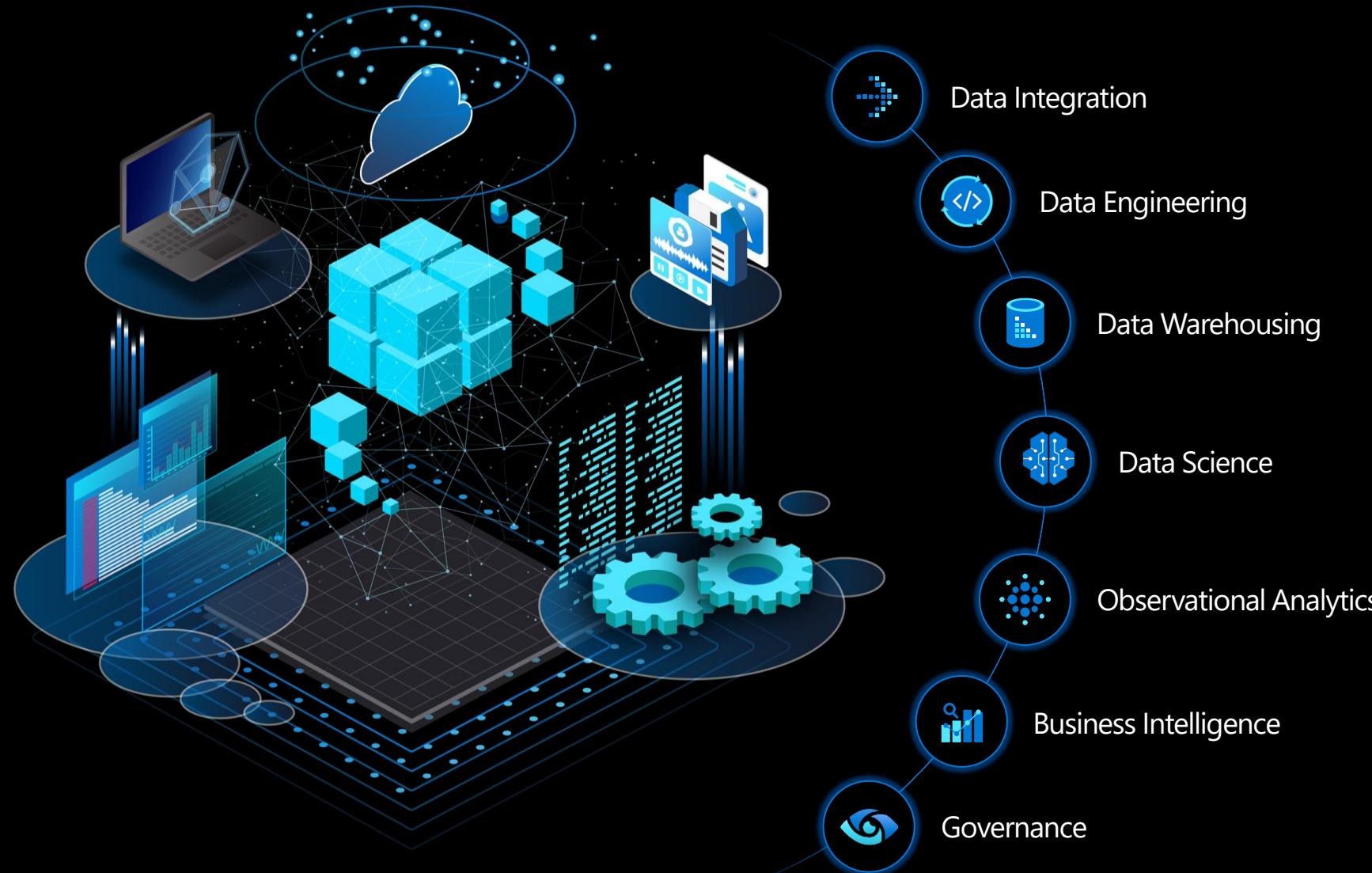
Operational databases

Analytics

Data governance



Processing de données dans Azure





Data Engineering

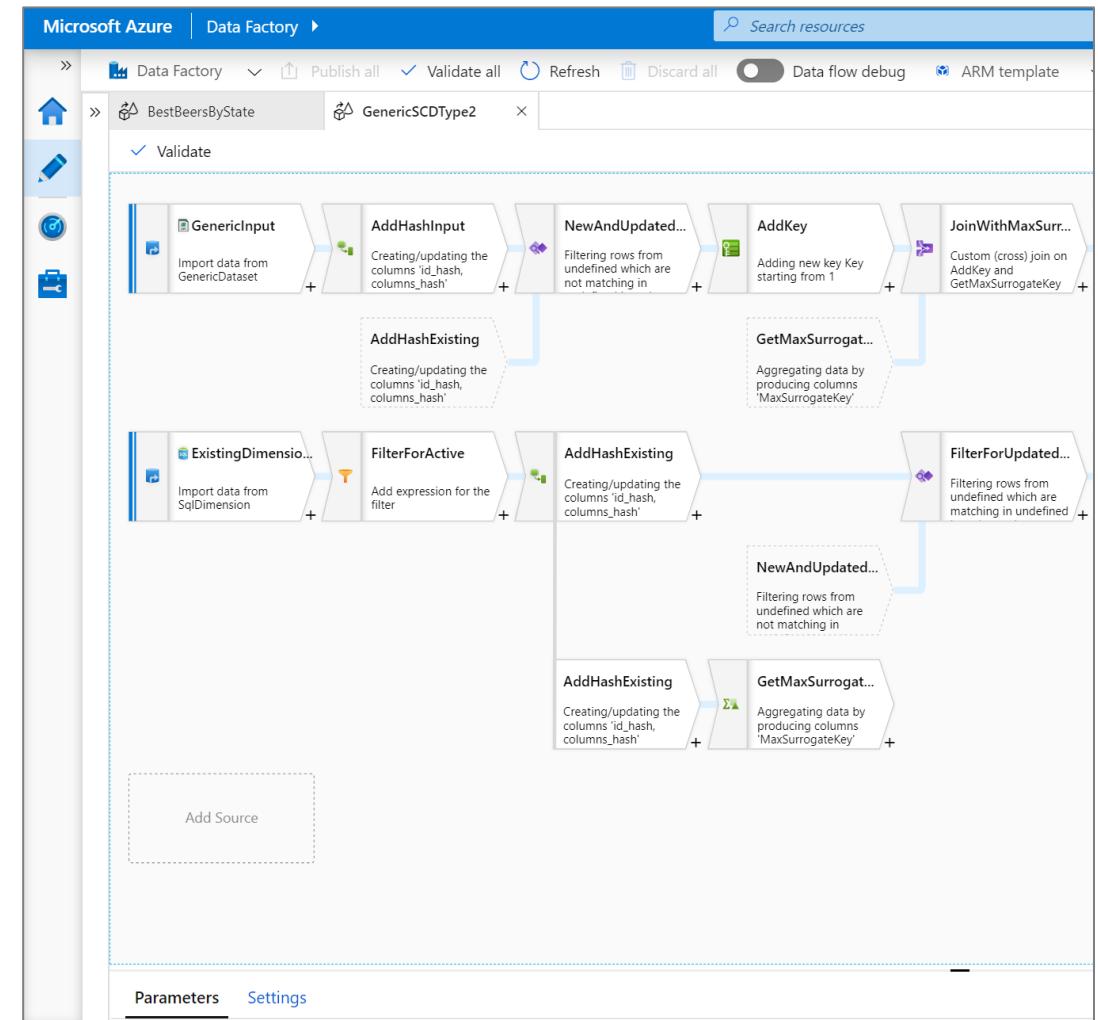


Mapping Data Flows

Code-free data transformation at scale with Azure Data Factory and
Synapse Analytics

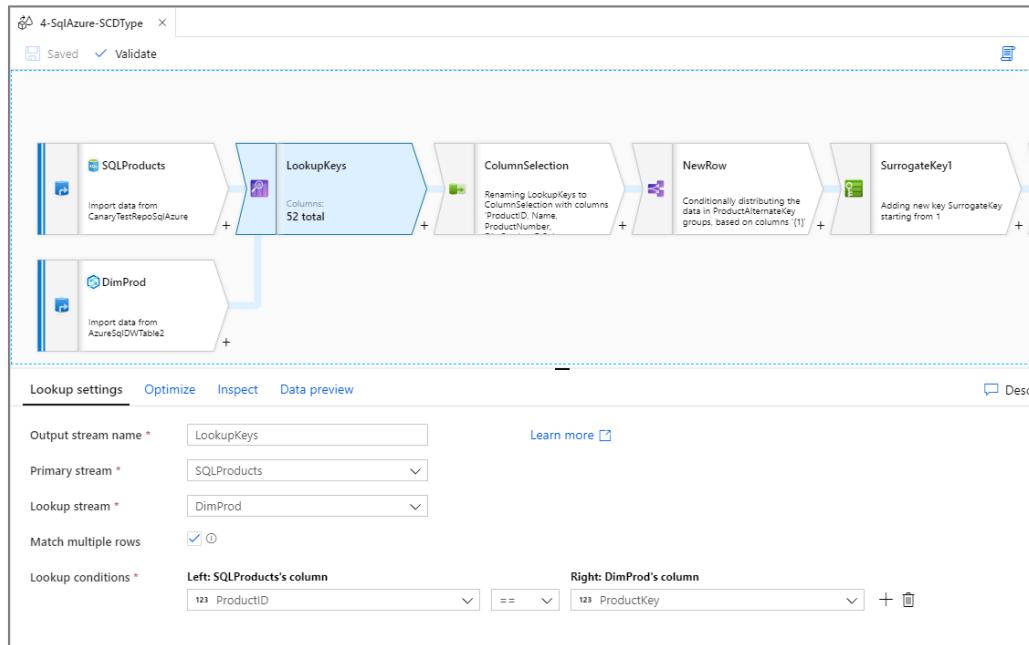
What are mapping data flows?

- Code-free data transformation at scale
- Serverless, scaled-out, ADF/Synapse-managed Apache Spark™ engine
- Resilient flows handle structured and unstructured data
- Operationalized as a data pipeline activity



Code-free data transformation at scale

- Intuitive UX lets you focus on building transformation logic
 - Data cleansing
 - Data validation
 - Data aggregation
- No requirement of knowing Spark, cluster management, Scala, Python, etc



vs

```
class ExampleJob extends JobSparkSQL {
    def run(args: Array[String]): Unit = {
        val transactions = spark.read()
            .option("header", true)
            .option("sep", ",")
            .option("inferSchema", true)
            .option("maxRows", 1000000)
            .option("mode", "PERMISSIVE")
            .csv(args(0))
        val t = transactions.rdd
        t.co-partitioned(1000000)
    }

    val schema = StructType(List(
        StructField("id", IntegerType, true),
        StructField("transactionId", IntegerType, true),
        StructField("productId", IntegerType, true),
        StructField("category", StringType, true),
        StructField("amount", DoubleType, true),
        StructField("date", DateType, true)
    ))

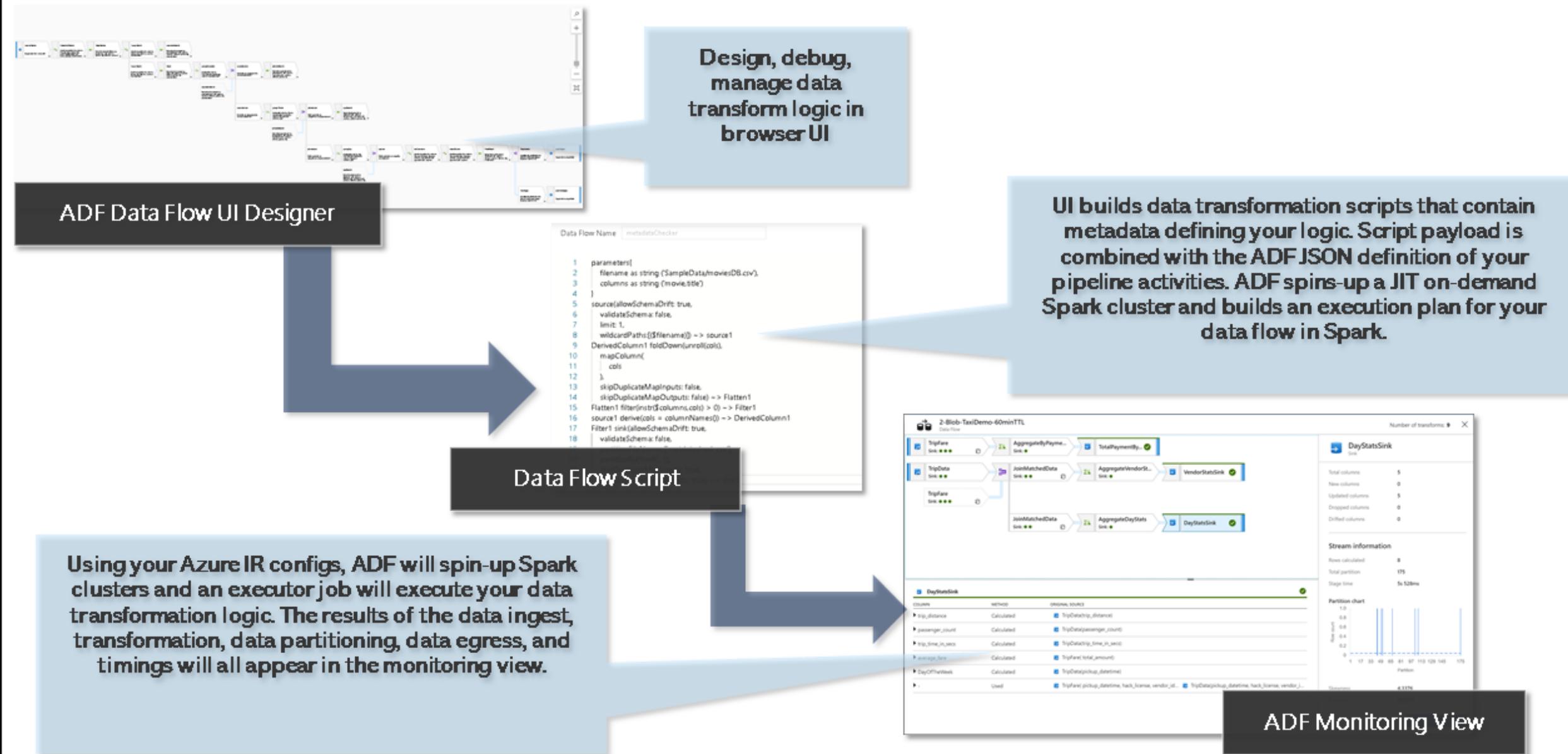
    val result = transactions.mapPartitions(partition: PartitionedDataset[Row] =>
        partition.map(row => {
            val transaction = row.as[ExampleJob.Transaction]
            val id = transaction.id
            val transactionId = transaction.transactionId
            val productId = transaction.productId
            val category = transaction.category
            val amount = transaction.amount
            val date = transaction.date
            val surrogateKey = transaction.surrogateKey
            val row = Row(id, transactionId, productId, category, amount, date, surrogateKey)
            Row(row)
        })
    )

    val resultDF = result.toDF(schema)
    resultDF.write().mode("Overwrite").parquet(args(1))
    return resultDF.parallelize(resultDF.rdd.toLocalIterator).map(_.as[ExampleJob.Transaction])
}

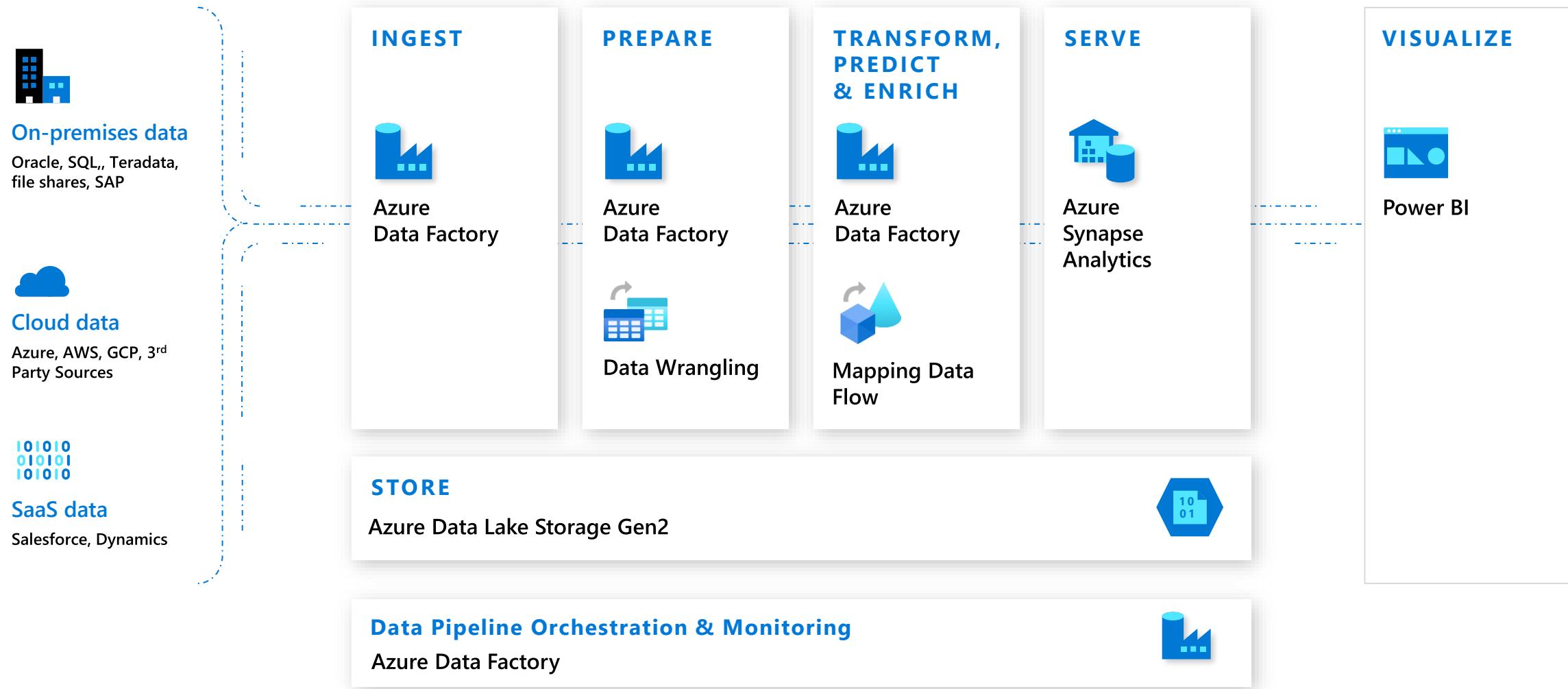
def processData(id: Int, transactionId: Int, productId: Int, category: String, amount: Double, date: Date, surrogateKey: Long): ExampleJob.Transaction = {
    ExampleJob.Transaction(id, transactionId, productId, category, amount, date, surrogateKey)
}

object ExampleJob {
    def main(args: Array[String]): Unit = {
        val transactionsPath = args(0)
        val outputDir = args(1)
        val rdd = new ParallelReader(transactionsPath, "parquet", "inferSchema", "true")
        val schema = new StructType(List(
            StructField("id", IntegerType, true),
            StructField("transactionId", IntegerType, true),
            StructField("productId", IntegerType, true),
            StructField("category", StringType, true),
            StructField("amount", DoubleType, true),
            StructField("date", DateType, true)
        ))
        val result = rdd.mapPartitions(partition: PartitionedDataset[Row] =>
            partition.map(row => {
                val transaction = row.as[ExampleJob.Transaction]
                val id = transaction.id
                val transactionId = transaction.transactionId
                val productId = transaction.productId
                val category = transaction.category
                val amount = transaction.amount
                val date = transaction.date
                val surrogateKey = transaction.surrogateKey
                val row = Row(id, transactionId, productId, category, amount, date, surrogateKey)
                Row(row)
            })
        )
        val resultDF = result.toDF(schema)
        resultDF.write().mode("Overwrite").parquet(outputDir)
    }
}
```

ADF Data Flows Service Architecture

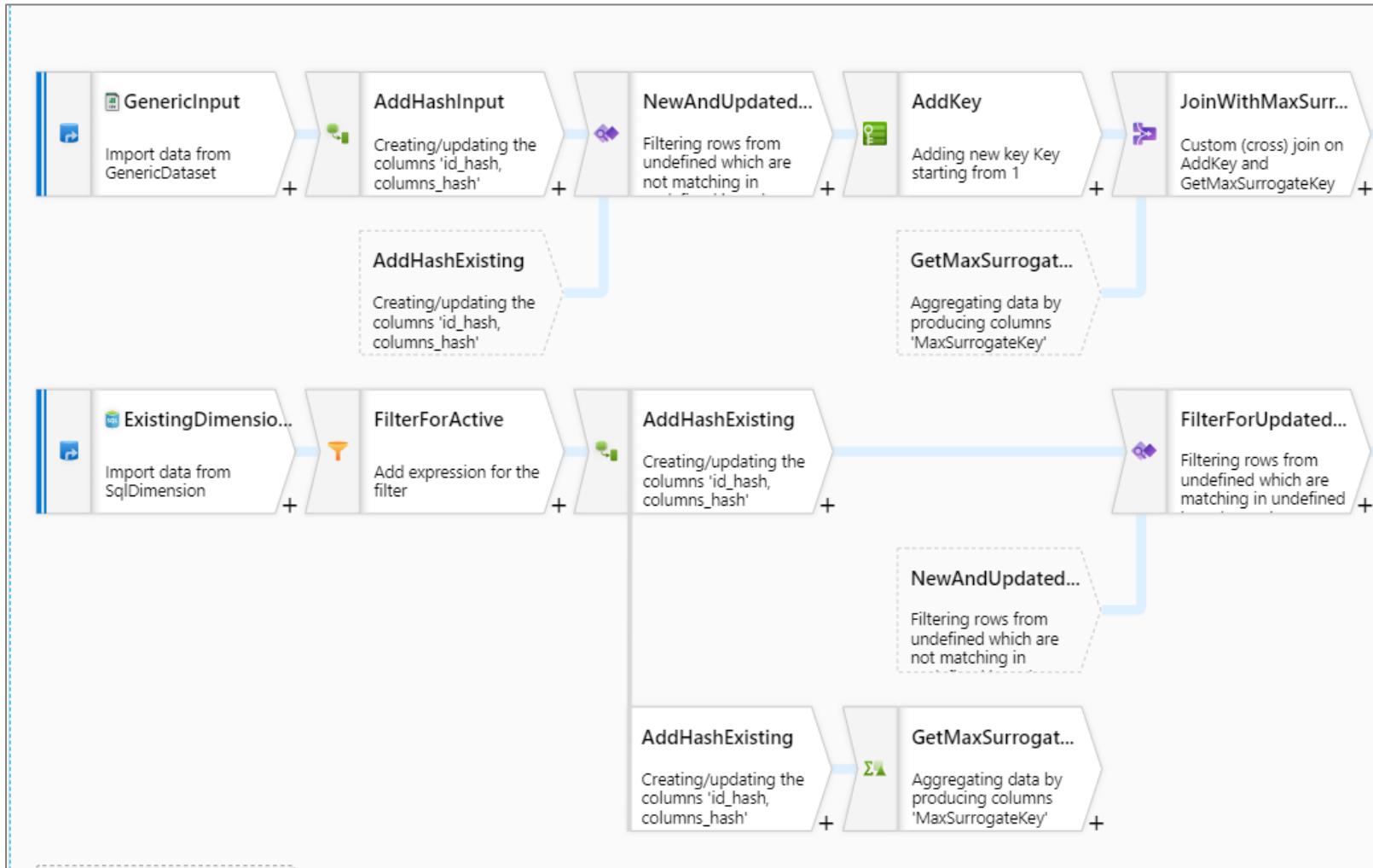


Modern Data Warehouse (MDW)

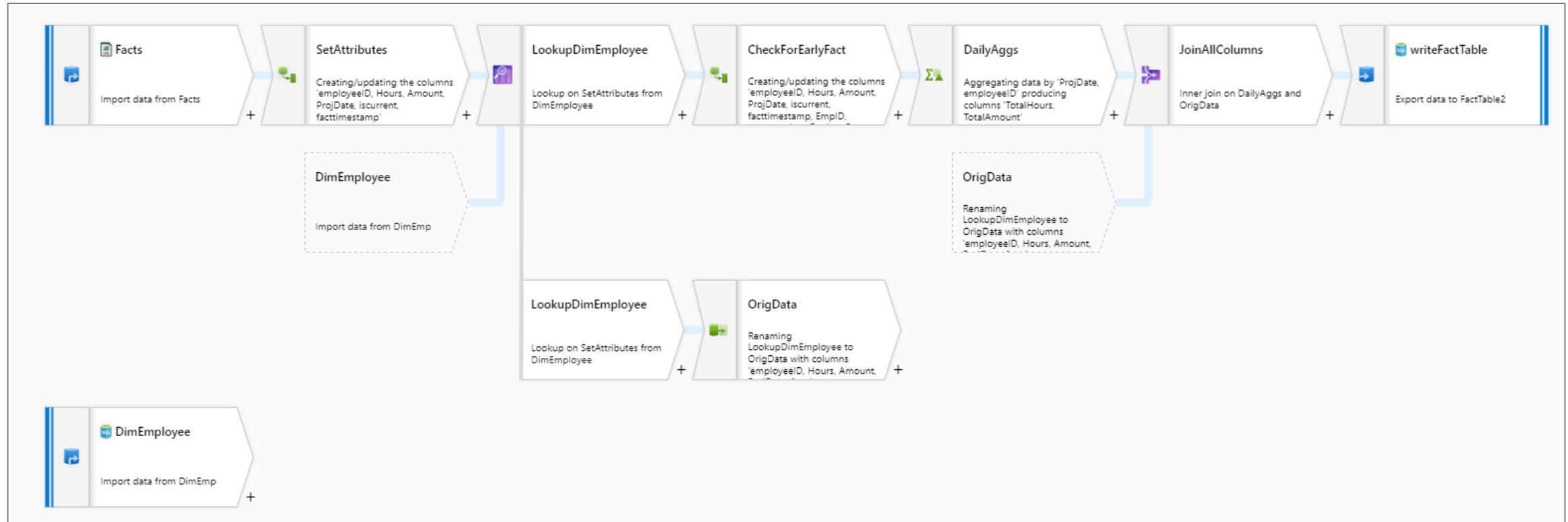


Common data flow scenarios

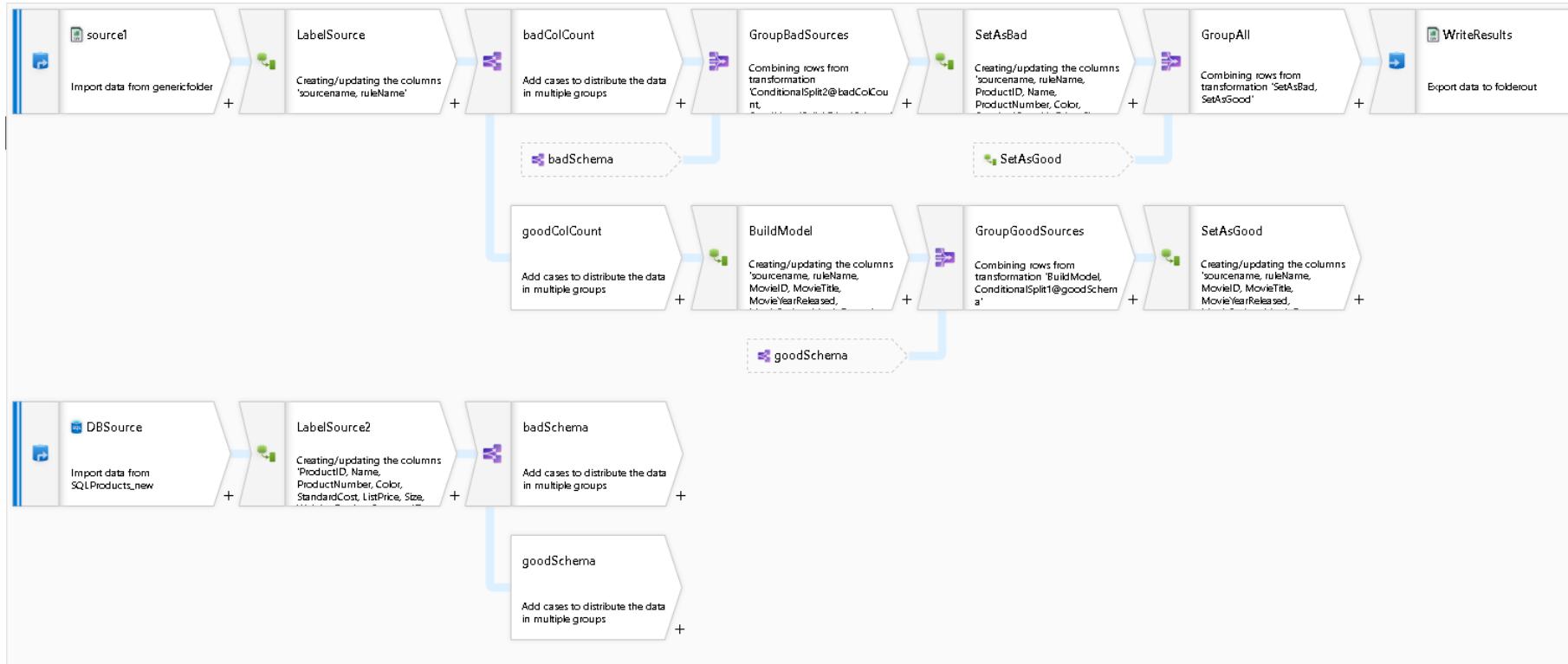
Slowly changing dimensions



Fact loading into a data warehouse

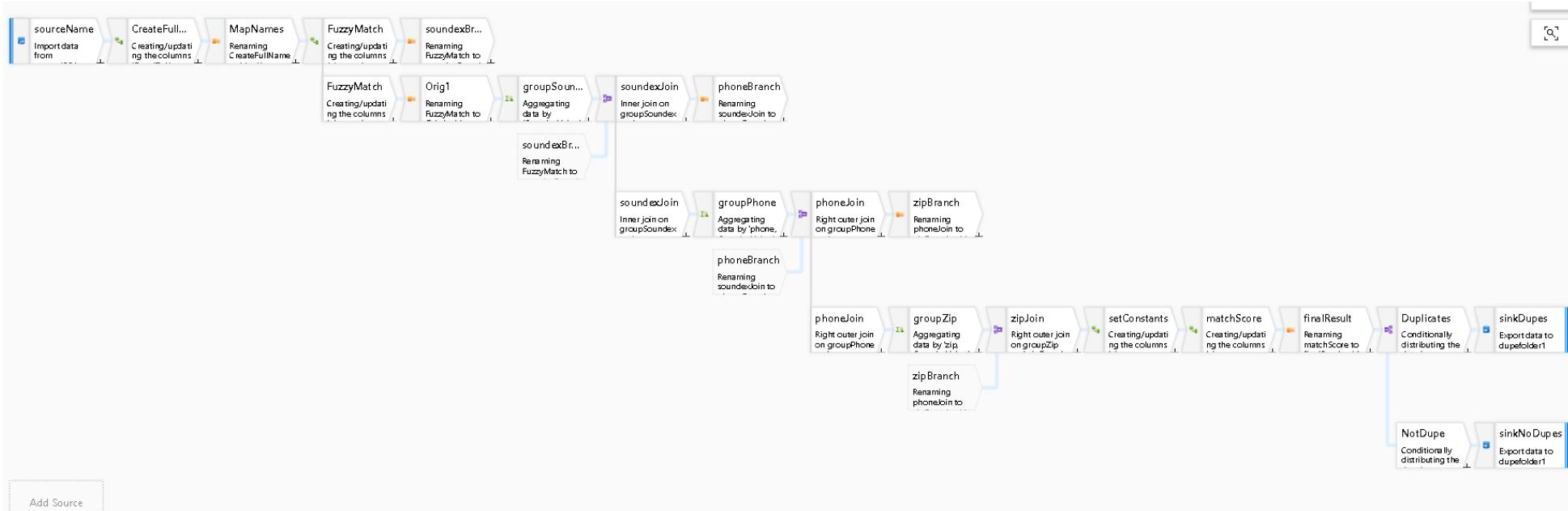


Metadata Validation Rules

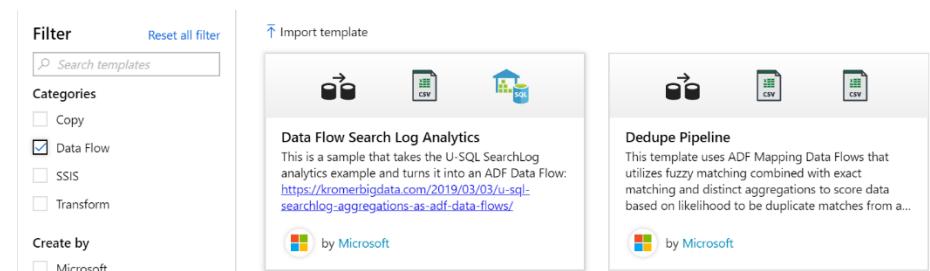


- Transform data conditionally based on metadata traits
- Create and manage metadata quality rules
- Manipulate column properties of source data
- https://www.youtube.com/watch?v=E_UD3R-VpYE

Data De-Duplication and Distinct Rows

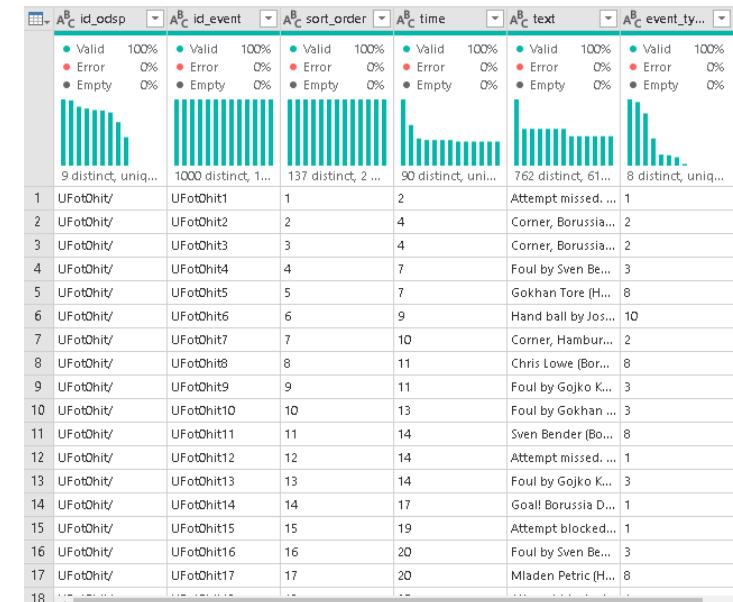


- Use this pattern to eliminate common rows from your data
- You pick a heuristic to use during duplicate matching
- You can tag rows and/or remove duplicate rows
- Use exact matching and/or fuzzy matching
- Available as pipeline template *Dedupe Pipeline*



Data Profiling

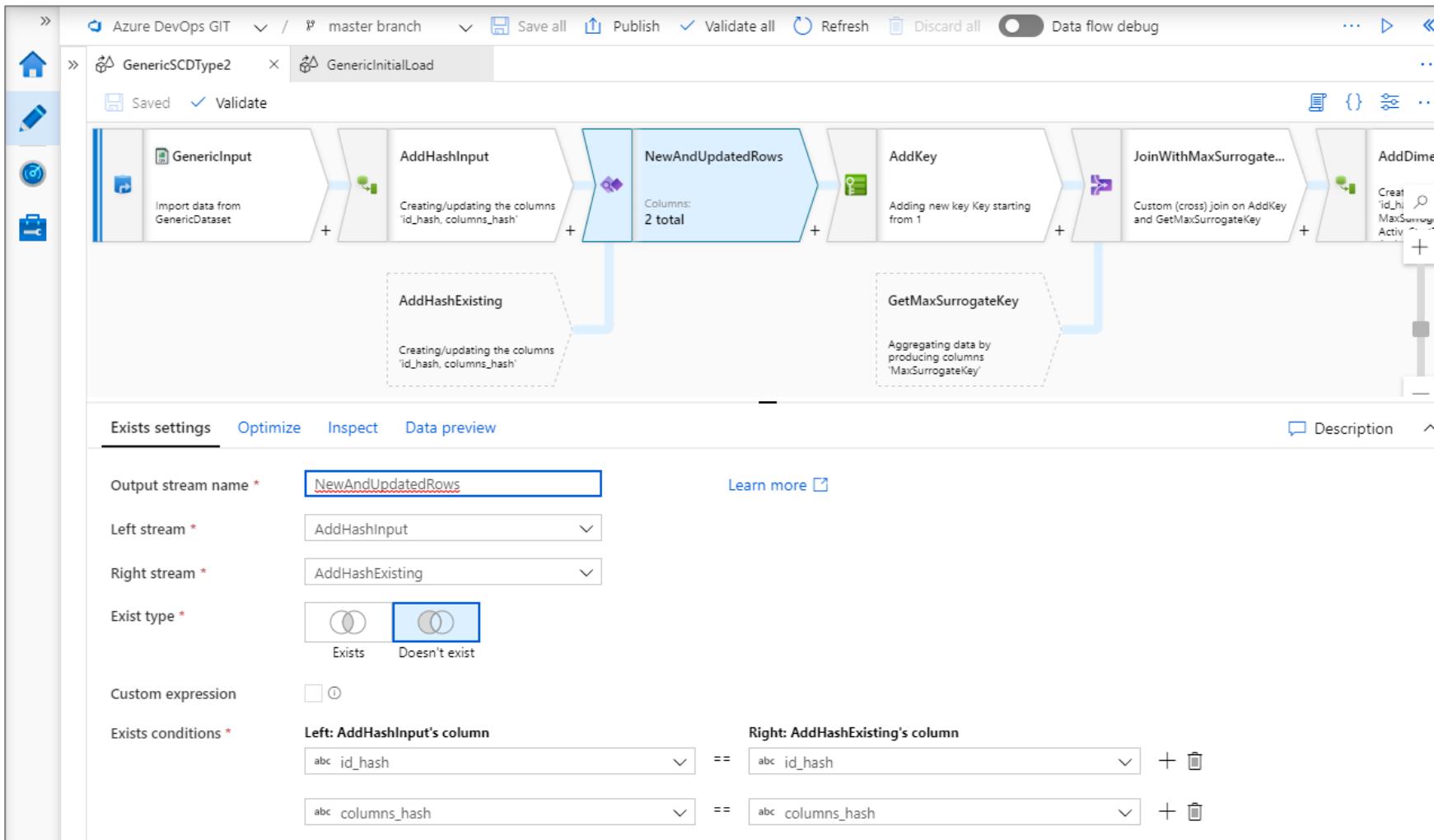
The screenshot shows the Azure Data Factory Data Preview interface. At the top, there are tabs for Source settings, Source options, Projection, Optimize, Inspect, and Data preview (which is selected). Below this, a summary bar shows: Number of rows (941009), INSERT 100, UPDATE 0, DELETE 0, UPSERT 0, LOOKUP 0, and TOTAL 941009. A 'Description' button is also present. The main area displays a table of data with columns: id_ode, event, sort_order, time, text, and event_type. The table contains 12 rows of football match events. To the right of the table is a detailed statistics panel for the 'event_type' column, showing: Count (726716, 77.2%, NA), % Data (167859, 17.8%, 12), 43475 (4.6%, 13), 2258 (0.2%, 14), 701 (0.1%, 15), 0 (0.0%, Remaining values), 0 (0.0%, Null), Not Null (941009), Null (0), Maximum Length (2), and Minimum Length (2).



- Summary statistics describing the shape, size, content of your data
- Helps you to understand what is inside your data
- As a Data Engineer, you have the responsibility to provide proper data for analytics and models
- For big data sets, use sampling / good enough approach
- View data statistics at any step in your data transformation
- Here is how to persist your data profile stats: <https://techcommunity.microsoft.com/t5/azure-data-factory/how-to-save-your-data-profiler-summary-stats-in-adf-data-flows/ba-p/1243251>

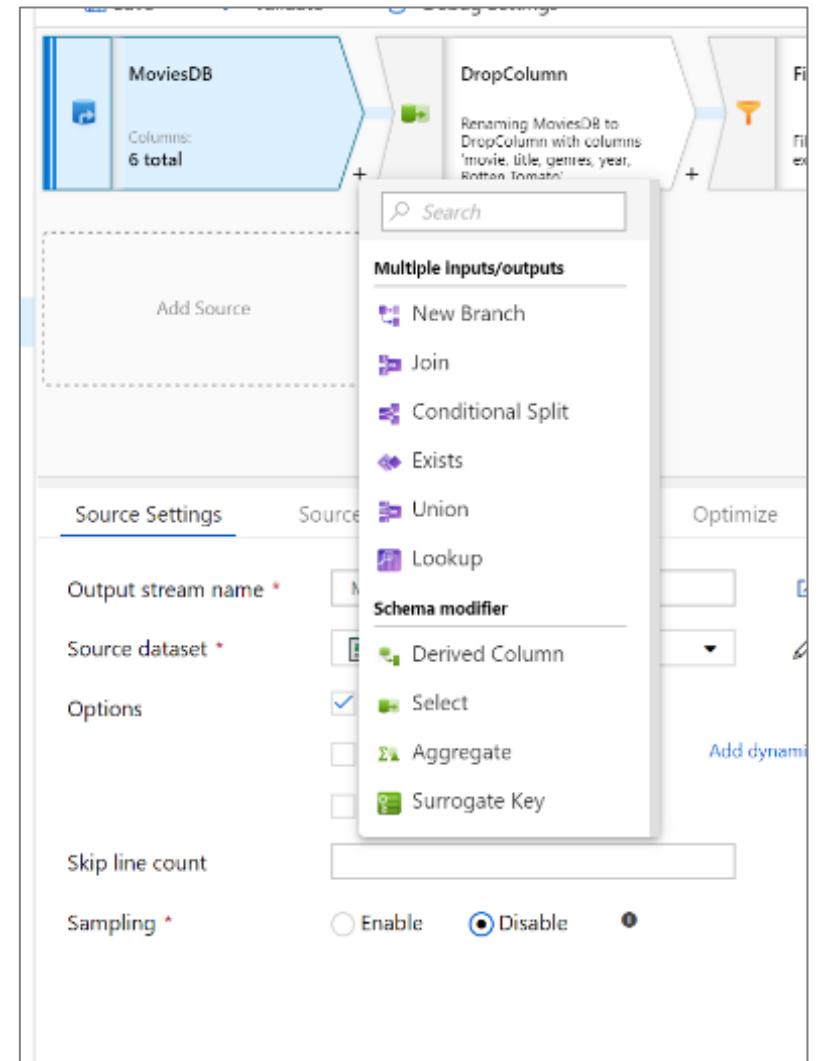
Authoring mapping data flows

Dedicated development canvas



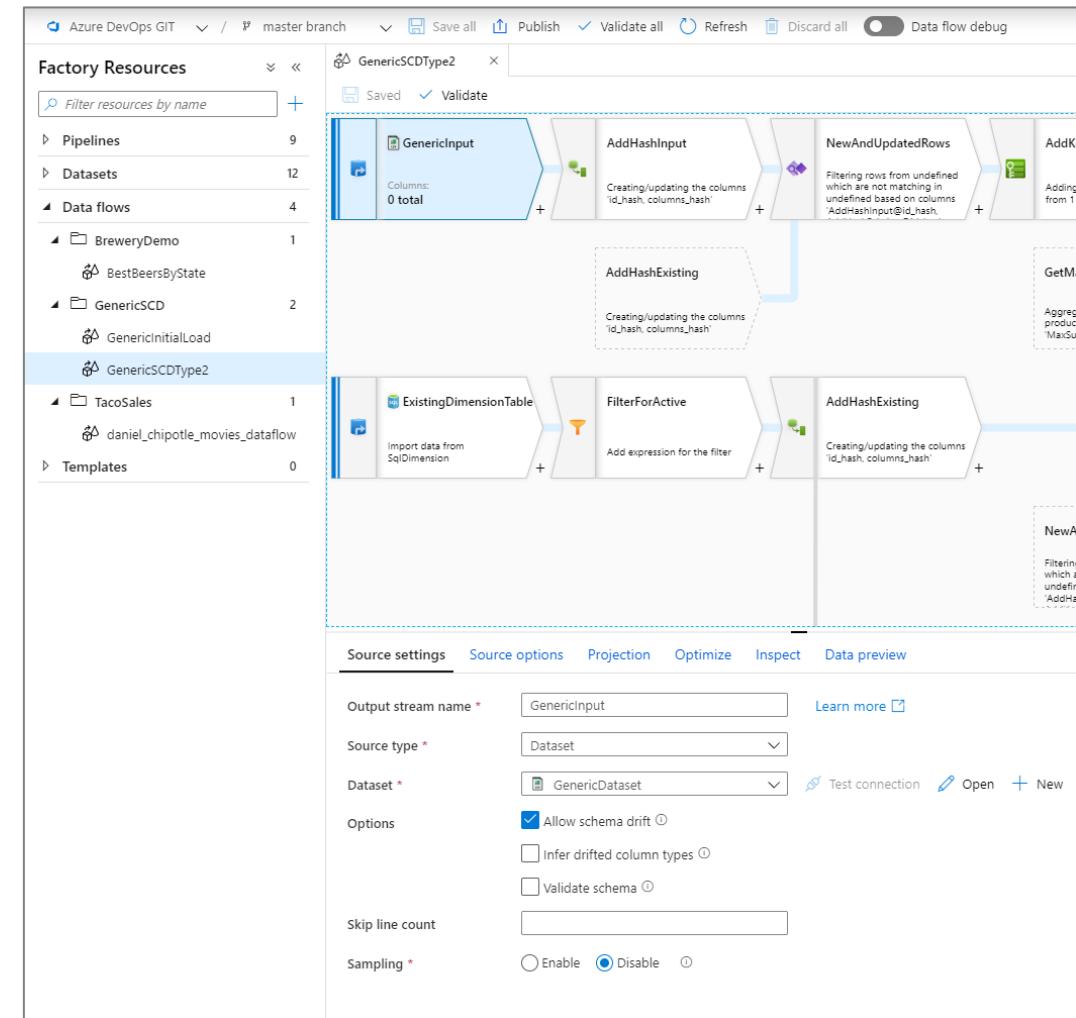
Building transformation logic

- Transformations: A 'step' in the data flow
 - Engine intelligently groups them at runtime
 - 19 currently available
- Core logic of data flow
 - Add/Remove/Alter Columns
 - Join or lookup data from datasets
 - Change number or order of rows
 - Aggregate data
 - Hierarchical to relational



Source transformation

- Define the data read by your data flow
 - Import projection vs generic
 - Schema drift
 - Connector specific properties and optimizations
- Min: 1, Max: ∞
- Define in-line or use dataset



Source: In-line vs dataset

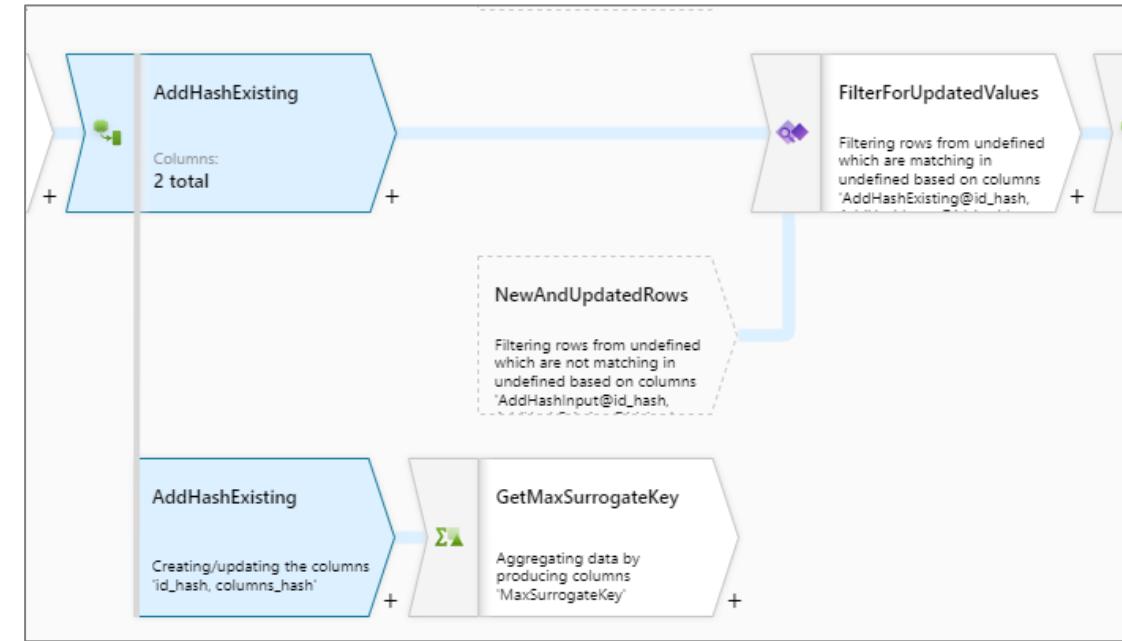
- Define all source properties within a data flow or use a separate entity to store them
- Dataset:
 - Reusable in other ADF activities such as Copy
 - Not based in Spark -> some settings overridden
- In-line
 - Useful when using flexible schemas, one-off source instances or parameterized sources
 - Do not need "dummy" dataset object
 - Based in Spark, properties native to data flow
- Most connectors only available in one

Supported connectors

- File-based data stores (ADLS Gen1/Gen2, Azure Blob Storage)
 - Parquet, JSON, DelimitedText, Excel, Avro, XML
 - In-line only: Common Data Model, Delta Lake
- SQL tables
 - Azure SQL Database
 - Azure Synapse Analytics (formerly SQL DW)
- Cosmos DB
- Coming soon: Snowflake
- **If not supported, ingest to staging area via Copy activity**
 - 90+ connectors supported natively

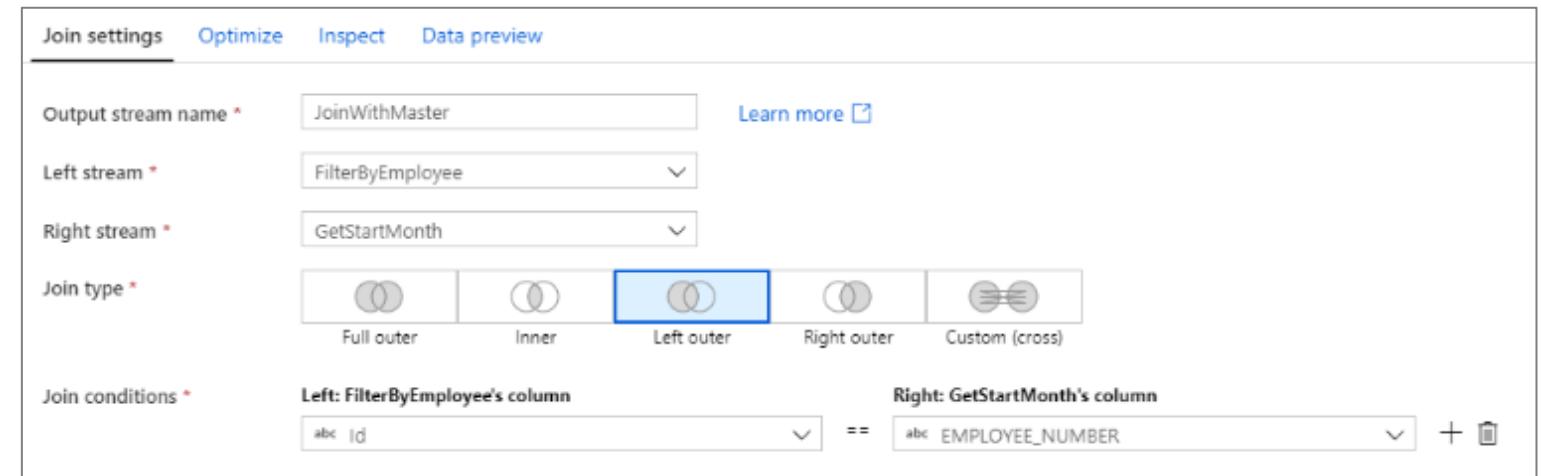
Duplicating data streams

- Duplicate data stream from any stage of your data flow
 - Select 'New branch'
- Operate on same data with different transformation requirements
 - Self-joins
 - Writing to different sinks
 - Aggregating in one branch



Joining two data streams together

- Use *Join transformation* to append columns from incoming stream to any stream in your data flow
 - Join types: full outer, inner, left outer, right outer, cross
 - SQL Join equivalent
- Match on computed columns or use non-equality conditions
- Broadcast small data streams to cache data and improve performance



Lookup transformation

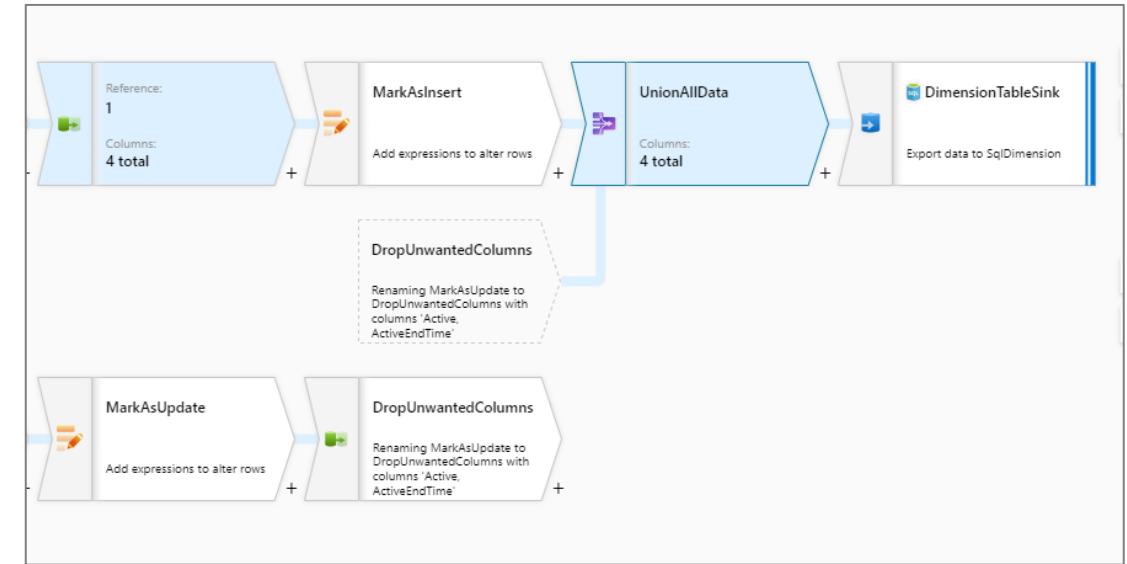
- Similar to left outer join, but with more functionality
 - All incoming rows are passed through regardless of match
- Matching conditions same as a join
- Multi or single row lookup
 - Match on all, first, last, or any row that meets join conditions
- *isMatch()* function can be used in downstream transformations to verify output

Exists transformation

- Check for existence of a value in another stream
 - SQL Exists equivalent
 - See if any row matches in a subquery, just like SQL
- Filter based on join matching conditions
- Choose **Exist** or **Not Exist** for your filter conditions
- Can specify a custom expressoin

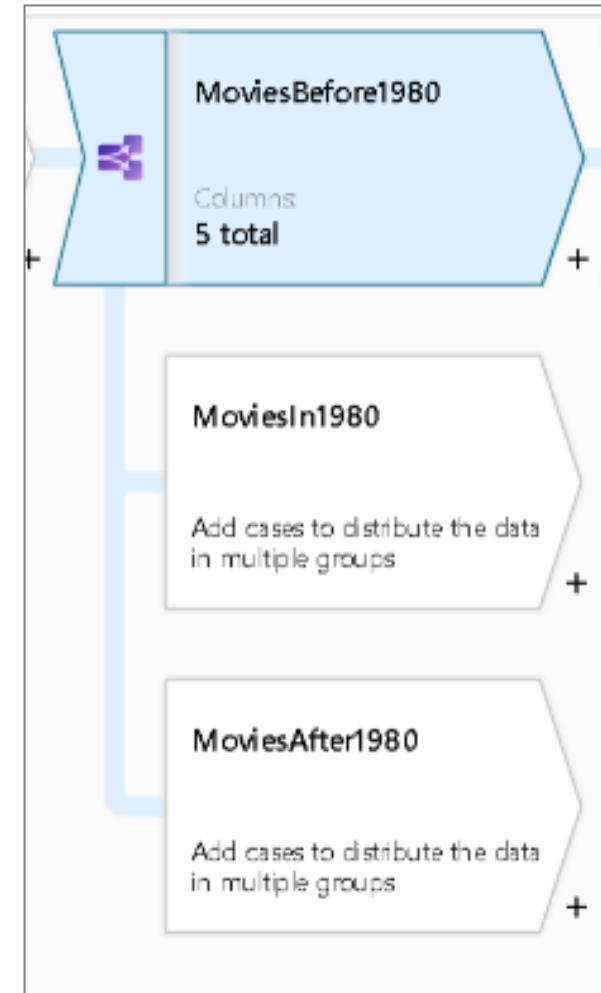
Union transformation

- Combine rows from multiple streams
- Add as many streams as needed
- Combine data based upon column name or ordinal column position
- Use cases:
 - Similar data from different connection that undergo same transformations
 - Writing multiple data streams into the same sink



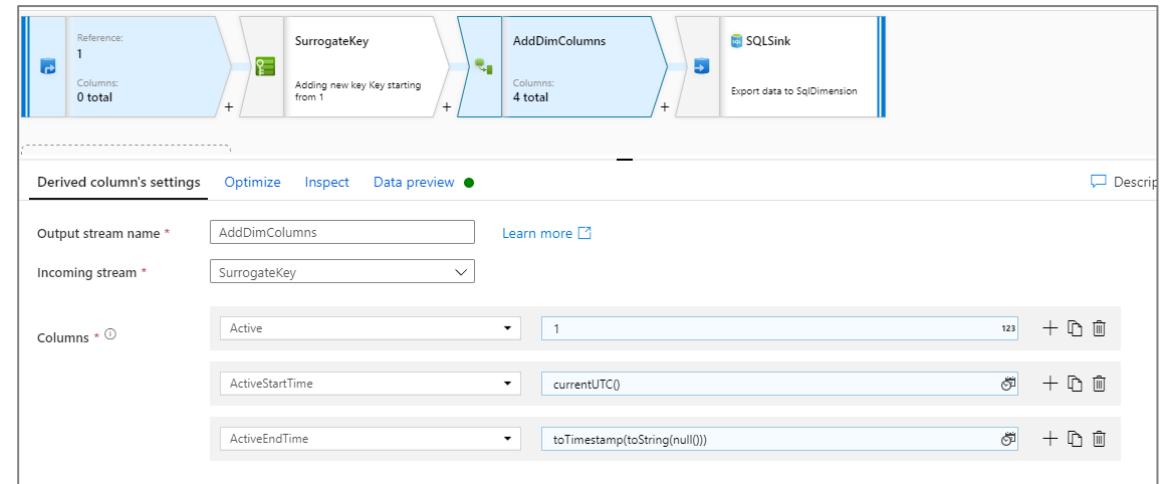
Conditional split

- Split data into separate streams based upon conditions
 - Use data flow expression language to evaluate boolean
- Use cases:
 - Sinking subset of data to different locations
 - Perform different calculations on data depending on a set of values



Derived column

- Transform data at row and column level using expression language
 - Generate new or modify existing columns
 - Build expressions using the expression builder
- Handle structured or unstructured data
 - Use column patterns to match on rules and regular expressions
 - Can be used to transform multiple columns in bulk
- Most heavily used transformation



Select transformation

- Metadata and column maintenance
 - SQL Select statement
- Alias or renames data stream and columns
- Prune unwanted or duplicate columns
 - Common after joins and lookups
- Rule-based mapping for flexible schemas, bulk mapping
 - Map hierachal columns to flat structure

The screenshot shows a user interface for a 'Select transformation'. At the top, there are tabs: 'Select settings' (selected), 'Optimize', 'Inspect', and 'Data preview'. Below these are fields for 'Output stream name' (set to 'DedupeColumns') and 'Incoming stream' (set to 'LookupBrewery'). Under 'Options', two checkboxes are checked: 'Skip duplicate input columns' and 'Skip duplicate output columns'. A button 'Auto mapping' is followed by a 'Reset' button and a '+ Add mapping' button. To the right, it says '15 mappings: 11 column(s) from the inputs left unmapped'. The main area displays a list of 'Input columns' from 'LookupBrewery's column' and their corresponding 'Name as' in the 'Name as' column. There are also '+' and '-' icons for each mapping.

LookupBrewery's column	Name as
abc beer_id	beer_id
1..2 look	look
1..2 smell	smell
1..2 taste	taste
1..2 feel	feel
1..2 overall	overall
1..2 score	score
abc Beers@name	name

Surrogate key transformation

- Generate incrementing key to use as a non-business key in your data
- To seed the starting value of your surrogate key, use derived column and a lookup from an existing table
 - Examples are in [documentation](#)
- Useful for generating keys for star schema dimension tables

Aggregate transformation

- Aggregate data into groups using aggregate function
 - Like SQL GROUP BY clause in a Select statement
 - Aggregate functions include *sum()*, *max()*, *avg()*, *first()*, *collect()*
- Choose columns to group by
 - One row for each unique group by column value
- Only columns used in transformation are in output data stream
 - Use self-join to append to existing data
- Supports pattern matching

The screenshot shows the Flink Web UI configuration for an aggregate transformation. The 'Aggregate settings' tab is selected. The 'Output stream name' is set to 'AggregateByBeer'. The 'Incoming stream' is 'ConvertTypes'. The 'Group by' field contains 'beer_id'. The 'Aggregates' section shows two aggregation rules:

- For columns matching 'type=='double':
 - Avg of 'abc': avg(\$\$) = 1.2
- For column 'reviewCount':
 - Count: count() = 121

Pivot and unpivot transformations

- Pivot row values into new columns and vice-versa
- Both are aggregate transformations that require aggregate functions
- If pivot key values not specified, all columns become drifted
 - Use map drifted quick action to add to schema quickly

Help graphic

The diagram illustrates a pivot transformation. On the left, under 'Before pivot transformation', there is a table with columns A, B, C, D, E, F. Row 1 has values 1 and 2 under columns D, E, F respectively. Row 2 has values 2 and 1 under columns D, E, F respectively. A legend indicates that grey shading represents 'GROUP BY' columns (A, B, C), green shading represents the 'PIVOT KEY' (D), and blue shading represents 'PIVOT COLUMNS' (E, F). An arrow points to the right, labeled 'After pivot transformation'. The resulting table has columns A, B, C, 1_e, 1_f, 2_e, 2_f. The values 1 and 2 are now in the 1_e and 2_e columns respectively, while the original values from E and F are moved to 1_f and 2_f.

Before pivot transformation

After pivot transformation

GROUP BY
Column: A, B, C

PIVOT KEY
Pivot column: D

PIVOT COLUMNS
Column name pattern: {Pivot key value}_{Pivot Column}
Column arrangement: Normal
Pivot Function: first(F)

Help Graphic

The diagram illustrates an unpivot transformation. On the left, under 'Before Unpivot transformation', there is a table with columns PO, Vendor, Apple, Pear. Rows #1 and #2 have values A and B under PO, and 1 and 3 under Apple, and 2 and 4 under Pear respectively. A legend indicates that grey shading represents 'GROUP BY' columns (PO, Vendor), green shading represents the 'UNPIVOT KEY' (Apple, Pear), and blue shading represents 'PIVOT COLUMNS' (Amount). An arrow points to the right, labeled 'After Unpivot transformation'. The resulting table has columns PO, Vendor, Fruit, Amount. Row #1 has values A, Apple, 1. Row #2 has values B, Apple, 2. Row #3 has values A, Pear, 3. Row #4 has values B, Pear, 4.

Before Unpivot transformation

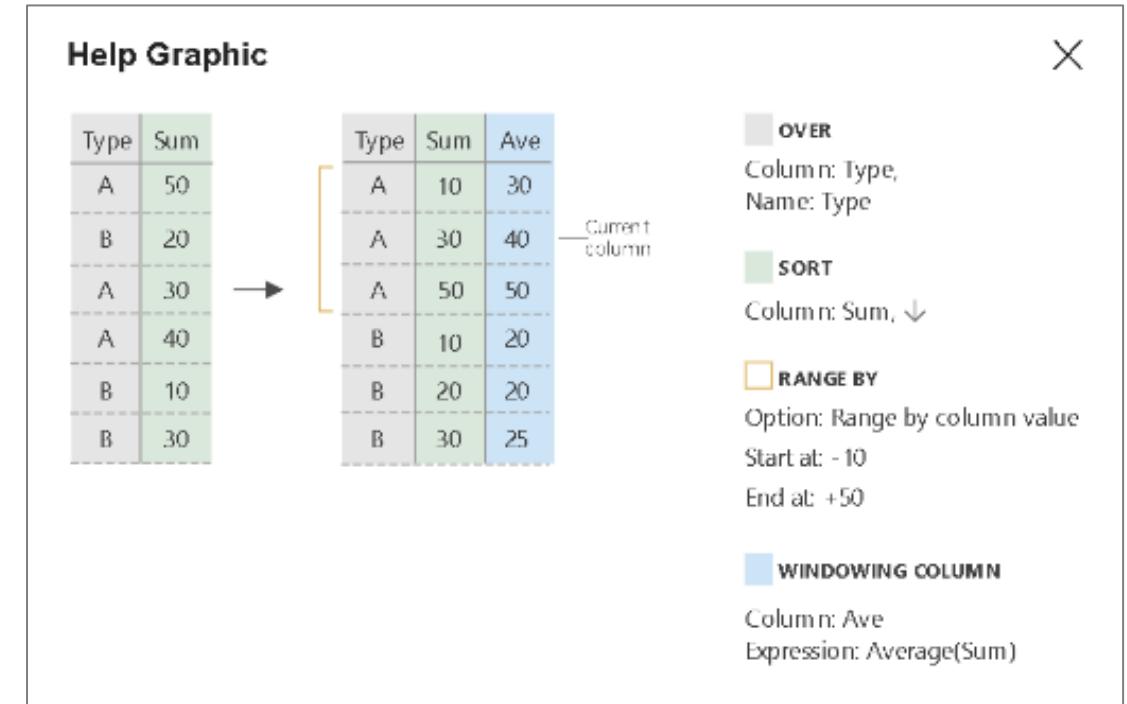
UNGROUP BY
Column: PO, Vendor

UNPIVOT KEY
Unpivot column name: Fruit
Unpivot column type: string
Option: Pick column names as values

PIVOT COLUMNS
Column arrangement: Normal
Column: Amount, Type: int

Window transformation

- Aggregate data across “windows” of data partitions
 - Used to compare a row of data against others in its ‘group’
- Group determined by group by columns, sorting conditions and range bounds
- Used for ranking rows in a group and getting lead/lag
- Sorting causes reshuffling of data
 - “Expensive” operation



Filter transformation

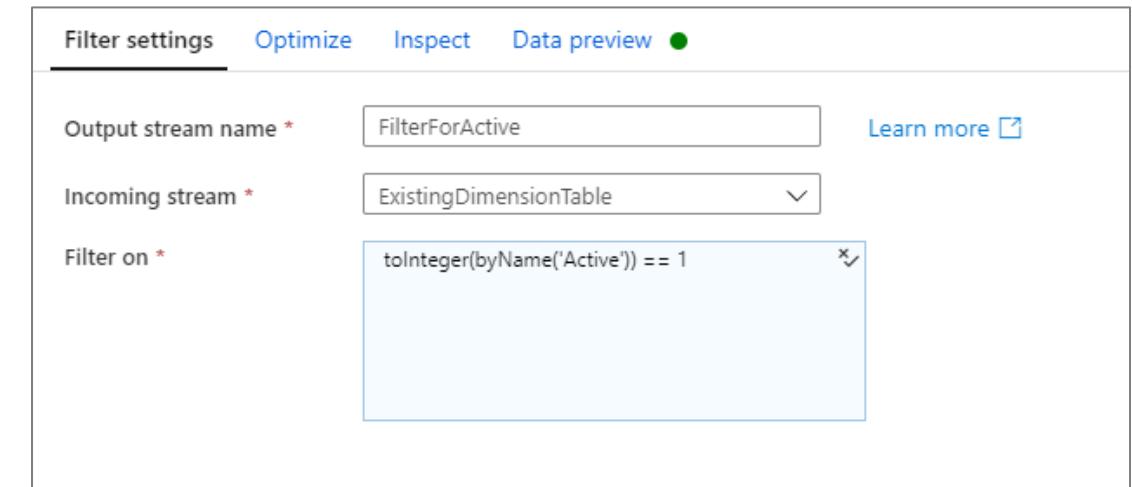
- Filter rows based upon an expression
 - Like SQL WHERE clause
- Expressions return true or false

Filter settings Optimize Inspect Data preview ●

Output stream name * [Learn more ▾](#)

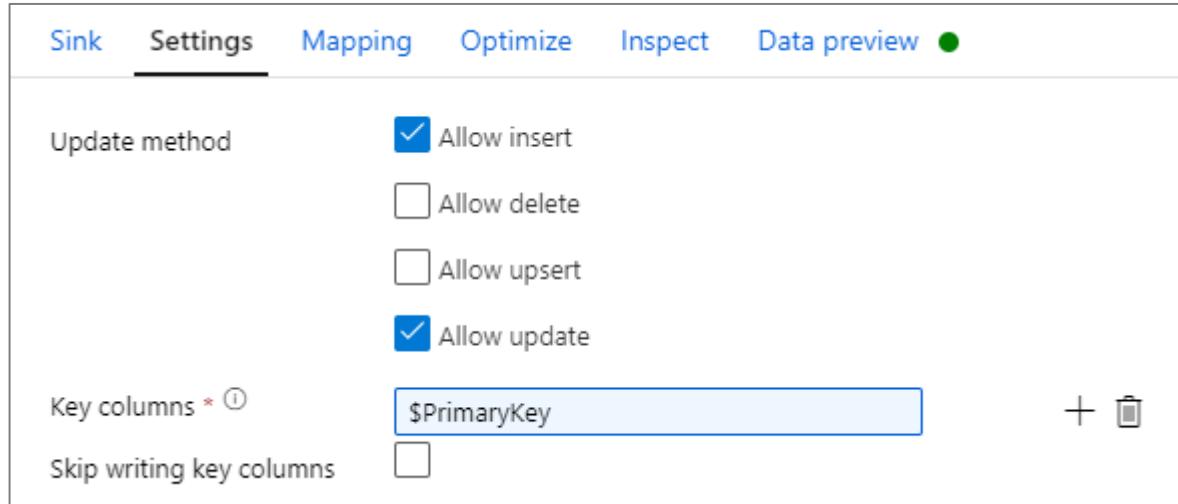
Incoming stream *

Filter on * [X](#)



Alter row transformation

- Mark rows as Insert, Update, Delete, or Upsert
 - Like SQL MERGE statement
 - Insert by default
- Define policies to update your database
 - Works with SQL DB, Synapse, Cosmos DB, and Delta Lake
- Specify allowed update methods in each sink



Flatten transformation

- Unroll array values into individual rows
 - One row per value
- Used to convert hierarchies to flat structures
- Opposite of collect() aggregate function

- Σ Aggregate
- 💡 Surrogate Key
- ᵀ Pivot
- ᵀ Unpivot
- 🕒 Window
- 📊 Rank

Formatters

- (Flatten)

- (Parse)

Row modifier

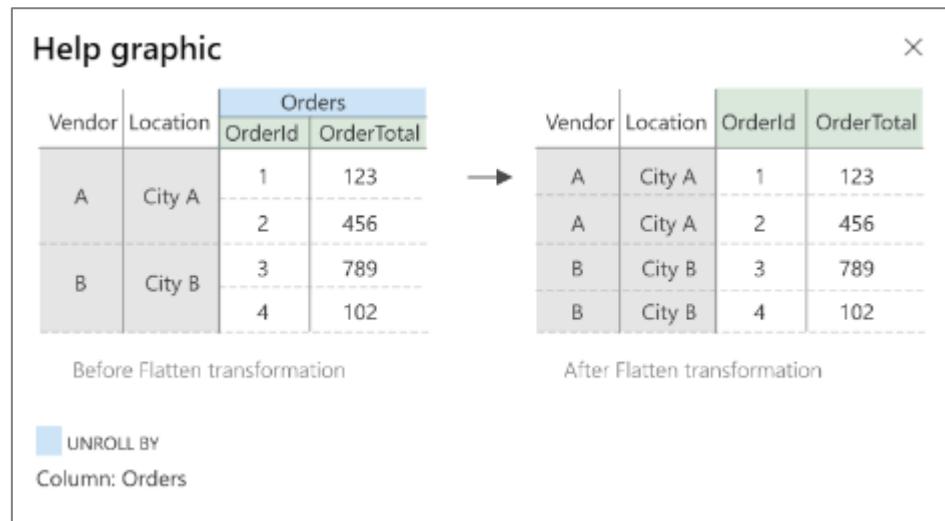
- Filter

- Sort

- Alter Row

Destination

- Sink



Rank transformation

- Rank data across an entire dataset
- Same as Rank() function in Window transformation, but scales better for a ranking an entire dataset
- For ranking of data partitions, use Window rank()
- For ranking entire dataset, use Rank transformation



Formatters



Row modifier

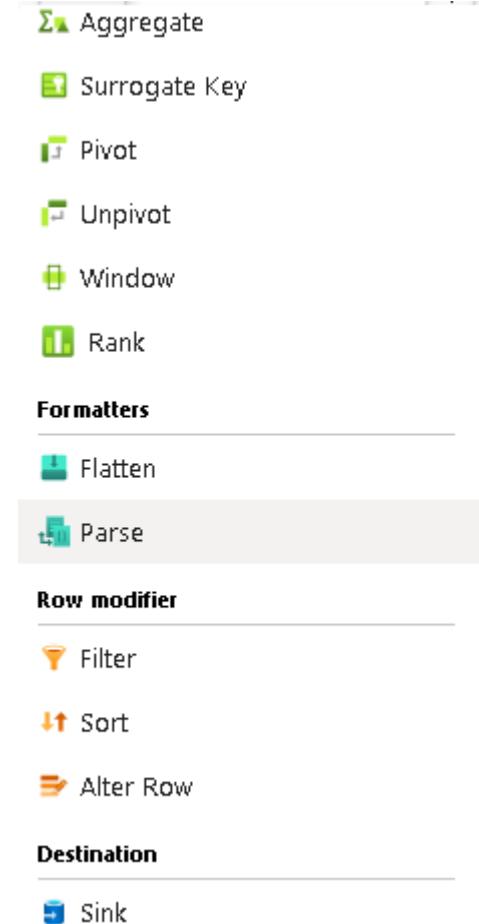


Destination



Parse transformation

- Parses string data from columns that are formatted text
- Currently supported formats: JSON, text delimited
- Ex: Turn plain text JSON strings from a source file in a single column into formatted JSON output



Sink transformation

- Define the properties for landing your data in your destination target data store
 - Define using dataset or in-line
- Can map columns similar to select transformation
 - Import schema definition from destination
- Set actions on destinations
 - Truncate table or clear folder, SQL pre/post actions, database update methods
- Choose how the written data is partitioned
 - Use current partitioning is almost always fastest
- **Note: Writing to single file can be very slow with large amounts of data**

Mapping data flow expression language

Visual expression builder

The screenshot shows the Visual expression builder interface. On the left is a sidebar with icons for Home, BusinessRules, Derived Columns, and a selected item, isPlayerGood. A blue callout box labeled "List of columns being modified" points to the sidebar. The main workspace has a header "Visual expression builder" and a sub-header "BusinessRules". It contains a "Column name *" input field with "isPlayerGood", an "Expression" input field with the placeholder code "{:playerName} is a {:playerRating} player because he averages {PTS} per game", and a "Save" button. A blue callout box labeled "Build expressions here with full auto-complete and syntax checking" points to the expression input field. Below the expression input is a toolbar with operators (+, -, *, /, ||, &&, !, ^, ==, ===, <=, !=, >, <, >=, <=, []). To the right of the expression input is a "Expression elements" section with tabs for All, Functions, Input schema, Parameters, and Locals. A blue callout box labeled "All available functions, fields, parameters ..." points to the "All" tab. Below the expression elements is a "Expression values" section with a search bar and a "Create new" dropdown. A list of available functions includes: ANY case(condition, ANY true_expression, ANY false_expression), 123 cbt(123 numeric_value), 123 ceil(123 numeric_value), ANY coalesce(ANY expression), [] collect(ANY expression), [] columnNames(abc stream name). At the bottom are buttons for "Data preview", "Save and finish", "Cancel", and "Clear contents". A blue callout box labeled "View results of your expression in the data preview pane with live, interactive results" points to the "Data preview" button.

Visual expression builder

BusinessRules

Derived Columns

+ Create new

playerName

isPlayerGood

Each column that matches type == 'do...' creates 1 column(s)

123 \$\$

List of columns being modified

Column name *

isPlayerGood

Expression

"{:playerName} is a {:playerRating} player because he averages {PTS} per game"

Save

+

-

*

/

||

&&

!

^

==

==

<=

!=

>

<

>=

<=

[]

Expression elements

All

Functions

Input schema

Parameters

Locals

Expression values

Filter by keyword

Create new

ANY case(condition, ANY true_expression, ANY false_expression)

123 cbt(123 numeric_value)

123 ceil(123 numeric_value)

ANY coalesce(ANY expression)

[] collect(ANY expression)

[] columnNames(abc stream name)

Data preview

Save and finish

Cancel

Clear contents

Build expressions here with full auto-complete and syntax checking

Expression reference documentation

All available functions, fields, parameters ...

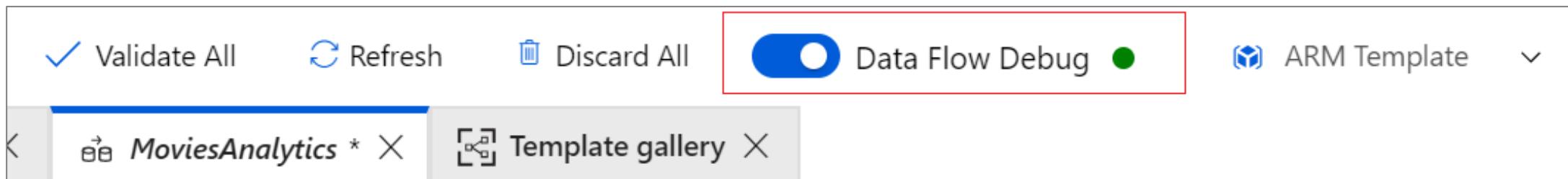
View results of your expression in the data preview pane with live, interactive results

Expression language

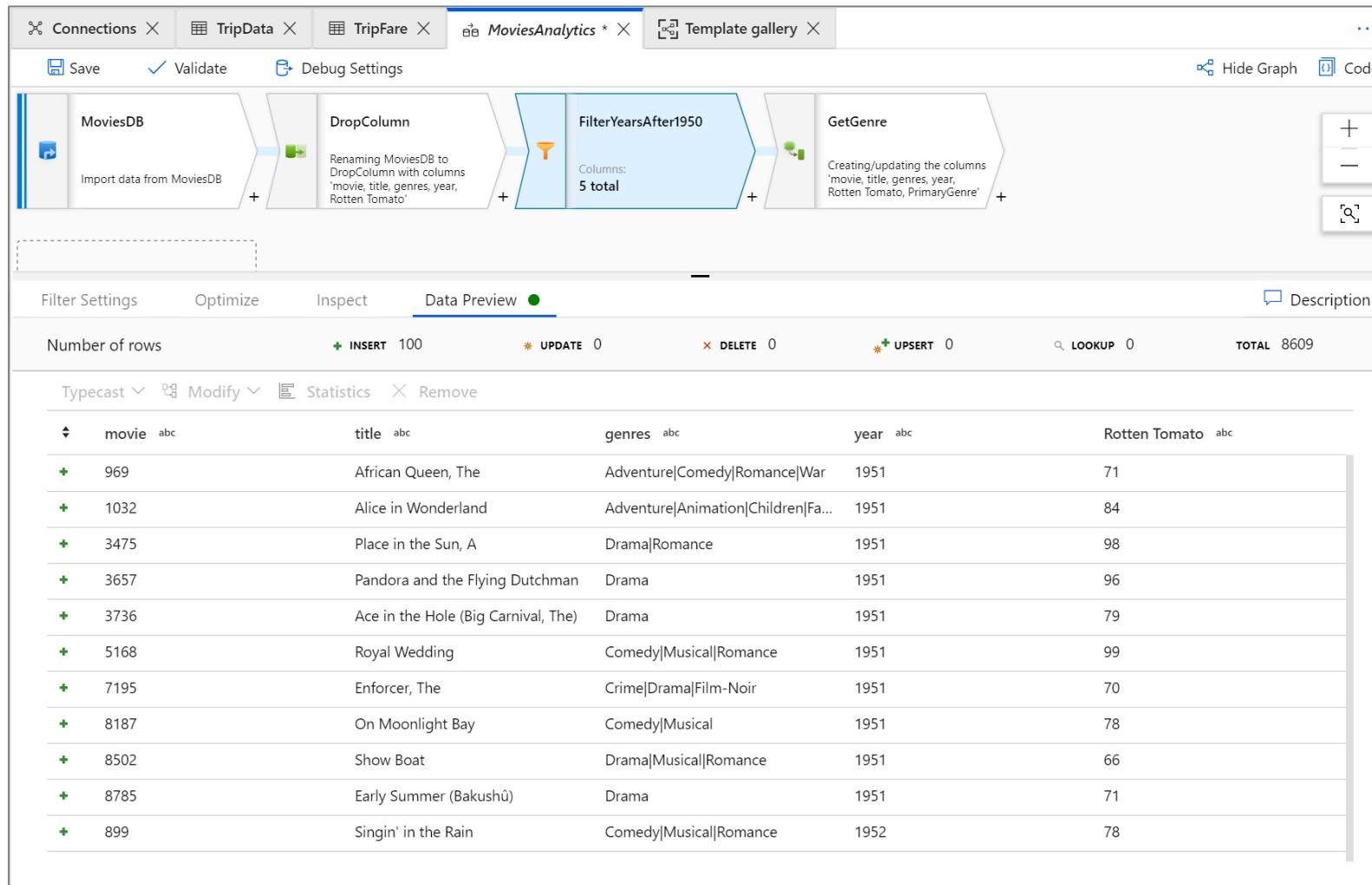
- Expressions are built using the data flow expression language
- Expressions can reference:
 - Built-in expression functions
 - Defined input schema columns
 - Data flow parameters
 - Literals
- Certain transformations have unique functions
 - Count(), sum() in Aggregate, denseRank() in Window, etc
- Evaluates to spark data types
- 100s of transformation, analytical, array, metadata, string, and mathematical functions available

Debug mode

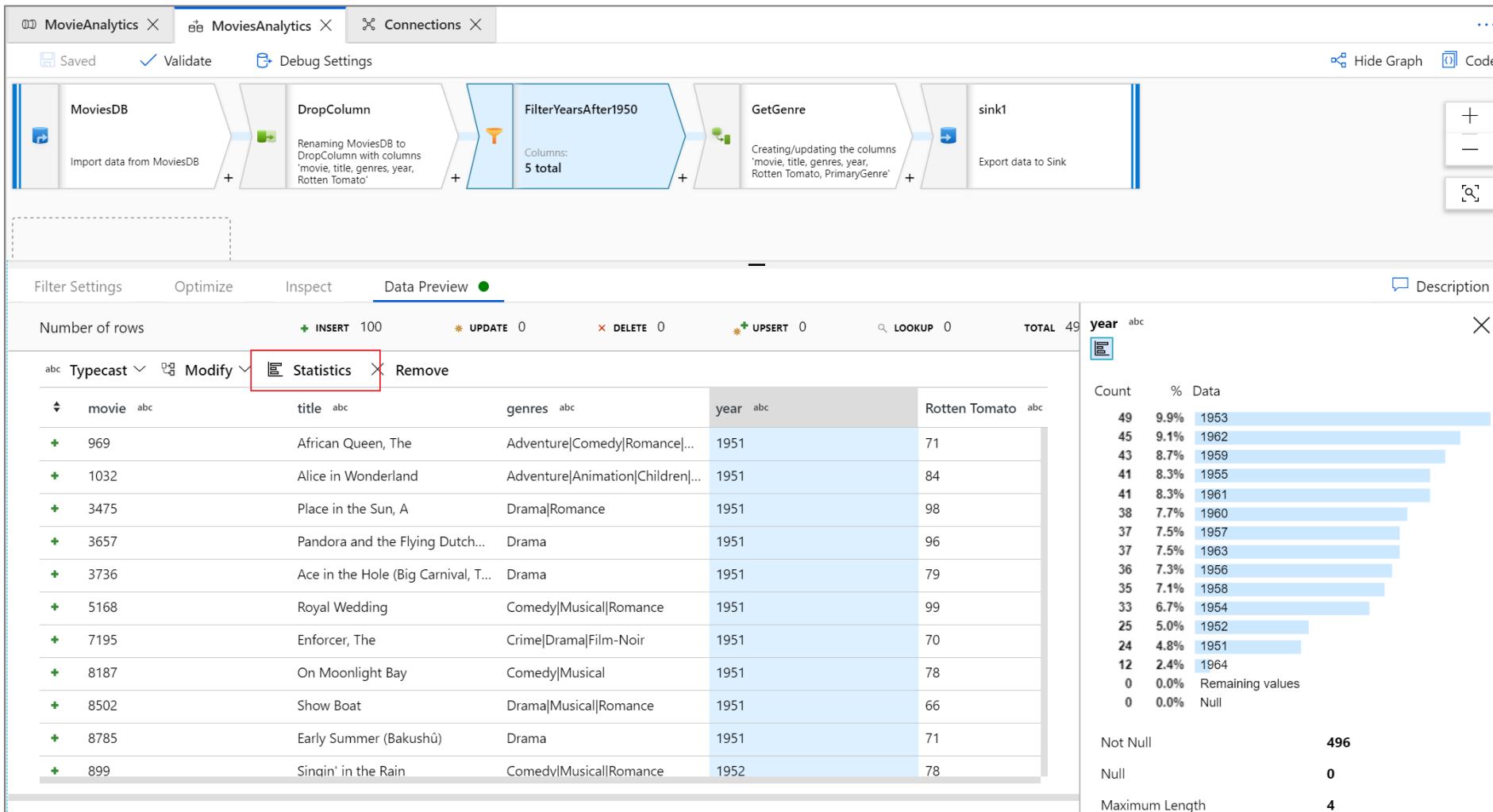
- Quickly verify logic during development on small interactive cluster
 - 4 core, 60-minute time to live
- Enables the following:
 - Get data preview snapshot at each transformation
 - Preview output of expression in expression builder
 - Run debug pipeline with no spin up
 - Import Spark projection of source schema
- **Rule of thumb:** If developing Data Flows, turn on right away
 - Initial 3-5-minute start up time



Debug mode: data preview



Debug mode: data profiling



Debug mode: expression output

Visual Expression Builder

Expression reference documentation [Save](#)

OUTPUT SCHEMA <> FUNCTIONS <> EXPRESSION FOR FIELD "PRIMARYGENRE"

abc PrimaryGenre

Filter... All Functions Input schema Parameters

abc movie
abc title

split(genres, '|')[1]

Data preview

Output: PrimaryGenre abc genres abc

Adventure	Adventure Comedy Romance War
Adventure	Adventure Animation Children Fantasy Musical
Drama	Drama Romance
Drama	Drama
Drama	Drama
Comedy	Comedy Musical Romance
Crime	Crime Drama Film-Noir
Comedy	Comedy Musical
Drama	Drama Musical Romance
Drama	Drama
Comedy	Comedy Musical Romance
Drama	Drama Romance
Drama	Drama Western
Drama	Drama

Refresh

The screenshot shows the Visual Expression Builder interface. In the top right, there's a link to 'Expression reference documentation' and a 'Save' button. The main area has three tabs: 'OUTPUT SCHEMA', 'FUNCTIONS', and 'EXPRESSION FOR FIELD "PRIMARYGENRE"'. The 'FUNCTIONS' tab is active, showing a search bar and a list of functions: 'All', 'Functions', 'Input schema', and 'Parameters'. Below this, there are two dropdown menus: 'abc movie' and 'abc title'. The 'EXPRESSION FOR FIELD "PRIMARYGENRE"' tab contains the expression 'split(genres, '|')[1]'. At the bottom, there's a 'Data preview' section with a table showing the output of the expression against a sample dataset of movies and their genres.

Parameterizing data flows

- Both dataset properties and data-flow expressions can be parameterized
 - Passed in via data flow activity
 - Can use data flow or pipeline expression language
 - Expressions can reference `$parameterName`
 - Can be literal values or column references

The screenshot shows the 'Parameters' tab of a data flow configuration interface. The tab has two buttons: '+ New' and 'Delete'. Below these buttons is a table with three columns: 'NAME', 'TYPE', and 'DEFAULT VALUE'. There are two rows in the table.

NAME	TYPE	DEFAULT VALUE
filterYear	123 integer	1950
triggerTime	abc string	

Referencing data flow parameters

Filter Settings Optimize Inspect Data Preview ●

Output stream name * FilterYearsAfter1950 [Documentation](#)

Incoming stream * DropColumn

Filter on * tolnteger(year) > \$filterYear 

Working with flexible schemas

Schema drift

- In real-world data integration solutions, source/target data stores change shape
 - Source data fields can change names
 - Number of columns can change over time
- Traditional ETL processes break when schemas drift
- Mapping data flow has built-in handling for flexible schemas
 - Patterns, rule-based mappings, byName(s) function, etc
 - Source: Read additional columns on top of what is defined in the source schema
 - Sink: Write additional columns on top of what is defined in the sink schema

Column pattern matching

- Match by name, type, stream, position

Aggregate settings Optimize Inspect Data preview ● Description

Output stream name * AggregateByBeer Learn more ↗

Incoming stream * ConvertTypes

Group by Aggregates

Grouped by: beer_id

Each column that matches `type=='double'` creates 1 column(s) ✘ + 📁 🗑

\$\$	abc	avg(\$\$)	1.2	+	✖
------	-----	-----------	-----	---	---

reviewCount ▾ count() 12l + 📁 🗑

The screenshot shows a Stream Processing configuration interface. At the top, tabs for 'Aggregate settings' (selected), 'Optimize', 'Inspect', 'Data preview', and 'Description' are visible. Below this, the 'Output stream name' is set to 'AggregateByBeer' and the 'Incoming stream' is 'ConvertTypes'. A 'Group by' section is followed by an 'Aggregates' section, which is currently active. Under 'Aggregates', it says 'Grouped by: beer_id'. A main aggregation step is shown with the condition 'Each column that matches type=='double''. This step creates one column and includes operations like 'avg(\$\$)' resulting in '1.2'. Below this, there's another row with 'reviewCount' and 'count()' operations, resulting in '12l'. Various icons for adding, deleting, and saving steps are present.

Rule-based mapping

Select settings Optimize Inspect Data preview ● Description ^

Output stream name * DropUnwantedColumns Learn more ↗

Incoming stream * MarkAsUpdate

Options Skip duplicate input columns ⓘ Skip duplicate output columns ⓘ

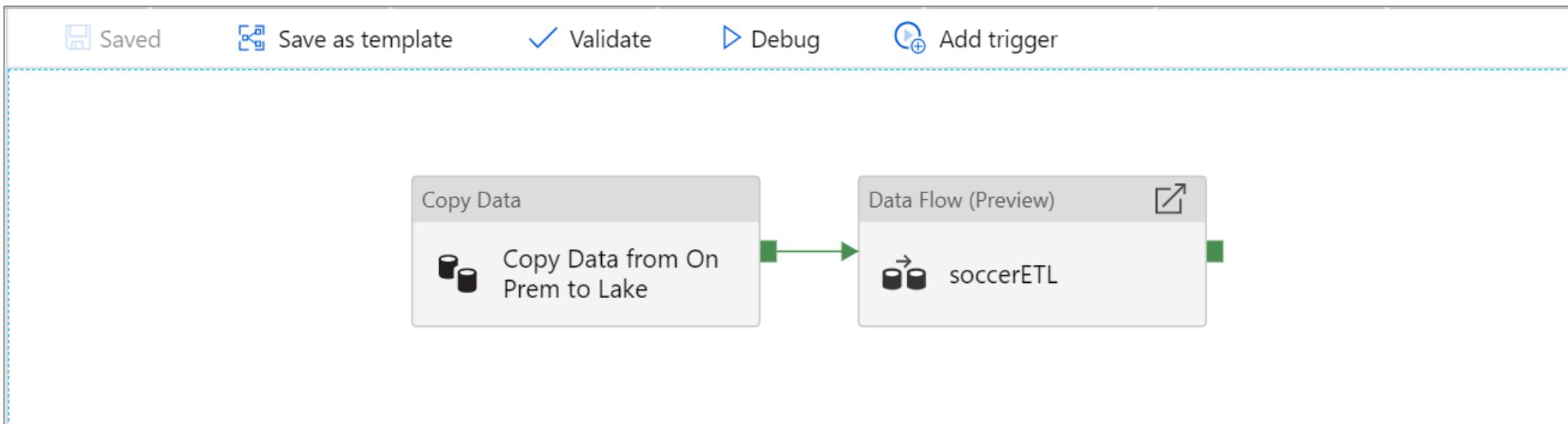
Input columns * Auto mapping ⓘ Reset Add mapping Delete 1 mappings: 2 column(s) from the inputs left unmapped ⓘ

MarkAsUpdate's column	Name as
!in(['id_hash','columns_hash','MaxSurrogateKey'],na... ✖️	\$\$ abc + ⚡

Operationalizing and monitoring data flows

Data flow activity

- Run as activity in pipeline
 - Integrated with existing ADF control flow, scheduling, orchestration, monitoring, C/ICD
- Choose which integration runtime (IR) to run on
 - # of cores, compute type, cluster time to live
- Assign parameters

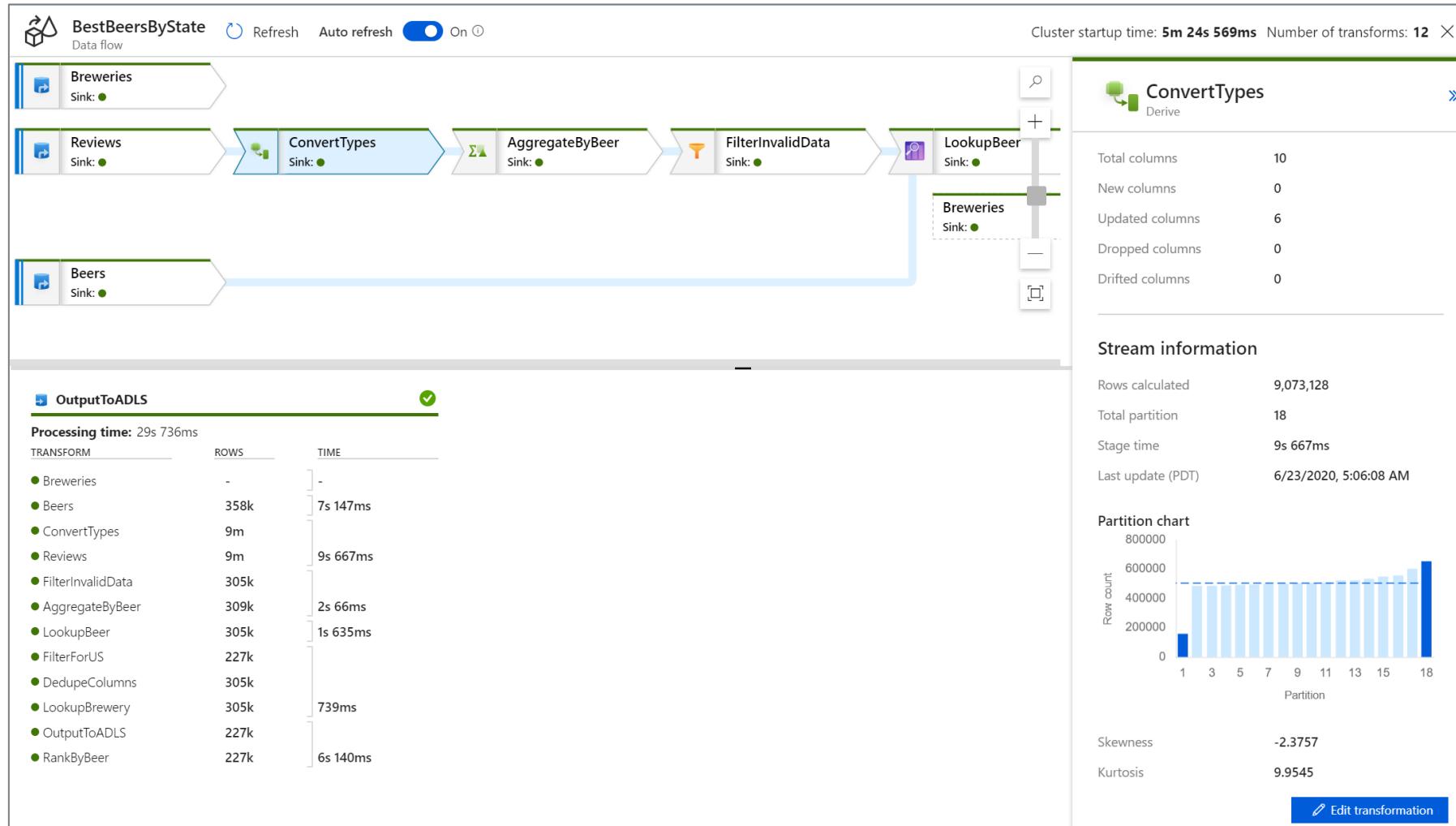


Data flow integration runtime

- Integrated with existing Azure IR
 - Choose compute type, # of cores, time to live
- **Time to live:** time a cluster is alive after last execution concludes
 - Minimal start up time for sequential data flows
- Parameterize compute type, # of cores if using *Auto Resolve*

Integration Runtimes					
Name	Actions	Type	Sub-type	Status	Region
Abhsh32CoresWorker		Azure	Public	Running	Auto Resolve
AbhshBasic		Azure	Public	Running	Auto Resolve
AutoResolveIntegrationRuntime		Azure	Public	Running	Auto Resolve
FourHourIR		Azure	Public	Running	Auto Resolve
Gen128		Azure	Public	Running	Auto Resolve

Monitoring data flows



Data flow security considerations

- All data stays inside VMs that run the Databricks cluster which are spun up JIT for each job
 - Azure Databricks attaches storage to the VMs for logging and spill-over from in-memory data frames during job operation. These storage accounts are fully encrypted and within the Microsoft tenant.
 - Each cluster is single-tenant and specific to your data and job. This cluster is not shared with any other tenant
- Data flow processes are completely ephemeral. Once a job is completed, all associated resources are destroyed
 - Both cluster and storage account are deleted
- Data transfers in data flows are protected using certificates
- Active telemetry is logged and maintained for 45 days for troubleshooting by the Azure Data Factory team

Data flow best practices and optimizations

Best practices – Lifecycle

1. Test your transformation logic using debug mode and data preview
 - Limited source size or use sample files
2. Test end-to-end pipeline logic using pipeline debug
 - Verify data is read/written correctly
 - Used as smoke test before merging your changes
3. Publish and trigger your pipelines within a Dev Factory
 - Test performance and cluster size
4. Promote pipelines to higher environments such as UAT and PROD using CI/CD
 - Increase size and scope of data as you get to higher environments

Best practices – Debug (Data Preview)

- Data Preview
 - Data preview is inside the data flow designer transformation properties
 - Uses row limits and sampling techniques to preview data from a small size of data
 - Allows you to build and validate units of logic with samples of data in real time
 - You have control over the size of the data limits under Debug Settings
 - If you wish to test with larger datasets, set a larger compute size in the Azure IR when switching on “Debug Mode”
 - Data Preview is only a snapshot of data in memory from Spark data frames. This feature does not write any data, so the sink drivers are not utilized and not tested in this mode.

Best practices – Debug (Pipeline Debug)

- Pipeline Debug
 - Click debug button to test your data flow inside of a pipeline
 - Default debug limits the execution runtime so you will want to limit data sizes
 - Sampling can be applied here as well by using the “Enable Sampling” option in each Source
 - Use the debug button option of “use activity IR” when you wish to use a job execution compute environment
 - This option is good for debugging with larger datasets. It will not have the same execution timeout limit as the default debug setting

Optimizing data flows

- Transformation order generally does not matter
 - Data flows have a Spark optimizer that reorders logic to perform as best as it can
 - Repartitioning and reshuffling data negates optimizer
- Each transformation has 'Optimize' tab to control partitioning strategies
 - Generally do not need to alter
- Altering cluster size and type has performance impact
- Four components
 1. Cluster startup time
 2. Reading from sources
 3. Transformation time
 4. Writing to sinks

Identifying bottlenecks

BestBeersByState Data flow Refresh Auto refresh On

Reviews Sink: ● ConvertTypes Sink: ● AggregateByBeer Sink: ● FilterInvalidData Sink: ● Beers Sink: ● Breweries Sink: ● Beers Sink: ●

OutputToADLS Processing time: 1m 35s 864ms

TRANSFORM	ROWS	TIME
Breweries	-	-
ConvertTypes	9m	-
Reviews	9m	1m 9s 447ms
Beers	358k	6s 379ms
FilterInvalidData	305k	-
AggregateByBeer	309k	2s 30ms
LookupBeer	305k	1s 779ms
FilterForUS	227k	-
DedupeColumns	305k	-
LookupBrewery	305k	955ms
OutputToADLS	227k	-
RankByBeer	227k	9s 790ms

1. Cluster startup time

Cluster startup time: 4m 12s 157ms Number of transforms: 12

ConvertTypes Derive

Total columns	10
New columns	0
Updated columns	6
Dropped columns	0
Drifted columns	0

Stream information

Rows calculated	9,073,128
Total partition	18
Stage time	1m 9s 447ms
Last update (PDT)	7/20/2020, 5:05:58 AM

Partition chart

4. Transformation stage time

Skewness -2.3757

Kurtosis 9.9545

- Sequential executions can lower the cluster startup time by setting a TTL in Azure IR
- Total time to process the stream from source to sink. There is also a post-processing time when you click on the Sink that will show you how much time Spark had to spend with partition and job clean-up. Write to single file and slow database connections will increase this time
- Shows you how long it took to read data from source. Optimize with different source partition strategies
- This will show you bottlenecks in your transformation logic. With larger general purpose and mem optimized IRs, most of these operations occur in memory in data frames and are usually the fastest operations in your data flow

Best practices - Sources

- When reading from file-based sources, data flow automatically partitions the data based on size
 - ~128 MB per partition, evenly distributed
 - Use current partitioning will be fastest for file-based and Synapse using PolyBase
 - Enable staging for Synapse
- For Azure SQL DB, use Source partitioning on column with high cardinality
 - Improves performance, but can saturate your source database
- Reading can be limited by the I/O of your source

Optimizing transformations

- Each transformation has its own optimize tab
 - Generally better to not alter -> reshuffling is a relatively slow process
- Reshuffling can be useful if data is very skewed
 - One node has a disproportionate amount of data
- For Joins, Exists and Lookups:
 - If you have a lot, memory optimized greatly increases performance
 - Can 'Broadcast' if the data on one side is small
 - Rule of thumb: Less than 50k rows
- Increasing integration runtime can speed up transformations
- Transformations that require reshuffling like Sort negatively impact performance

Best practices - Sinks

- SQL:
 - Disable indexes on target with pre/post SQL scripts
 - Increase SQL capacity during pipeline execution
 - Enable staging when using Synapse
- File-based sinks:
 - Use current partitioning allows Spark to create output
 - Output to single file is a very slow operation
 - Combines data into single partition
 - Often unnecessary by whoever is consuming data
 - Can set naming patterns or use data in column
 - Any reshuffling of data is slow
- Cosmos DB
 - Set throughput and batch size to meet performance requirements

Clear the folder

File name option * Default Pattern Per partition As data in column Output to single file

Output to single file * Enter file name here

Quote All

Update method Allow insert
 Allow delete
 Allow upsert
 Allow update

Table action None Recreate table Truncate table

Batch size

Pre SQL scripts

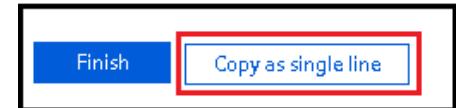
Post SQL scripts

Azure Integration Runtime

- Data Flows use JIT compute to minimize running expensive clusters when they are mostly idle
 - Generally more economical, but each cluster takes ~4 minutes to spin up
 - IR specifies what cluster type and core-count to use
 - Memory optimized is best, compute optimized doesn't generally work for production workloads
- When running Sequential jobs utilize *Time to Live* to reuse cluster between executions
 - Keeps cluster alive for TTL minutes after execution for new job to use
 - Maximum one job per cluster
- **Rule of thumb:** start small and scale up

Data flow script

Data flow script (DFS)

- DFS defines the logical intent of your data transformations
- Script is bundled and marshalled to Spark cluster as a job for execution
- DFS can be auto-generated and used for programmatic creation of data flows
- Access script behind UI via “Script” button
- Click “Copy as Single Line” to save version of script that is ready for JSON
- <https://docs.microsoft.com/en-us/azure/data-factory/data-flow-script>

Data flow script (DFS)

Movie Analytics: Split and aggregate genres , then unpivot column values to rows

```
source(output(  
    movie as string,  
    title as string,  
    genres as string,  
    year as string,  
    Rating as integer,  
    {Rotten Tomato} as string  
,  
    allowSchemaDrift: true,  
    validateSchema: false) ~> source1  
  
Aggregate1 unpivot(output(  
    avg_genre as string,  
    avgRating as double  
,  
    ungroupBy(year),  
    lateral: false,  
    ignoreNullPivots: true) ~> Unpivot1  
  
source1 aggregate(groupBy(year),  
    avg_comedy = round(avglf(split(lower(genres), '|')[1]=='comedy',Rating),1,2),  
    avg_drama = round(avglf(split(lower(genres), '|')[1]=='drama',Rating),1,2),  
    avg_action = round(avglf(split(lower(genres), '|')[1]=='action',Rating),1,2)) ~> Aggregate1  
  
Unpivot1 sort(desc(year, true)) ~> Sort1  
Sort1 sink(allowSchemaDrift: true,  
    validateSchema: false) ~> sink1
```



- Syntax: `input_name transform_type(properties) ~> stream_name`

Source projection (1)

Source properties

Unpivot transformation (3)

Aggregate Transformation (2)

Sort (4)

Sink (5)

Data flow script (DFS)

Eliminate duplicate rows and only keep distinct values

```

source(output(
    movie as integer,
    title as string,
    genres as string,
    year as string,
    Rating as string,
    {Rotten Tomato} as string
),
allowSchemaDrift: true,
validateSchema: false) ~> MoviesCSV
MoviesCSV aggregate(groupBy(movie),
    each(match(name!='movie'), $$ = first($$)) ~> DistinctRows
DistinctRows aggregate(rowcount_agg = count(1)) ~> RowCountDistinct
MoviesCSV select(mapColumn(
    movie,
    title,
    genres,
    year,
    Rating,
    {Rotten Tomato}
),
skipDuplicateMapInputs: false,
skipDuplicateMapOutputs: false) ~> OriginalData
OriginalData aggregate(rowcount_orig = count(1)) ~> RowCountOrig
DistinctRows sink(allowSchemaDrift: true,
    validateSchema: false,
    partitionBy('hash', 1)) ~> OutputDistinctData
  
```

Source projection (1)

Source properties

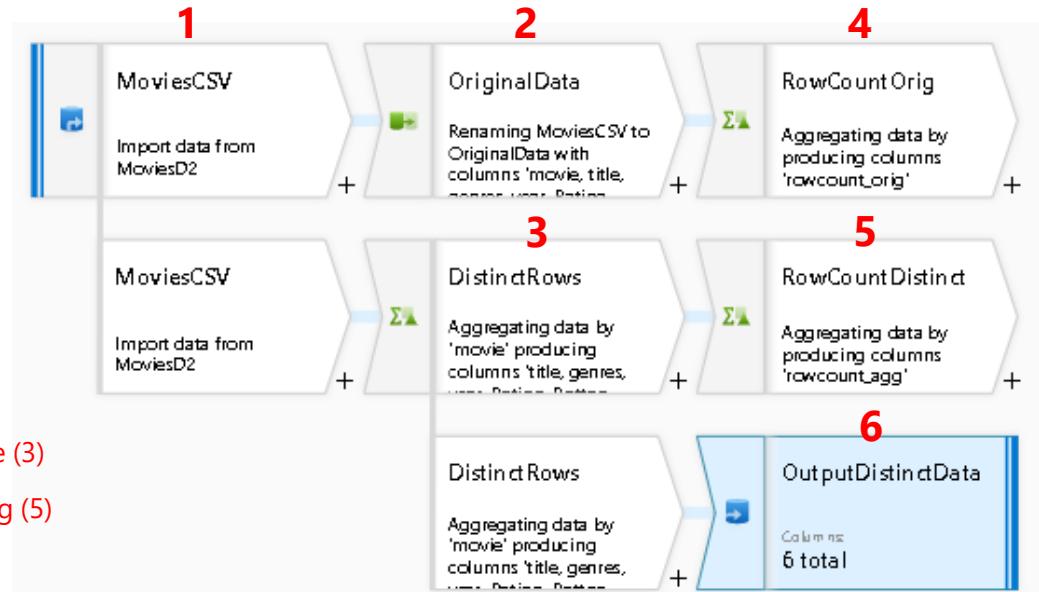
Distinct Aggregate (3)

Row Count Agg (5)

Select transformation mappings (2)

Select properties (2)

Sink transformation (6)



- `~> name_of_transform`
- New branch does not require any script element

ETL Migrations

ETL Tool Migration Overview

- Migrating from an existing large enterprise ETL installation to ADF and data flows requires adherence to a formal methodology that incorporates classic SDLC, change management, project management, and a deep understanding of your current data estate and ETL requirements.
- Successful migration projects require project plans, executive sponsorship, budget, and a dedicated team to focus on rebuilding the ETL in ADF.
- For existing on-prem ETL estates, it is very important to learn basics of Cloud, Azure, and ADF generally before taking this Data Flows training.

Sponsorship

- Migration to another platform requires buy-in and sponsorship at the executive level and budget for this effort
- A team must be dedicated to this work
- That team can be internal or external, but a commitment to fund milestone-based project work is necessary

Discovery

- How many mappings are in the current ETL tool?
- Understand the high-level objectives, business rules, and logic of each mapping
- List out all required data source, sinks, and sample data for testing
- Capture current schedules and data SLAs
- Capture all required developer, user, and tester credentials as well as required credentials to access data sources and destinations
- Think about how you want to categorize the ETL mappings for migration
 - By complexity
 - By workflow
 - By business unit

Training

- On-prem to Cloud, Azure general training, ADF general training, Data Flows training
- A general understanding of the difference between legacy client/server on-prem ETL architectures and cloud-based Big Data processing is required
- ADF and Data Flows execute on Spark, so learn the fundamentals of the difference between row-by-row processing on a local server and batch/distributed computing on Spark in the Cloud

Execution

- Start with the top 10 mission-critical ETL mappings and list out the primary logical goals and steps achieved in each
- Use sample data and debug each scenario as new pipelines and data flows in ADF
- UAT each of those 10 mappings in ADF using sample data
- Lay out end-to-end project plan for remaining mapping migrations
- Plan the remainder of the project into quarterly calendar milestones
- Expect each phase to take around 3 months
- Majority of large existing ETL infrastructure modernization migrations take 12-18 months to complete

ETL System Integrator Partners

- Bitwise Global
 - <https://www.bitwiseglobal.com/webinars/automated-etl-conversion-to-adf-for-accelerated-data-warehouse-migration-on-azure/>
- Next Pathway
 - <https://blog.nextpathway.com/next-pathway-adds-ground-breaking-capability-to-translate-informatica-and-datastage-etl-pipelines-to-the-cloud-in-latest-version-of-shift>

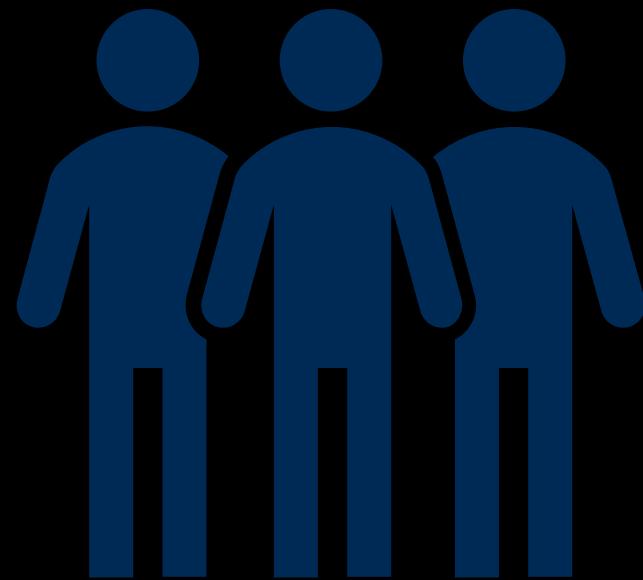


Azure HPC/Batch

The left side of the image features a dark blue background with numerous thin, glowing blue lines radiating from a central point on the left. Interspersed among these lines are small, bright blue dots of varying sizes, creating a sense of depth and motion.

HPC and the Cloud

HPC in the cloud: Market snapshot



47%

of HPC customers are testing workloads in the cloud

Drivers

- Faster time-to-solution
- Integration of AI
- Broad awareness

Challenges

- Complex workflows
- Software licensing
- Long sales cycle
- Perceived costs

New workloads

- Predictive analytics
- Autonomous driving
- Deep learning

By 2022, global HPC market expected to grow to

~\$45B

at a compound annual growth rate of 7%

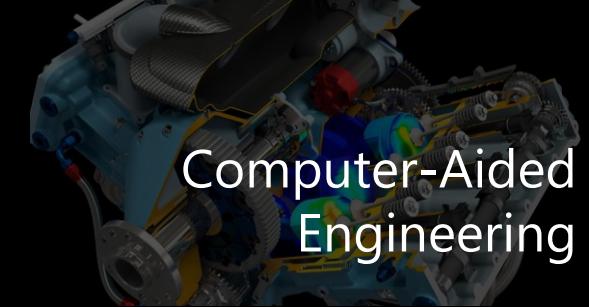
Innovation without boundaries



Discrete
Manufacturing



Genomics



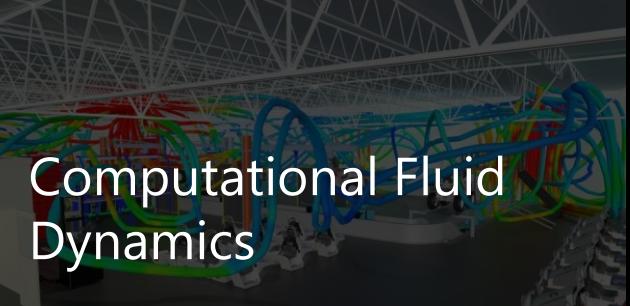
Computer-Aided
Engineering



Media
Rendering



Digital
Twins



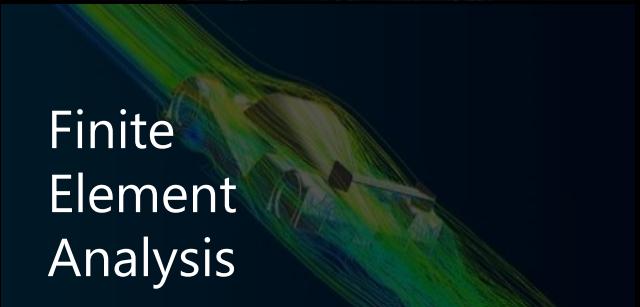
Computational Fluid
Dynamics



Climate
Modeling



Healthcare



Finite
Element
Analysis



Oil & Gas

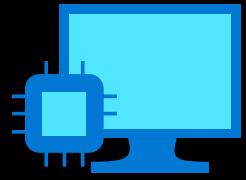


Design
Simulation



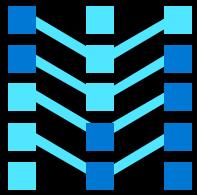
Financial
Services

~~HPC in the Cloud~~ A cloud built for HPC



Purpose-built HPC

A full range of CPU and GPU capabilities that help applications scale to 80K+ cores



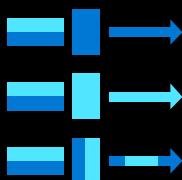
Fast, Secure Networking

Fast InfiniBand interconnects as well as edge-to-cloud connectivity



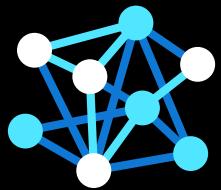
High Performing Storage

A range of storage capabilities to support simple-to-complex storage needs



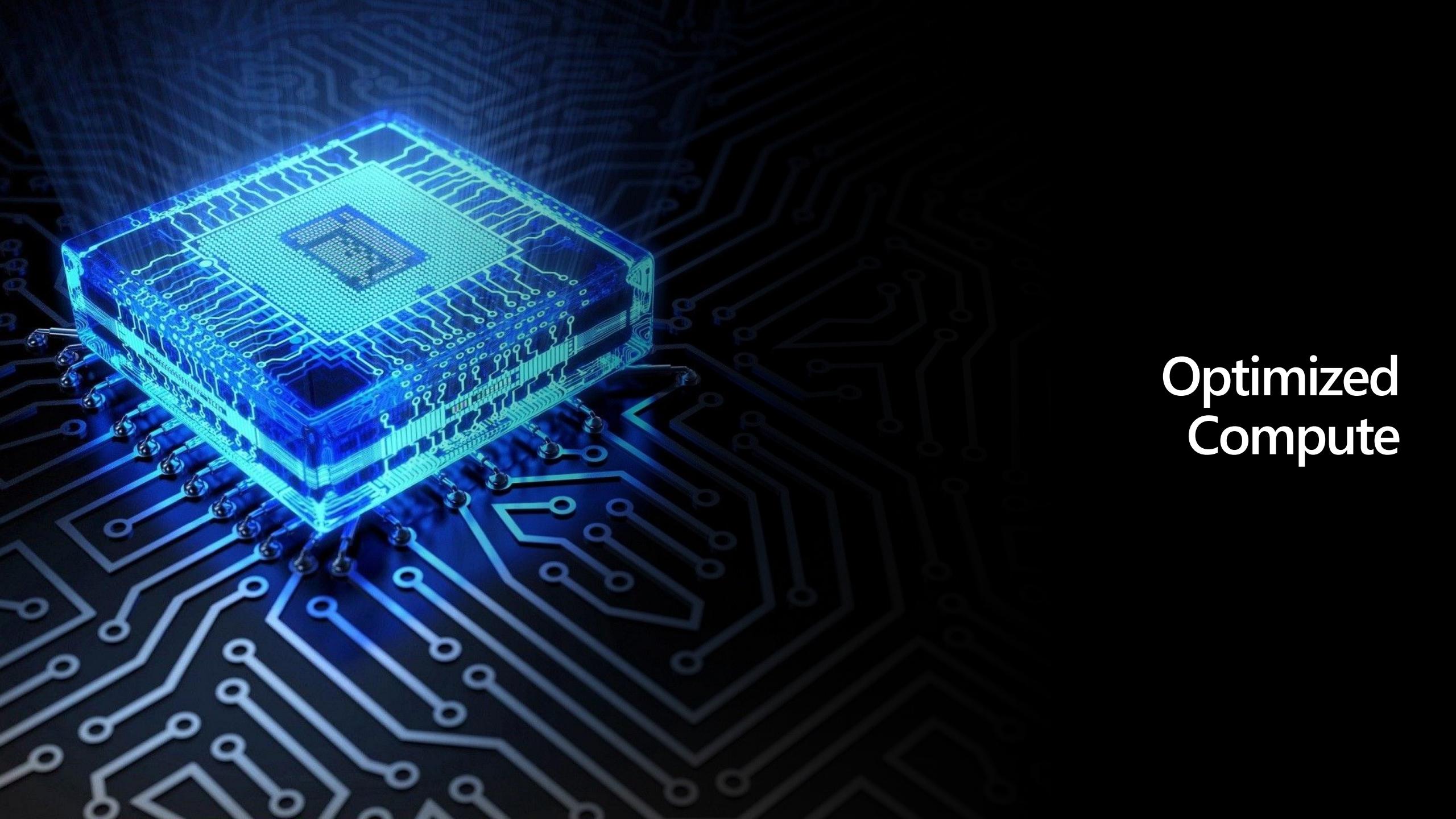
Workload Orchestration

End-to-end workflow agility using known, familiar tools & processes



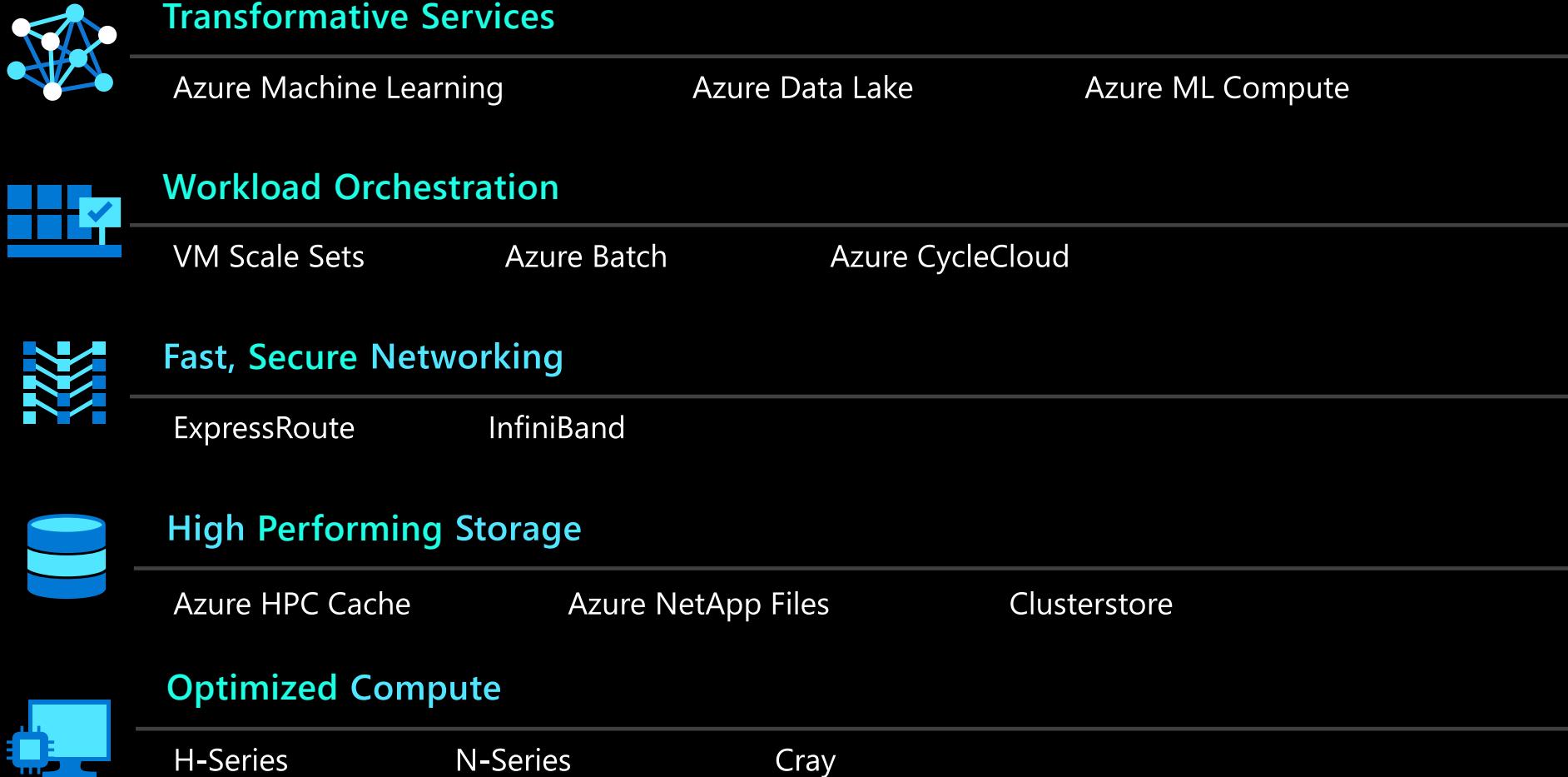
Intelligence Services

AI, machine learning, and deep learning at supercomputer scale

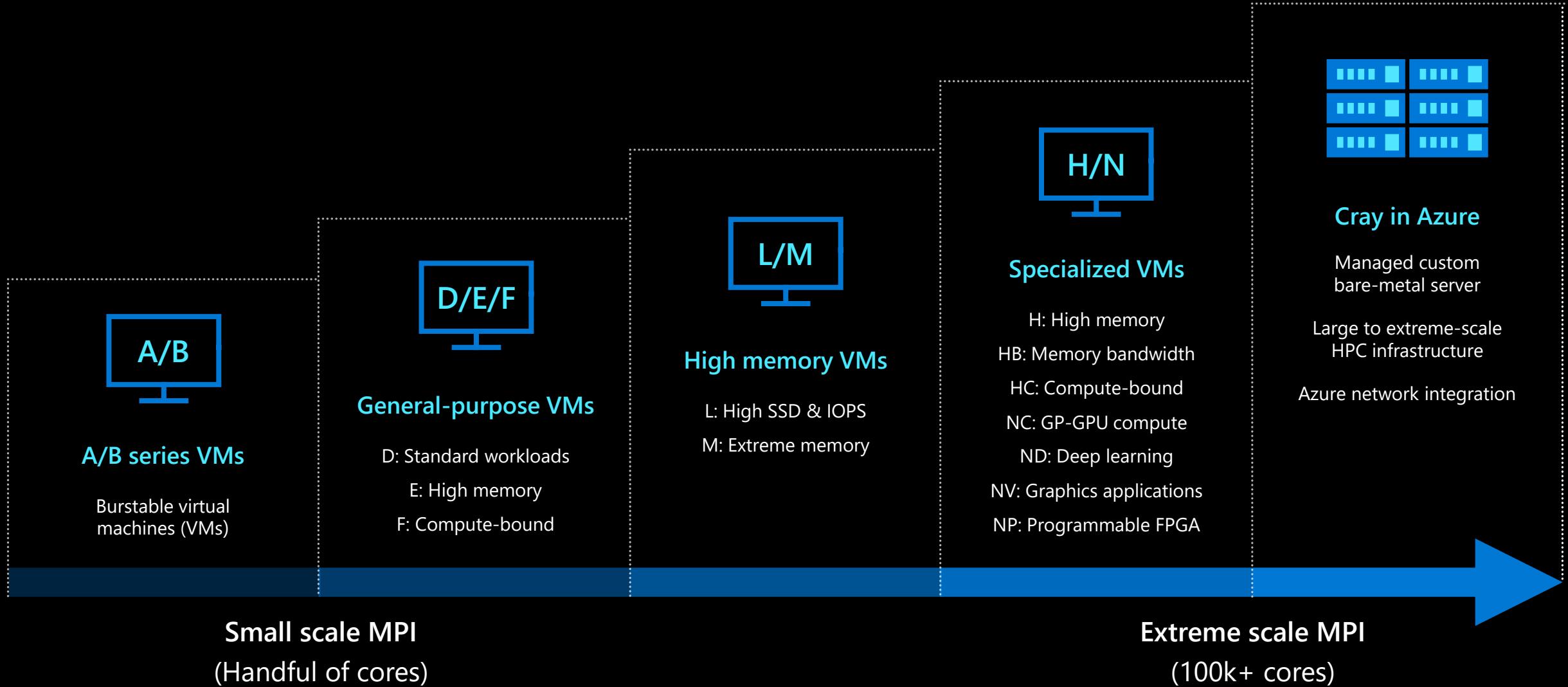
A glowing blue 3D rendering of a computer processor (CPU) is centered against a dark blue background featuring a complex, glowing blue circuit board pattern.

Optimized
Compute

HPC Resource Stack on Azure



Solve any HPC, AI workload—at any scale



High-Performance Computing VMs (H-Series)

					
Workload Optimized	Memory Bandwidth	Memory Bandwidth	Memory Bandwidth	Dense Compute	Large-Memory HPC
CPU	AMD EPYC 3 rd Gen w/ 3D V-Cache "Milan-X"	AMD EPYC 2 nd Gen "Rome"	AMD EPYC 1 st Gen "Naples"	Intel Xeon Platinum 1 st Gen "Skylake"	Intel Xeon E5 v3 "Haswell"
Cores/VM	120	120	60	44	16
TeraFLOPS/VM (FP64)	4T	4 TF	.9 TF	2.6 TF	.7 TF
Memory Bandwidth	350 GB/sec	350 GB/s	263 GB/s	191 GB/s	80 GB/s
Memory	448 GB total	4 GB/core, 480 total	4 GB/core, 240 total	8 GB/core, 352 GB	14 GB/core, 224 GB
Local Disk	2 * 960 GB NVMe	960 GB NVMe	700 GB NVMe		2 TB SATA
InfiniBand	200 Gb HDR	200 Gb HDR	100 Gb EDR		56 Gb FDR
Network	50 GbE	50 GbE	50 GbE		16 GbE

*All cores, non-AVX, peak Boost/Turbo frequencies

NVIDIA in Azure

Broad offering with the N-series

GPUs from M60 to V100 to support ML & DL training & inference, HPC across industries and Graphics workloads

Integration with key Microsoft solutions

like AML and ONNX for end-to-end acceleration for the most complex workloads, reducing burden of producing world-class solutions

Versatility across applications and frameworks

to support engineering on one platform and solution for the entire team, regardless of the preferences of particular engineers

All NVIDIA acceleration software housed in NGC

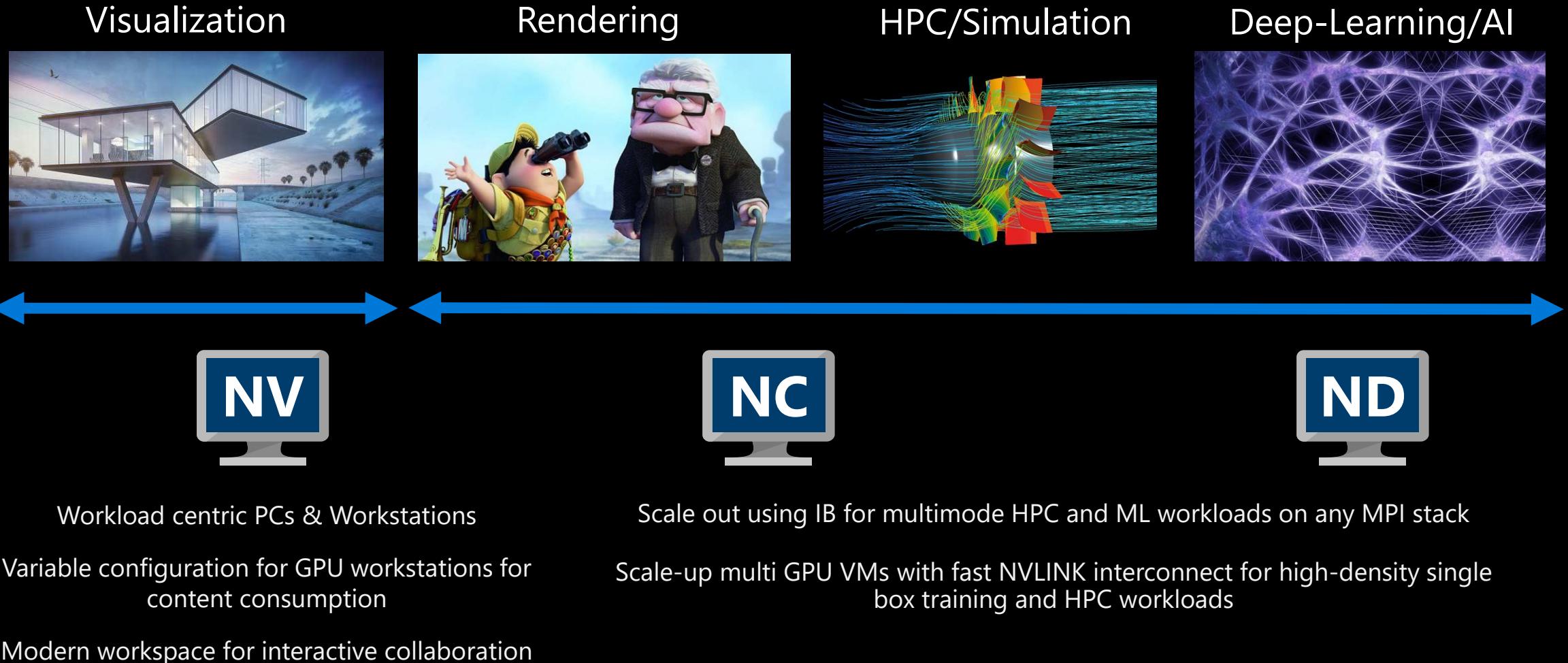
via Azure marketplace to ensure users have the easiest experience packaging what they need to achieve best performance in no time.



Optional InfiniBand interconnect

enables scale-up performance

GPU Computing



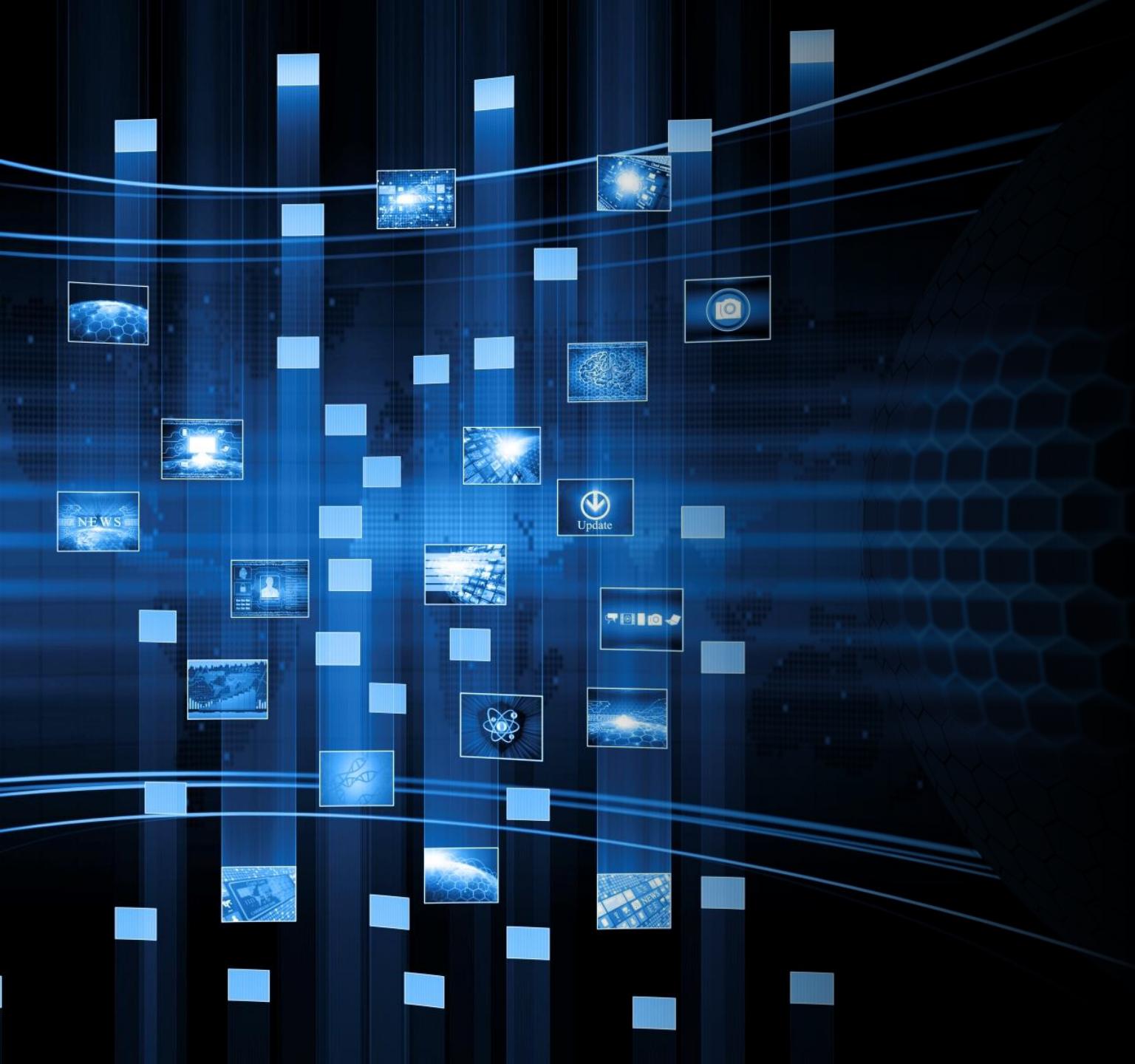
GPUs for Deep Learning (ND Series)

	 ND	 NDv2	 ND A100 v4	 NDM A100 v4
CPU Cores	6,12,24	40	96	96
GPU	1x, 2x, or 4x P40 GPUs	8x V100 32 GB (NVLink) GPUs	8x A100 40 GB GPUs	8x A100 80 GB GPUs
Memory	12/224/448 GB	672 GB	900 GB	1900 GB
Local Disk	736/1474/2948 GiB SSD	2948 GiB SSD	6 TB SSD	6.4 TB SSD
Network	Azure Network + InfiniBand EDR	Azure Network + InfiniBand EDR + NVLink GPU Interconnect	Azure Network + InfiniBand EDR + NVLink GPU Interconnect	Azure Network + InfiniBand EDR + NVLink GPU Interconnect



deep learning

Workload Orchestration



Azure Batch

Batch pools

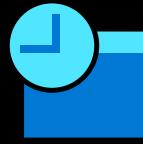
- Configure and create VMs to cater for any scale: tens to thousands
- Automatically scale the number of VMs to maximize utilization
- Easy low-priority and VM sizing, suited to your application

Batch jobs and tasks

- Task is a unit of execution; task = application command line (EXE, BAT, CMD, PS1, etc.)
- Jobs are created and tasks are submitted to a pool. Next, tasks are queued and assigned to VMs
- Any application, any execution time; run applications unchanged
- Automatic detection and retry of frozen or failing tasks



Azure Batch Capabilities



Job scheduling

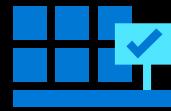
Supports both EP & MPI jobs

Run > 1 task in parallel per node

Detect and retry failed tasks

Task dependencies

Job prep and cleanup tasks



Rich app management

Get apps from blobs, Batch app packages, package managers, custom VM images

Docker container images



Choice of VMs

Windows or Linux

Standard or custom images

Windows pool can use AHUB

Use low-priority VMs



Monitoring

VM monitoring and auto-recover

Metrics and logs available via Portal and API



Developer access

.NET, Java, Node.js, Python, REST

PowerShell, x-plat Azure CLI

Azure Portal, Batch Labs x-plat client UI

What is Azure Batch?

Service / Solution

...manage and

...get a

...move ta
ou

Now how do I...

isks?

...scale up and down?

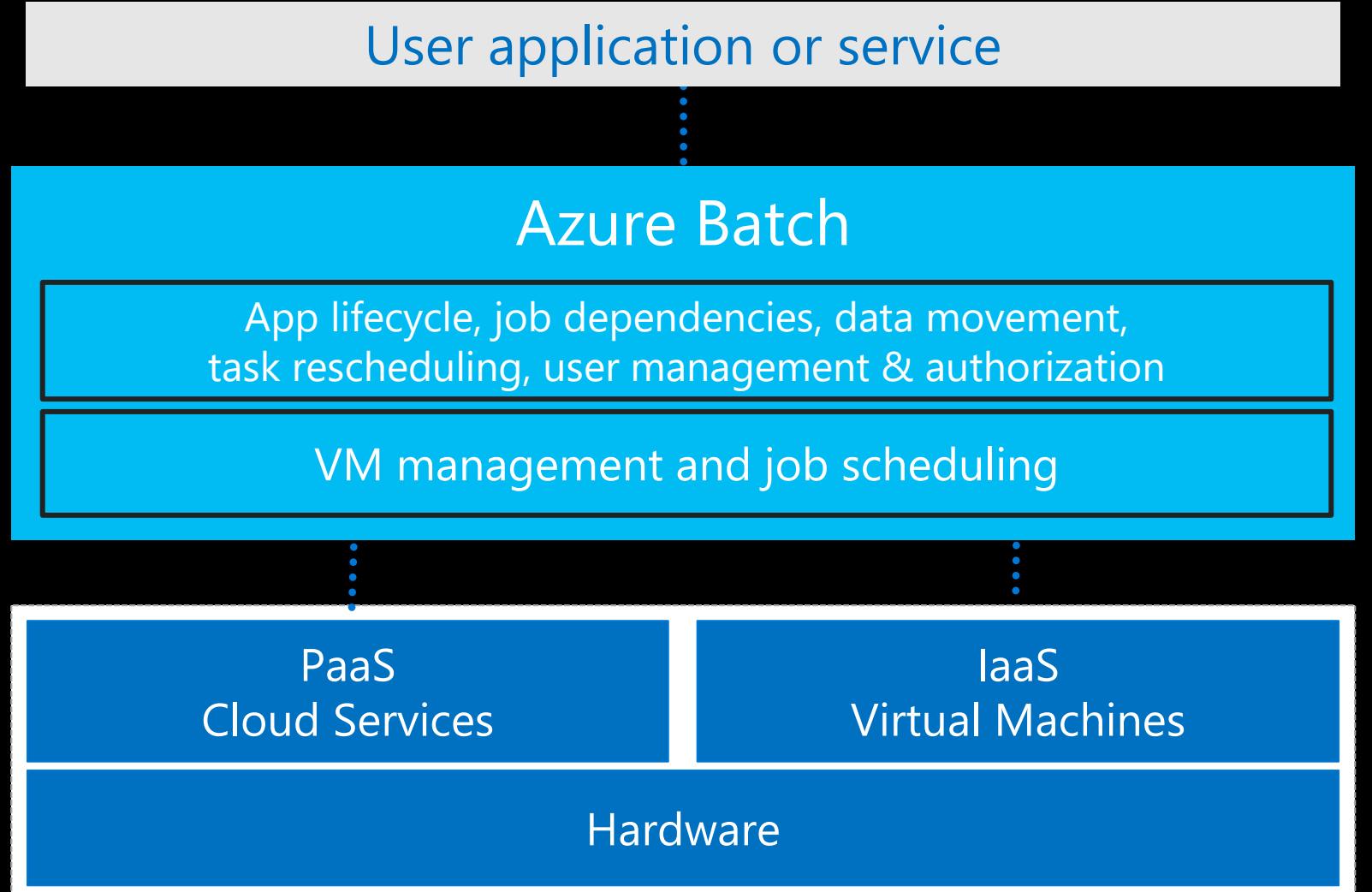
task?

PaaS
Cloud Services

IaaS
VM / VMSS

Hardware

What is Azure Batch?



Provided by the cloud platform

- Don't worry about the "plumbing"
- Focus on the workload/app
- Access higher-level capabilities
- Minimize the required cloud or Azure experience

Azure Batch Capabilities



Job scheduling

- Supports both embarrassingly parallel and tightly-coupled MPI jobs
- Run > 1 task in parallel per node
- Detect and retry failed tasks
- Can set max execution time for jobs and tasks
- Task dependencies
- Job prep and cleanup tasks



Rich app management

- Get apps from blobs, Batch app packages, package managers, custom VM images
- Docker container images



Choice of VMs

- Windows or Linux
- Standard or custom images
- Windows pool can use AHUB
- Use low-priority VMs



Monitoring

- VM monitoring and auto-recover
- Metrics and logs available via Portal and API



Access via APIs, CLIs, and UIs

- .NET, Java, Node.js, Python, REST
- PowerShell, Azure CLI
- Azure Portal, Batch Explorer client UI

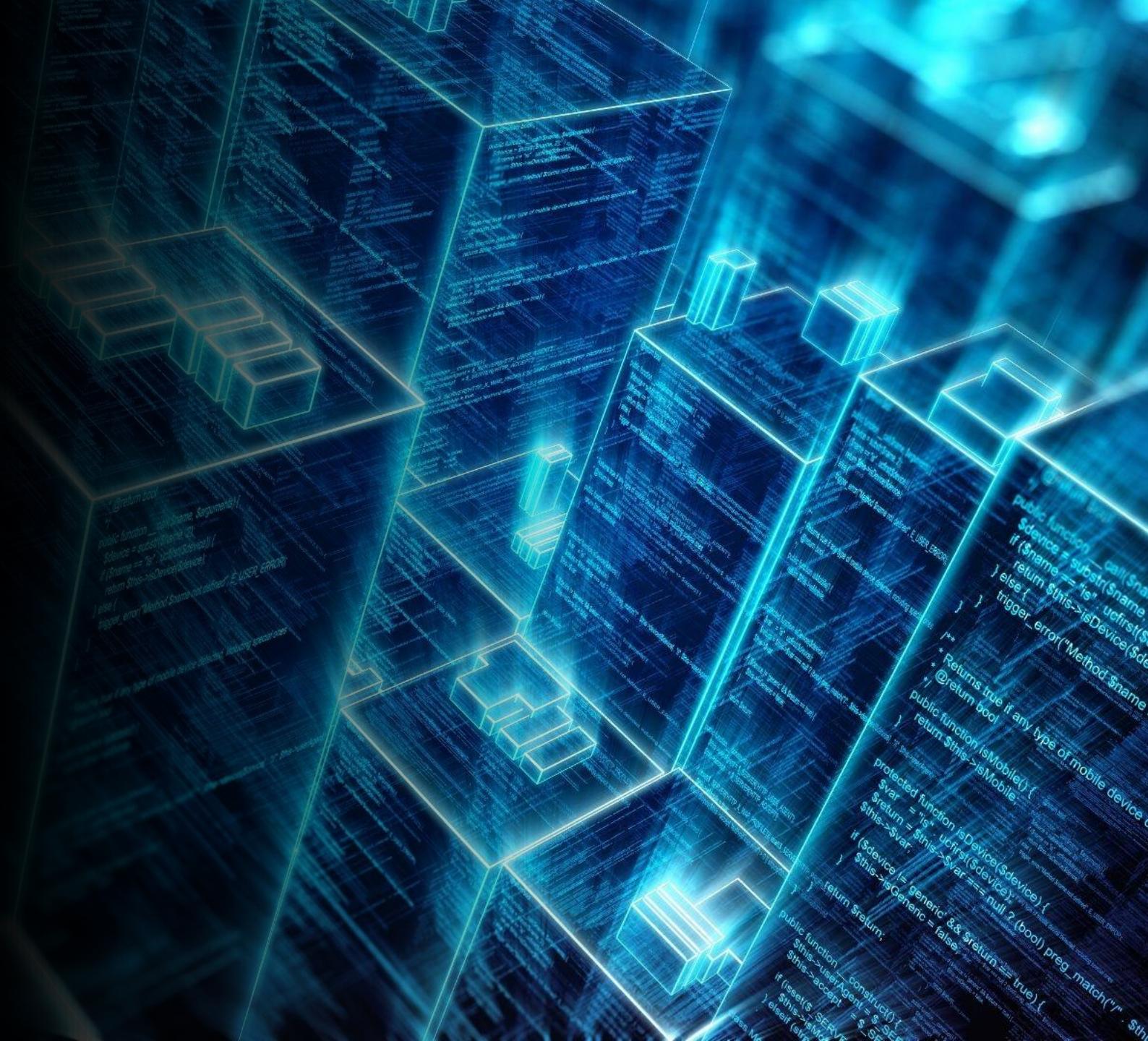
Batch Explorer

<https://github.com/Azure/BatchExplorer>

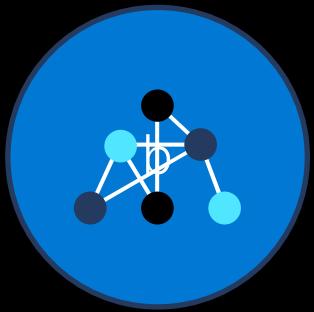
The screenshot displays the Azure Batch Explorer interface. On the left, a sidebar navigation bar includes links for Dash, Jobs, Pools, Packages, Data, and Profile. The main area shows a 'Batch Labs' dashboard for a job named '7c2060ad-163a-4bee-b97c-a099b2bc26dc'. The dashboard indicates 'No current background tasks' and features a 'Job statistics' chart. The chart shows 8808 Queued tasks, 94 Running tasks, and 0 Failed tasks, with a success rate of 15.77% and 1667 Succeeded tasks. Below the chart, a summary table provides details for 16 tasks.

ID	State	Created	Started	Completed	Exit code
1	completed	36 minutes ago	36 minutes ago	34 minutes ago	0
100001	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100002	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100003	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100004	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100005	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100006	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100007	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100008	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100009	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100010	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100011	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100012	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100013	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100014	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100015	completed	36 minutes ago	35 minutes ago	35 minutes ago	0
100016	completed	36 minutes ago	35 minutes ago	34 minutes ago	0

Incorporating Intelligence



Azure AI



Machine learning

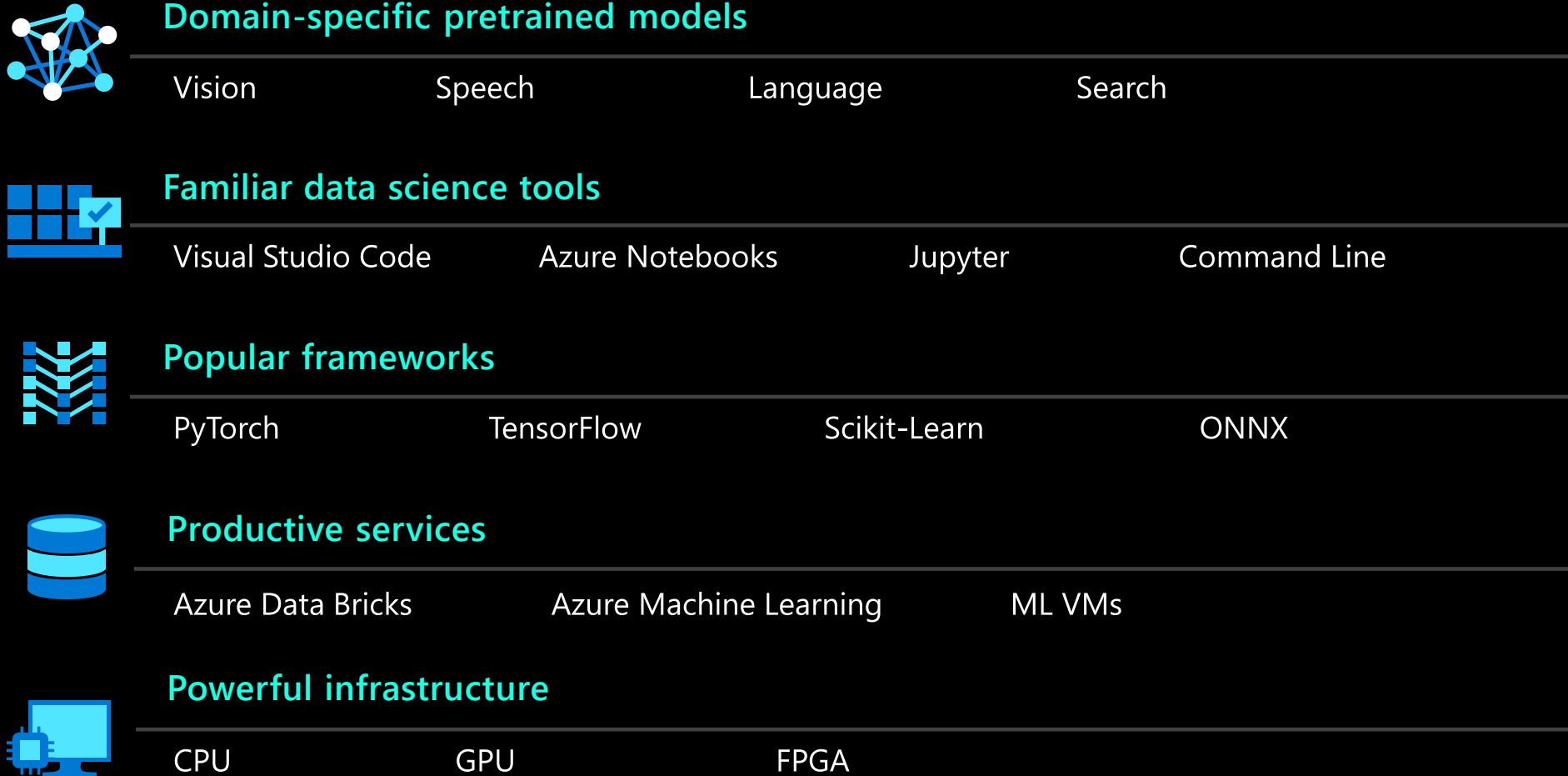


AI apps & agents



Knowledge mining

Machine Learning on Azure



Azure Machine Learning Compute

Addresses the challenges of training Machine Learning models

Challenges with training models

- Managing dependencies
- Securing access
- Scheduling jobs
- Feeding data
- Provisioning VMs
- Scaling resources
- Handling failures
- Gathering results



Solutions with Azure ML Compute

Integrates Batch AI to provide managed compute clusters and job scheduling

Work with clusters of GPU or CPU to run experiments in parallel and at scale to reduce training time

Enables data scientists to easily train, test, and score Deep Learning and other AI/ML models without managing infrastructure

Use with Azure Machine Learning service and integrate with your data science tools and pipeline



Azure Synapse - Spark

Data Engineering with Spark

Code-first Data Engineering

PySpark, Scala, SQL and C# languages supported

Author multiple languages in a single notebook

Analyze & transform data from the data warehouse, data lake, and real-time operational data from one place

Synapse Spark uses Spark 3.2 runtime, which includes Delta Lake 1.0

```
from pandas.tseries.frequencies import to_offset
from azureml.core._vendor.automl.client.core.common import metrics
from matplotlib import pyplot as plt
from automl.core.common import constants

def align_outputs(y_predicted, X_trans, X_test, y_test, target_column_name,
                  predicted_column_name='predicted',
                  horizon_colname='horizon_origin'):

    if (horizon_colname in X_trans):
        df_fcst = pd.DataFrame({predicted_column_name: y_predicted,
                               horizon_colname: X_trans[horizon_colname]})

    else:
        df_fcst = pd.DataFrame({predicted_column_name: y_predicted})

    # y and X outputs are aligned by forecast() function contract
    df_fcst.index = X_trans.index

    # align original X_test to y_test
    X_test_full = X_test.copy()
    X_test_full[target_column_name] = y_test

    # X_test_full's index does not include origin, so reset for merge
    df_fcst.reset_index(inplace=True)
    X_test_full = X_test_full.reset_index().drop(columns='index')
    together = df_fcst.merge(X_test_full, how='right')

    # drop rows where prediction or actuals are nan
    clean = together[[target_column_name,
                      predicted_column_name]].notnull().all(axis=1)]
    return(clean)

X_test[time_column_name] = pd.to_datetime(X_test[time_column_name])
df_all = align_outputs(y_predictions, X_trans, X_test, y_test, target_column_name)

# use automl metrics module
```

Spark 3.3

Enables developers can leverage the latest innovations in the Spark ecosystem

Pandas (Koalas) integration

A highly popular and flexible library with broad industry adoption

Adaptive Query Execution (AQE) enabled by default

Significant improvements in query performance out-of-the-box

Small Query execution improvements

Small queries run faster due to reduced initialization overhead

What is Apache Spark?

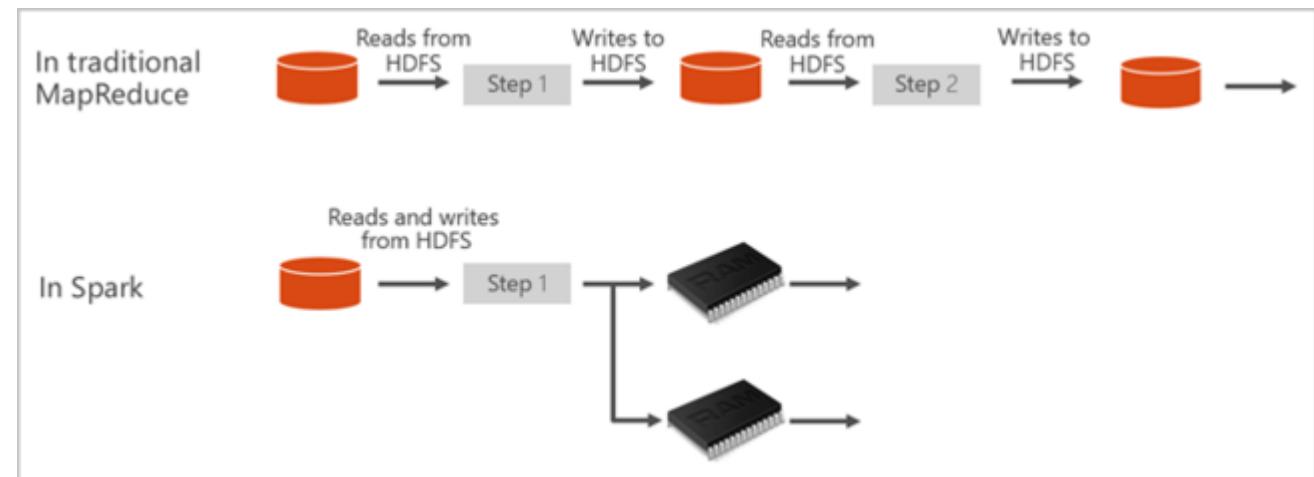
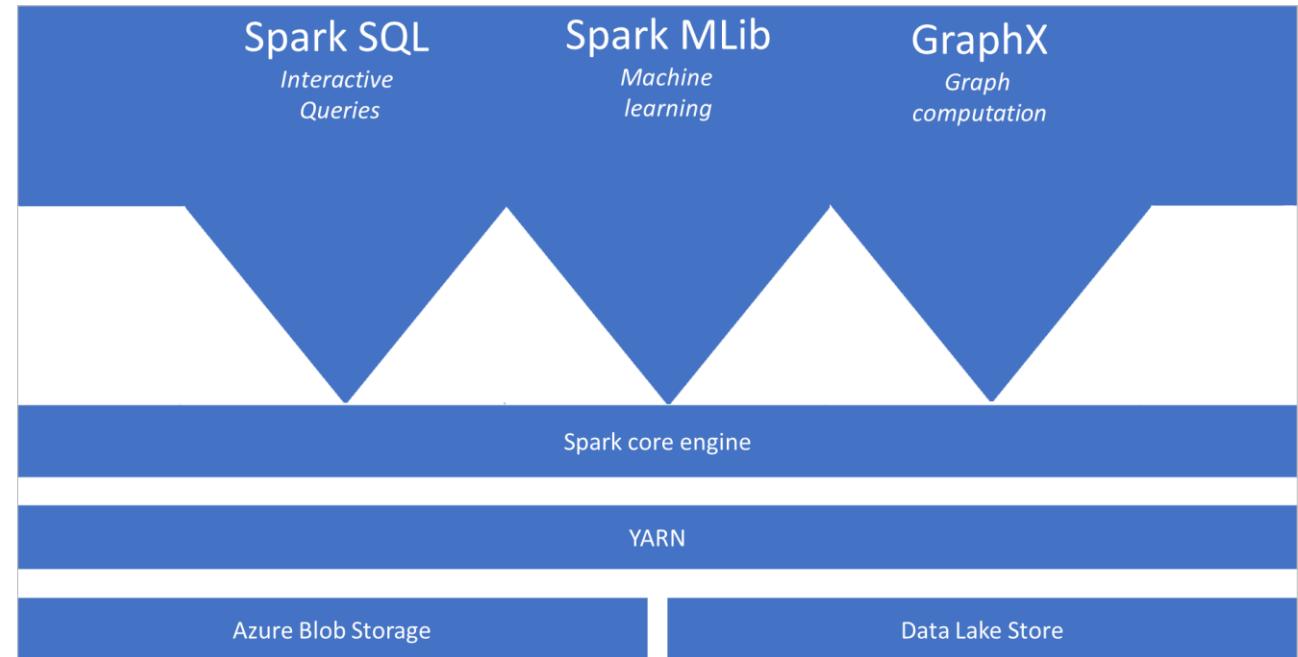
Parallel processing framework

In-memory processing engine to boost the performance of big-data analytic applications

Much faster than disk-based applications

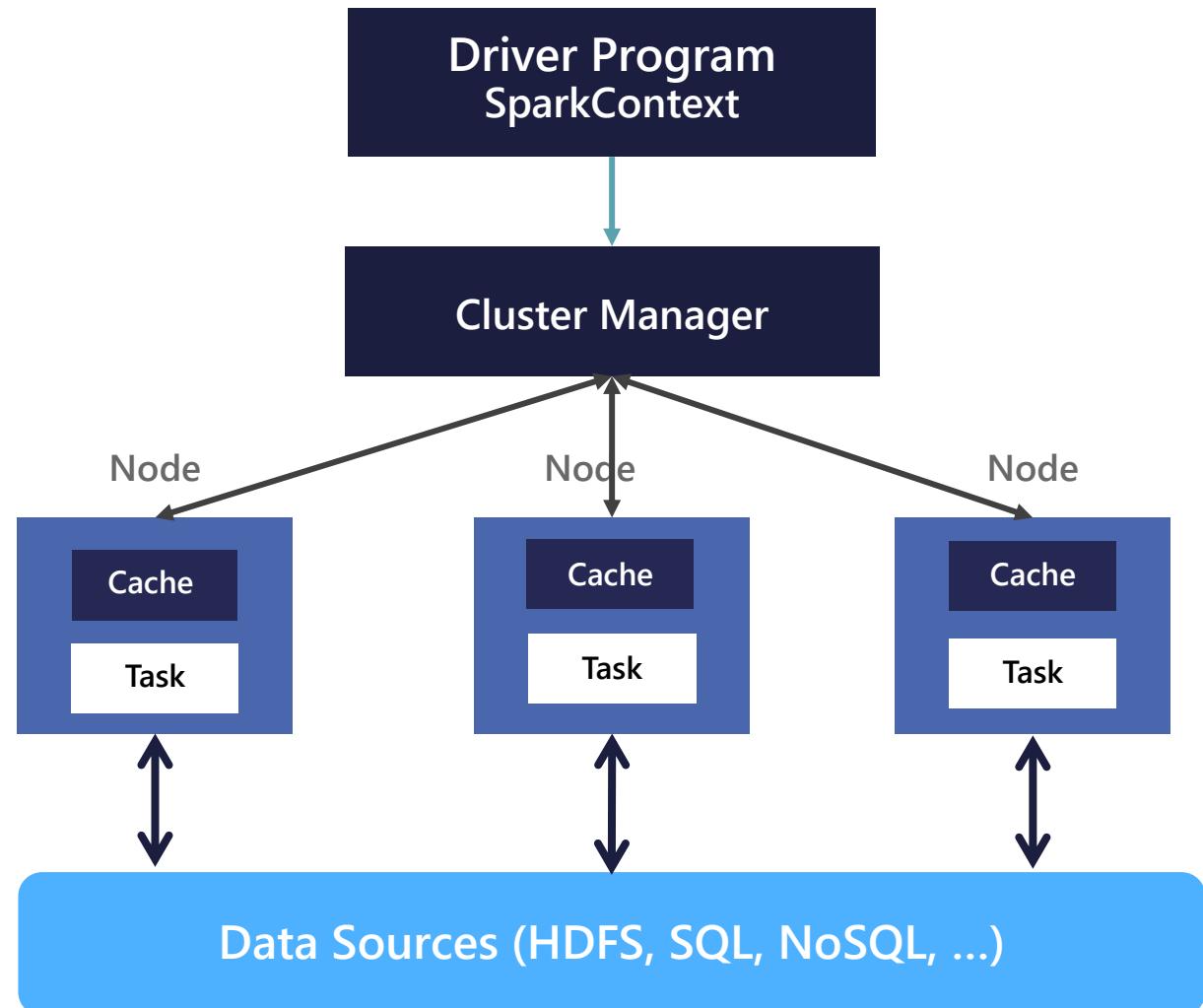
Integrates with multiple programming languages

Supports many workloads: Data Engineering, SQL, ML, Graphs., etc.



General Spark Cluster Architecture

- 'Driver' runs the user's 'main' function and executes the various parallel operations on the worker nodes.
- The results of the operations are collected by the driver
- The worker nodes read and write data from/to Data Sources including HDFS.
- Worker node also cache transformed data in memory as RDDs (Resilient Data Sets).
- Worker nodes and the Driver Node execute as VMs in public clouds (like Azure).
- These resources come in the form of *executors* on the worker nodes (logical blocks of CPU/memory). These executors are what read and write data for your application.
- The *SparkContext* breaks down the application code into *tasks* that are operated on in parallel by the executors.



Motivation for Spark Pool in Synapse

Speed and efficiency

- Spark instances start in approximately 2 minutes for fewer than 60 nodes and approximately 5 minutes for more than 60 nodes.
- The instance shuts down, by default, 5 minutes after the last job executed unless it is kept alive by a notebook connection.

Ease of creation

- Quickly create a new Spark pool in Azure Synapse using:
- Azure portal
 - Azure PowerShell
 - Synapse Analytics .NET SDK.

Ease of use

Synapse Analytics includes a custom notebook derived from [Nteract](#). You can use these notebooks for interactive data processing and visualization.

Motivation for Spark Pool in Synapse

REST APIs	Spark in Azure Synapse Analytics includes Apache Livy , a REST API-based Spark job server to remotely submit and monitor jobs.
Support for Azure Data Lake Storage Generation 2	Spark pools in Azure Synapse can use: <ul style="list-style-type: none">• Azure Data Lake Storage Generation 2• BLOB storage.
Integration with third-party IDEs	Azure Synapse provides an IDE plugin for JetBrains' IntelliJ IDEA that is useful to create and submit applications to a Spark pool.

Motivation for Spark Pool in Synapse

Pre-loaded Anaconda libraries

- Spark pools in Azure Synapse come with Anaconda libraries pre-installed.
- Providing close to 200 libraries for machine learning, data analysis, visualization, etc.

Scalability

- Can have Auto-Scale enabled, so that pools scale by adding or removing nodes as needed.
- Spark pools can be shut down with no loss of data since all the data is stored in Azure Storage or Data Lake Storage.

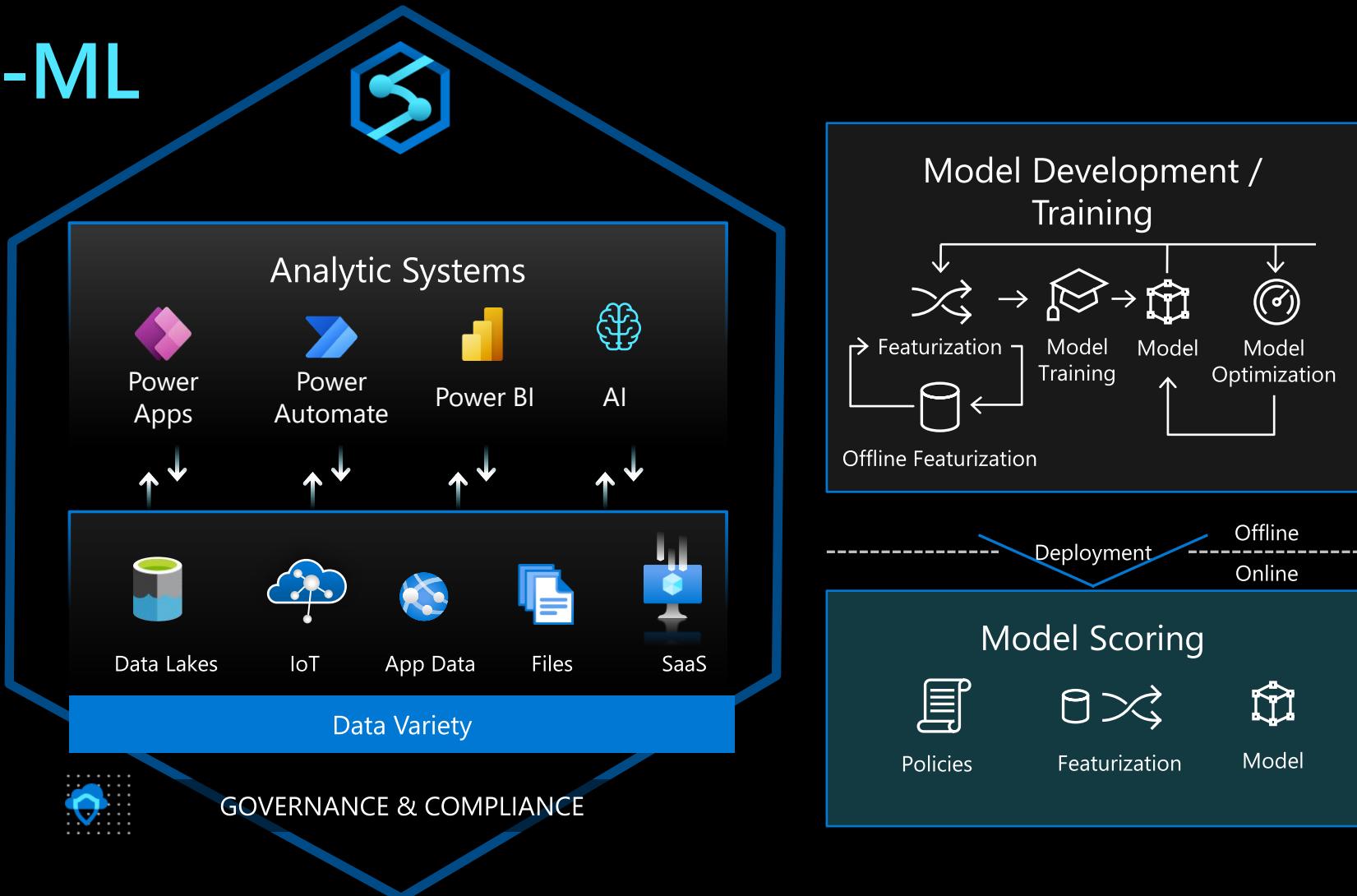


Modeling
capabilities

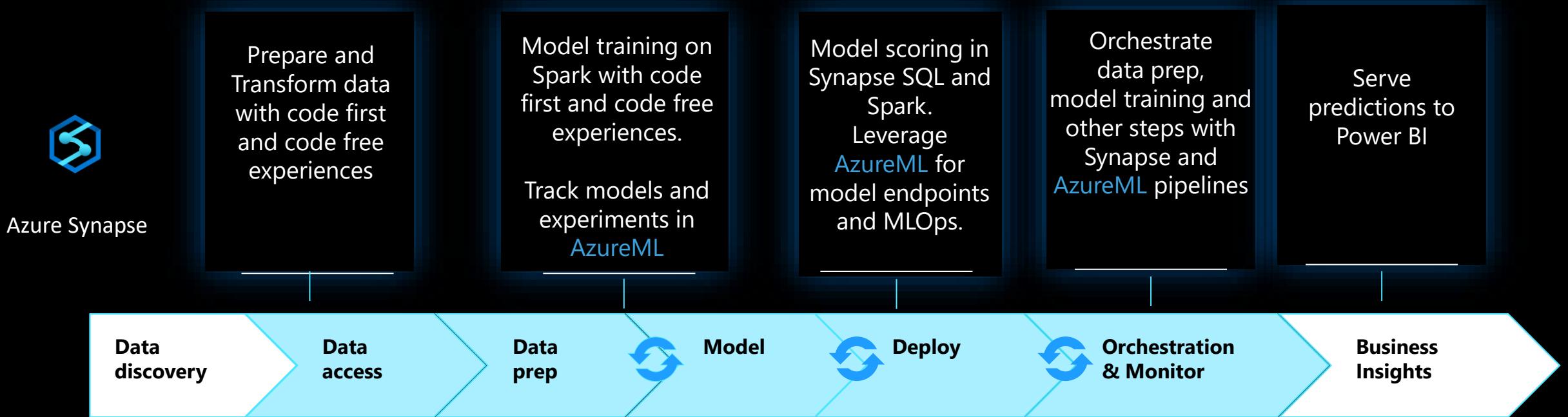
Why Synapse Enterprise Grade-ML

Productizing AI

- Train in the cloud within the Hub
- Scoring with operational systems
- Governance everywhere (models, lineage)
- Ethical AI
- Control over deployment
- Deployment across Apps, BI, Processes
- Exchange of models (ONNX)
- Enabling Reinforcement learning



Synapse & Azure ML: Supporting the full Data & AI lifecycle



[Data wrangling with Apache Spark pools \(preview\) - Azure Machine Learning | Microsoft Learn](#)

[MachineLearningNotebooks/how-to-use-azureml/azure-synapse at master · Azure/MachineLearningNotebooks \(github.com\)](#)

Synapse Apache Spark in Azure Machine Learning

Link workspace and attach compute

Link and attach compute via ARM, SDK and UI

Authentication via Identity

Permission control based on Role-Based Access Control (RBAC)

Leverage Spark magic for data prep

Run spark magic on cell to do data prep

Data access via HDFS data loading or Azure ML dataset

Pick up the data via dataset for training

Orchestration

Orchestrate data prep and ML steps via Azure ML pipeline

The screenshot shows the Microsoft Azure Machine Learning Studio interface. On the left, there's a sidebar with various options like Home, Notebooks, Automated ML, Designer, Assets, Datasets, Experiments, Pipelines, Models, Endpoints, Compute, Datastores, Data Labeling, and Linked Services. The main area is titled "Notebooks" and shows a list of files under "Files". One file, "synapse-big-dataprep-demo...", is currently selected. The right side of the screen shows a Jupyter notebook interface with the following content:

```
[49]: 1 # show help
      2 %synapse ?
```

...
%synapse [-s SUBSCRIPTION_ID] [-r RESOURCE_GROUP] [-w WORKSPACE_NAME]
[-f CONFIG_FILE] [-c COMPUTE_TARGET]
[--driver-memory DRIVER_MEMORY] [--driver-cores DRIVER_CORES]
[--executor-memory EXECUTOR_MEMORY]
[--executor-cores EXECUTOR_CORES] [-n NUM_EXECUTORS]
[-t SESSION_TIMEOUT] [--start-timeout START_TIMEOUT]
[-e ENVIRONMENT] [--environment-version ENVIRONMENT_VERSION]
[command [command ...]]

Magic to execute spark remotely against a Synapse Spark pool.

Sub commands:
start: Start a Livy session against target AML Synapse compute.
You can set spark config in the magic body by json format. e.g.

1. Start Spar Session

Spark Pool - Configuration

- **Nodes**
 - Apache Spark pool instance consists of one head node and two or more worker nodes with a minimum of three nodes in a Spark instance.
 - The head node runs additional management services such as Livy, Yarn Resource Manager, Zookeeper, and the Spark driver.
 - All nodes run services such as Node Agent and Yarn Node Manager.
 - All worker nodes run the Spark Executor service.

Spark Pool - Configuration

- **Node Sizes**

- Node sizes can be altered after pool creation, instance may need to be restarted.

Size	vCore	Memory
Small	4	32 GB
Medium	8	64 GB
Large	16	128 GB
XLarge	32	256 GB
XXLarge	64	432 GB
XXX Large (Isolated Compute)	80	504 GB

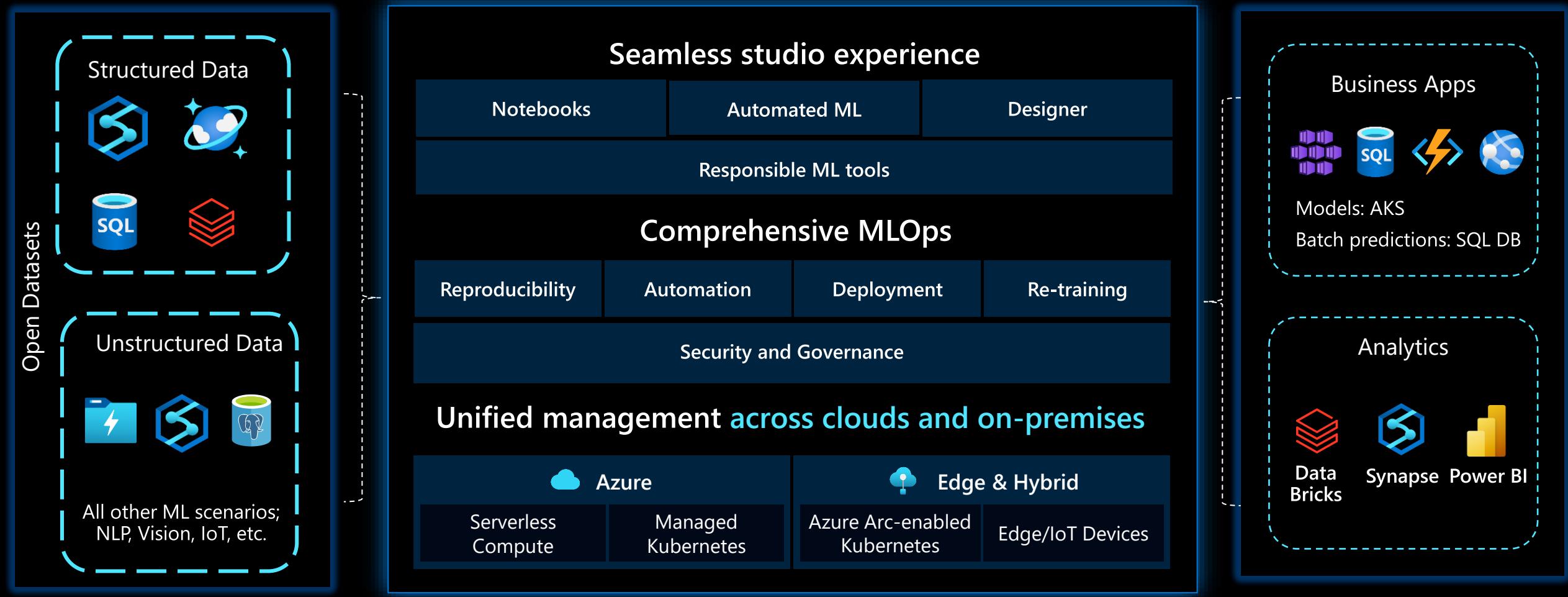
ML Platform with Azure Machine Learning



Machine Learning Platform Architecture



Azure ML



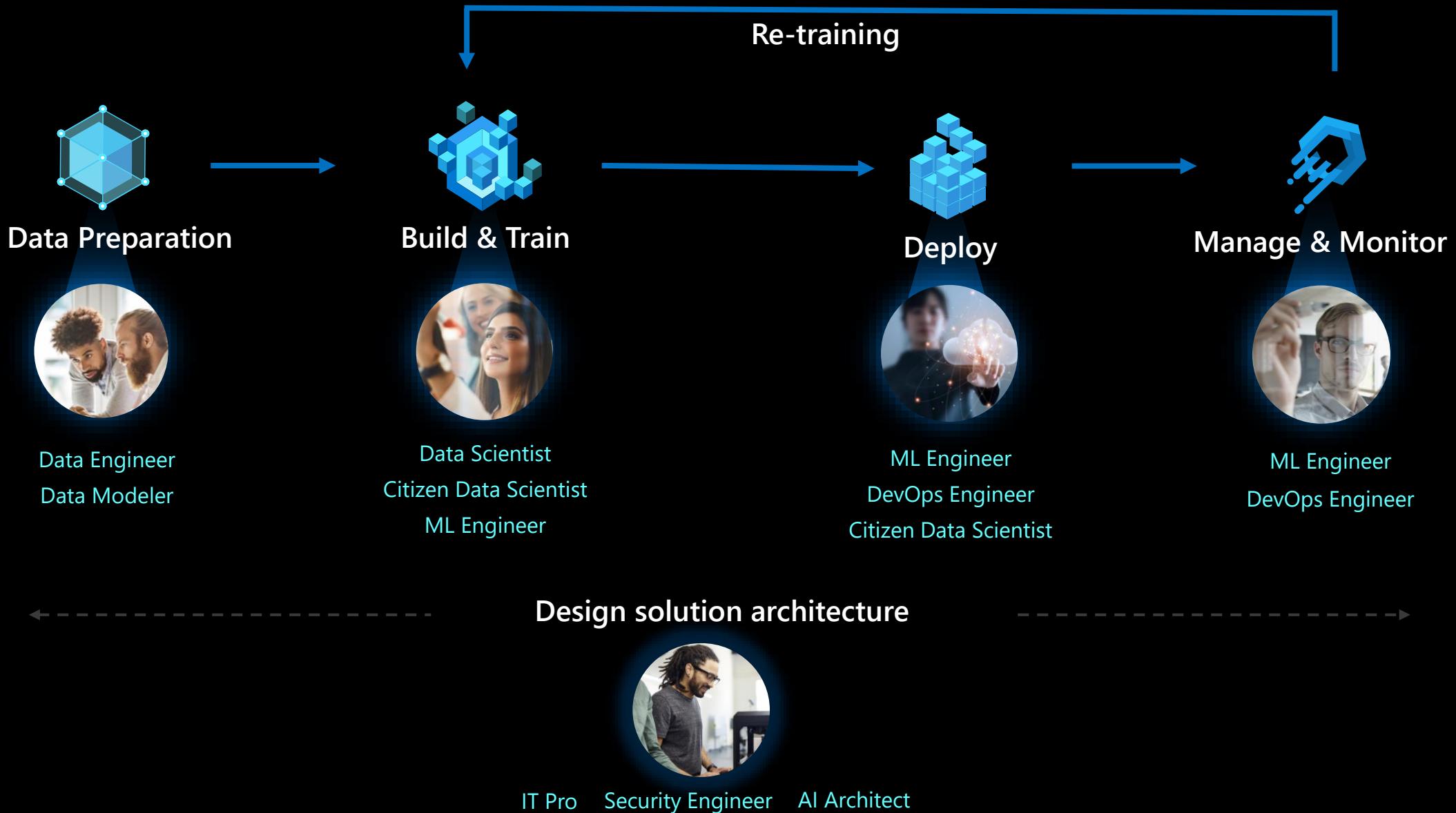
Prepare Data

Build & Train

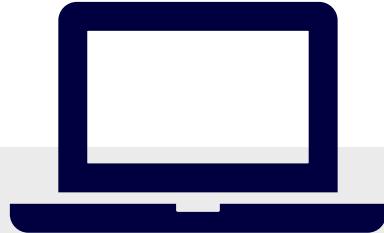
Deploy

Manage & Monitor

ML Lifecycle workflows



Compute Options for Experiment Runs



Local Compute

- Compute where the control code for the experiment is running
- Often a development workstation or Azure Machine Learning compute instance



Compute Cluster

- Cloud-based cluster managed in an Azure Machine Learning workspace
- Starts, stops, and scales on-demand



Attached Compute

- Azure compute resource outside of a workspace
- For example:
 - Virtual Machine
 - Azure Databricks
 - Azure HDInsight

What are Compute Targets in Azure ML?

A *compute target* is a designated compute resource or environment where you run your training script or host your service deployment.

In a typical model development lifecycle, you might:

1. Start by developing and experimenting on a small amount of data. At this stage, use your [local environment](#), such as a local computer or cloud-based virtual machine (VM), as your compute target.
2. Scale up to larger data, or do [distributed training](#) by using one of these [training compute targets](#).
3. After your model is ready, deploy it to a web hosting environment with one of these [deployment compute targets](#).

Training Compute Targets

Training targets	Automated machine learning	Machine learning pipelines	Azure Machine Learning designer
Local computer	Yes		
Azure Machine Learning compute cluster	Yes	Yes	Yes
Azure Machine Learning compute instance	Yes (through SDK)	Yes	Yes
Azure Machine Learning Kubernetes	Yes	Yes	Yes
Remote VM	Yes	Yes	
Apache Spark pools (preview)	Yes (SDK local mode only)	Yes	
Azure Databricks	Yes (SDK local mode only)	Yes	
Azure Data Lake Analytics		Yes	
Azure HDInsight		Yes	
Azure Batch		Yes	

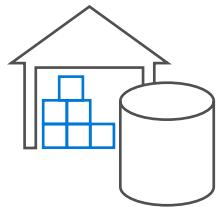
Inference Compute Targets

When performing inference, Azure Machine Learning creates a Docker container that hosts the model and associated resources needed to use it. This container is then used in a compute target

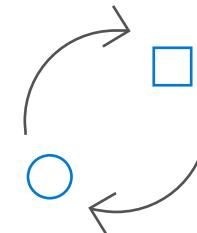
Compute target	Used for	GPU support	Description
Local web service	Testing/debugging		Use for limited testing and troubleshooting. Hardware acceleration depends on use of libraries in the local system.
Azure Machine Learning endpoints	Real-time inference Batch inference	Yes	Fully managed computes for real-time (managed online endpoints) and batch scoring (batch endpoints) on serverless compute.
Azure Machine Learning Kubernetes	Real-time inference Batch inference	Yes	Run inferencing workloads on on-premises, cloud, and edge Kubernetes clusters.
Azure Container Instances (SDK/CLI v1 only)	Real-time inference Recommended for dev/test purposes only.		Use for low-scale CPU-based workloads that require less than 48 GB of RAM. Doesn't require you to manage a cluster. Supported in the designer.

AZURE DATA FACTORY

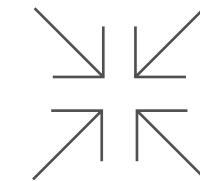
Fully-managed data integration service in the cloud



Flexible
Data integration



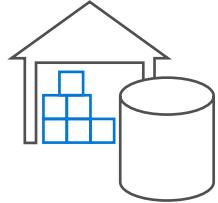
Hybrid
Data orchestration



Data movement
As-a-service

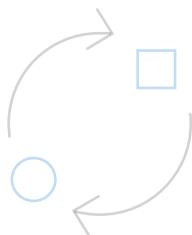


Security and compliance



Flexible
Data integration

Modernize your data warehouse with Azure big data and advanced analytics services such as HDInsight and Data lake Analytics



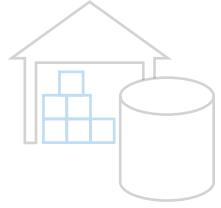
Hybrid
Data orchestration

Build custom data-driven SaaS applications unique to your customer data using your language of choice



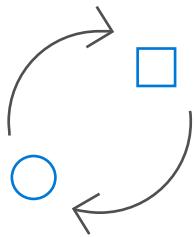
Data movement
As-a-service

Bring together all your sources of data to understand your customers and drive impactful business decisions



Flexible
Data integration

Orchestrate your data pipeline wherever your data lives – in cloud or in self-hosted environment



Hybrid
Data orchestration

Meet your security and compliance needs while taking advantage of truly hybrid integration capabilities



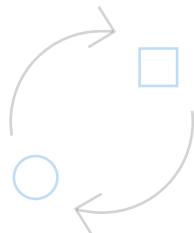
Data movement
As-a-service

Execute your SQL Server Integration Services (SSIS) packages in the cloud



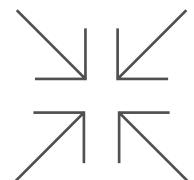
Flexible
Data integration

Accelerate integration with managed data movement as-a-service



Hybrid
Data orchestration

Improve your TCO with 30+ natively supported connectors across 18 global points of presence

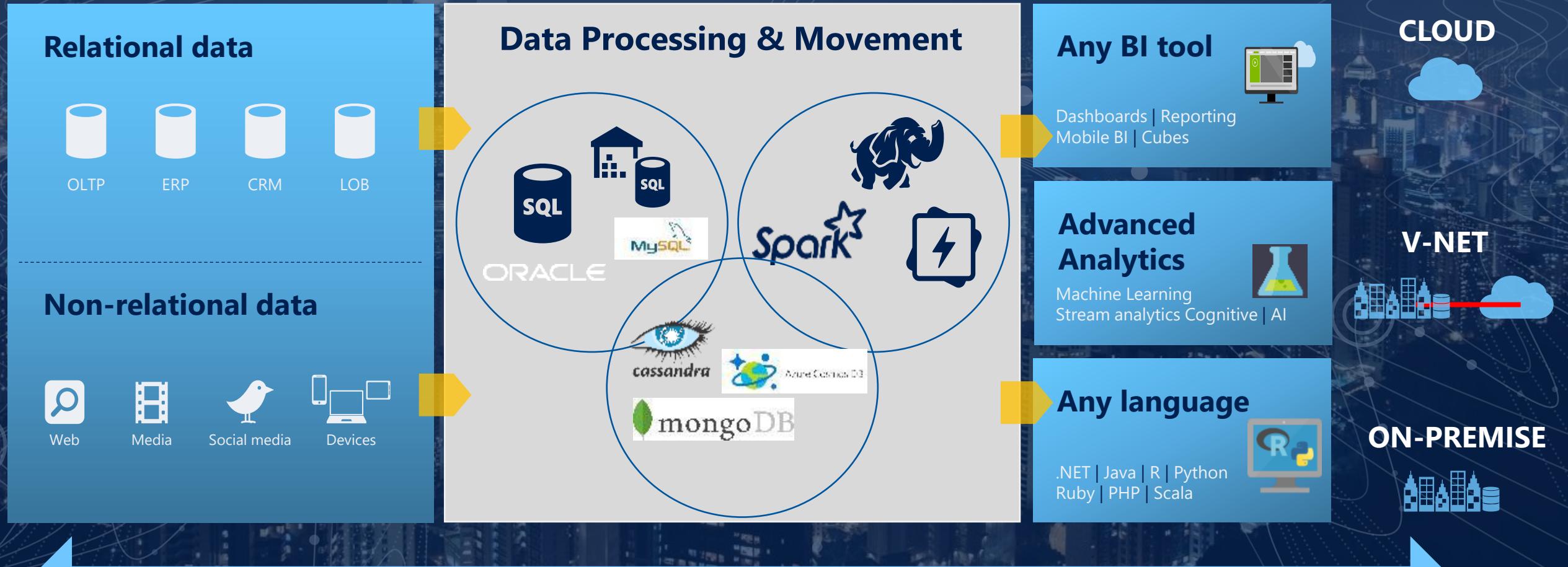


Data movement
As-a-service

Elastic data movement at scale

Serverless data movement with no infrastructure to manage

HYBRID DATA INTEGRATION AT SCALE



AZURE DATA FACTORY ORCHESTRATES DATA PIPELINE ACTIVITY WORKFLOW & SCHEDULING

ADF: Cloud-First Data Integration Objectives

- Consume hybrid disparate data
 - On-prem + Cloud
 - Grow ADF ecosystem of structured, un-structured, semi-structured data connectors
- Calculate and format data for analytics
 - Transform, aggregate, join, normalize
 - Separate data flow (transformation) from control flow (orchestration)
- Address large-scale Big Data requirements
 - Scale-up or Scale-out data movement and transformation
 - Support multiple processing engines
- Operationalize
 - Support flexible scheduling and triggering mechanism for broad range of use cases
 - Manage & monitor multiple pipelines (via Azure Monitor & OMS)
 - Support secure VNET environments
- Enable SSIS package execution
 - Execute SSIS packages in ADF Integration Runtime

ADF: Cloud-First Data Integration Scenarios

Lift and Shift to the Cloud

- Migrate on-prem DW to Azure
- Lift and shift existing on-prem SSIS packages to cloud
- No changes needed to migrate SSIS packages to Cloud service

DW Modernization

- Modernizing DW arch to reduce cost & scale to needs of big data (volume, variety, etc)
- Flexible wall-clock and triggered event scheduling
- Incremental Data Load

Build Data-Driven, Intelligent SaaS Application

- C#, Python, PowerShell, ARM support

Big Data Analytics

- Customer profiling, Product recommendations, Sentiment Analysis, Churn Analysis, Customized offers, customer usage tracking, customized marketing
- On-demand Spark cluster support

Load your Data Lake

- Separate control-flow to orchestrate complex patterns with branching, looping, conditional processing

ADF V2 Improvements

- Integration Runtimes (IR) replace DMG, provide data movement and activity dispatch on-prem or in the cloud
- Supports resources within virtual networks
- Integration Runtime includes SSIS option to lift & shift SSIS packages to the Cloud
- Separation of “control flow” & “data flow” capabilities for more flexible pipeline management
 - Looping, conditionals, dependencies, parameters
- Python SDK
- On-Demand Spark support
- Flexible pipeline scheduling with wall-clock and triggered executions
- Expanded use cases: From primarily time window-oriented pipelines, to trigger-based on-demand for more flexible ETL and data integration orchestrations
- (Coming in GA) UX pipeline and data transformation builder code-free experience

New ADF V2 Concepts

Concept	Description	Sample
Control Flow	Orchestration of pipeline activities that includes chaining activities in a sequence, branching, conditional branching based on an expression, parameters that can be defined at the pipeline level and arguments passed while invoking the pipeline on demand or from a trigger. Also includes custom state passing and looping containers, i.e. For-each, Do-Until iterators.	{ "name":"MyForEachActivityName", "type":"ForEach", "typeProperties":{ "isSequential":true}, "items": "@pipeline().parameters.mySinkDatasetFolderPathCollection", " "activities": [{ "name":"MyCopyActivity", "type":"Copy", "typeProperties": ...
Runs	A Run is an instance of the pipeline execution. Pipeline Runs are typically instantiated by passing the arguments to the parameters defined in the Pipelines. The arguments can be passed manually or properties created by the Triggers.	POST <a href="https://management.azure.com/subscriptions/<subId>/resourceGroups/<resourceGroupName>/providers/Microsoft.DataFactory/factories/<dataFactoryName>/pipelines/<pipelineName>/createRun?api-version=2017-03-01-preview">https://management.azure.com/subscriptions/<subId>/resourceGroups/<resourceGroupName>/providers/Microsoft.DataFactory/factories/<dataFactoryName>/pipelines/<pipelineName>/createRun?api-version=2017-03-01-preview
Activity Logs	Every activity execution in a pipeline generates activity start and activity end logs event	
Integration Runtime	Replaces DMG as a way to move & process data in Azure PaaS Services, self-hosted or on prem or IaaS Works with VNets Enables SSIS package execution	
Scheduling	Flexible Scheduling Wall-clock scheduling Event-based triggers	"type": "ScheduleTrigger", "typeProperties": { "recurrence": { "frequency": <<Minute, Hour, Day, Week, Year>>, "interval": <<int>>, // optional, how often to fire (default to 1) "startTime": <<datetime>>, "endTime": <<datetime>>, "timeZone": <<default UTC>> "schedule": { // optional (advanced scheduling specifics) "hours": [<<0-24>>], "weekDays": ": [<<Monday-Sunday>>], "minutes": [<<0-60>>], "monthDays": [<<1-31>>], "monthlyOccurrences": [{ "day": <<Monday-Sunday>>, "occurrence": <<1-5>>

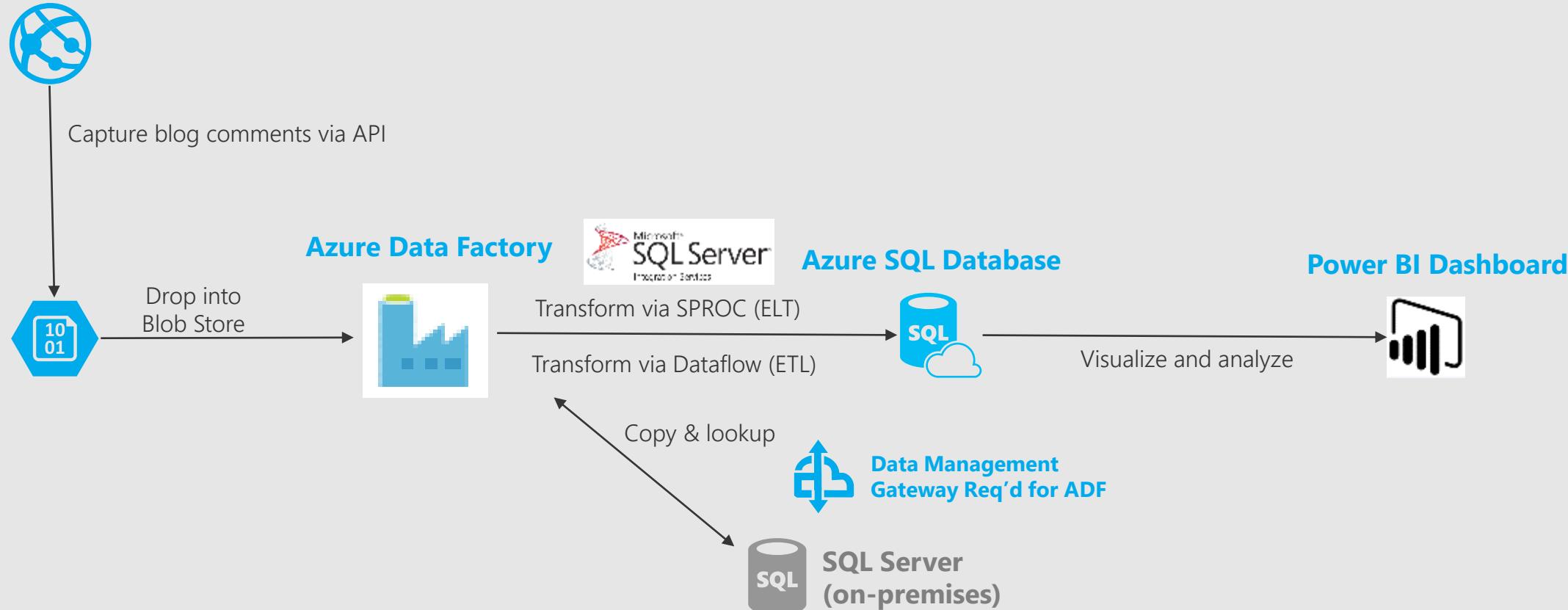
New ADF V2 Concepts

Concept	Description	Sample
On-Demand Execution	Instantiate a pipeline by passing arguments as parameters defined in a pipeline and execute from script / REST / API.	<code>Invoke-AzureRmDataFactoryV2PipelineRun -DataFactory \$df -PipelineName "Adfv2QuickStartPipeline" -ParameterFile .\PipelineParameters.json</code>
Parameters	<p>Name-value pairs defined in the pipeline. Arguments for the defined parameters are passed during execution from the run context created by a Trigger or pipeline executed manually. Activities within the pipeline consume the parameter values.</p> <p>A Dataset is a strongly typed parameter and a reusable/referenceable entity. An activity can reference datasets and can consume the properties defined in the Dataset definition</p> <p>A Linked Service is also a strongly typed parameter containing the connection information to either a data store or a compute environment. It is also a reusable/referenceable entity.</p>	<p>Accessing parameters of other activities Using expressions <code>@parameters("{name of parameter}")</code></p> <p><code>@activity("{Name of Activity}").output.RowsCopied</code></p>
Incremental Data Loading	Leverage parameters and define your high-water mark for delta copy while moving dimension or reference tables from a relational store either on premises or in the cloud to load the data into the lake	

Patterns & Scenarios for ADF V2

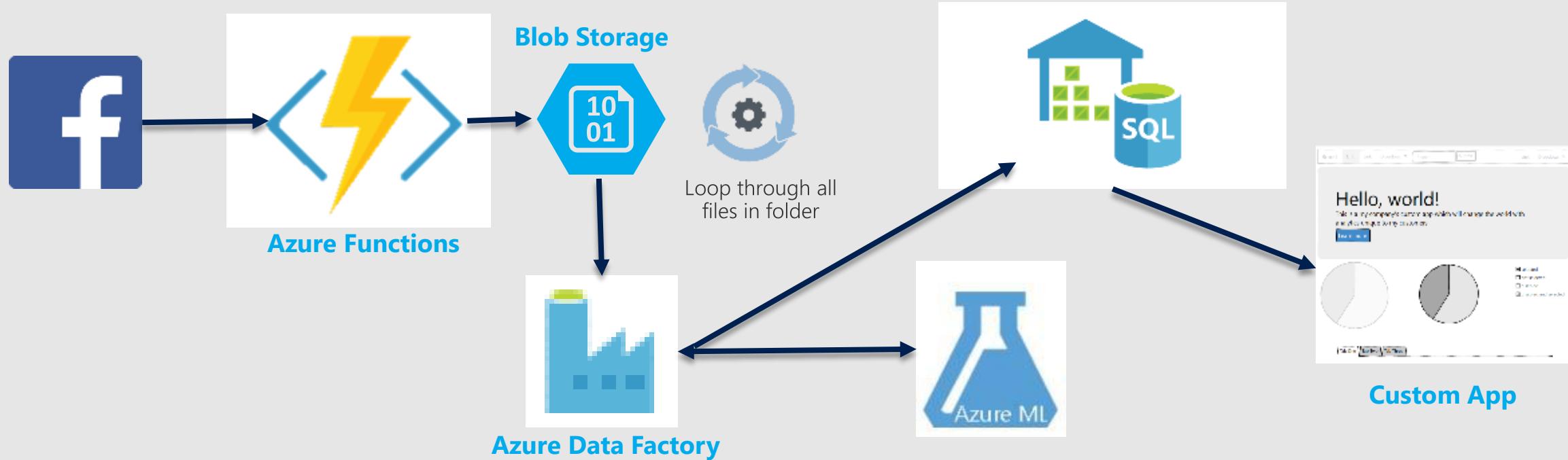
Hybrid Data Integration Pattern 1:

Analyze blog comments

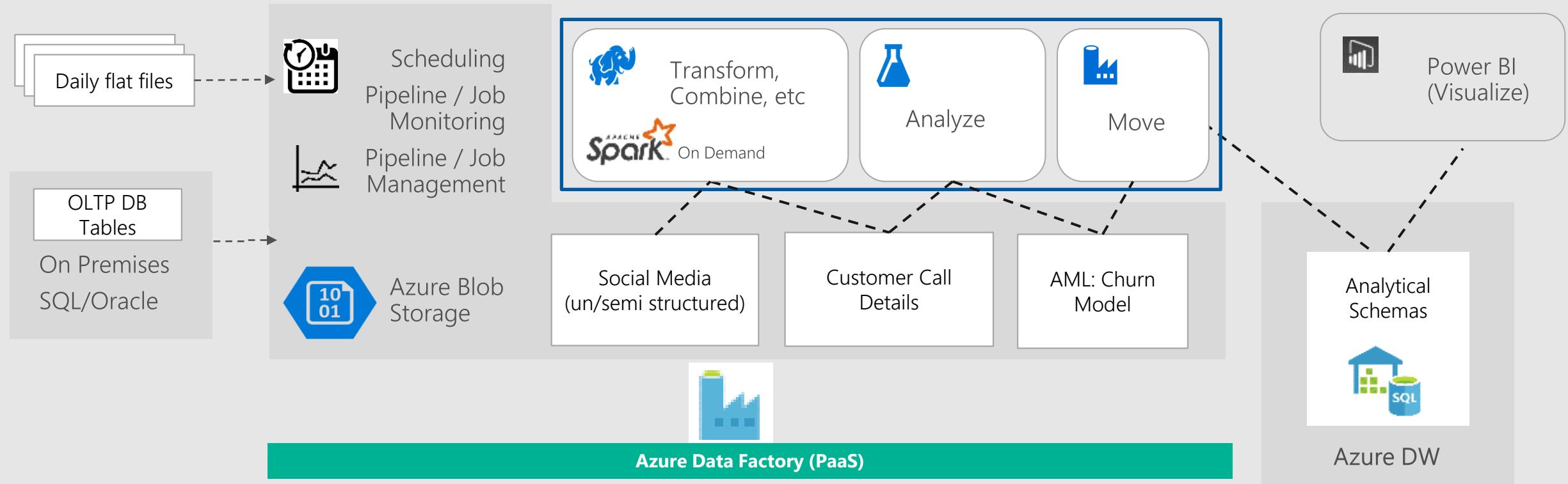


Hybrid Data Integration Pattern 2:

Data-Driven SaaS App: Sentiment Analysis w/Machine Learning

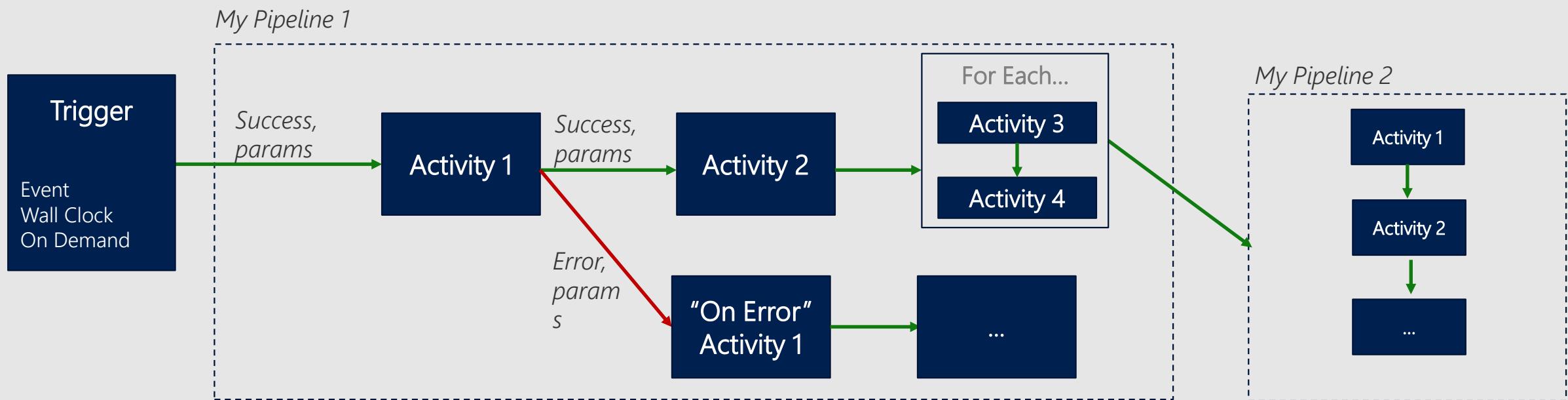


Hybrid Data Integration Pattern 3: Modern Data Warehouse



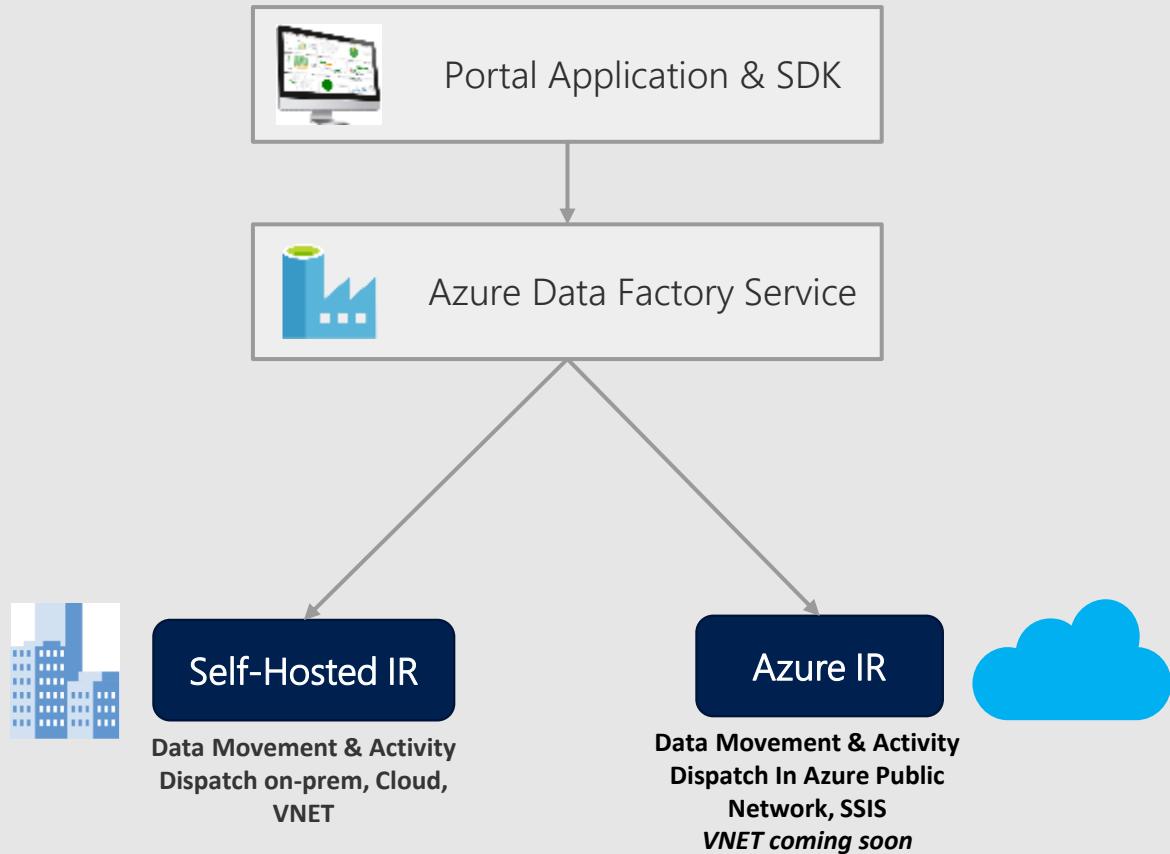
ADFv2: Control Flow

Coordinate pipeline activities into finite execution steps to enable looping, conditionals and chaining while separating data transformations into individual data flows



Integration Runtime

ADF Integration Runtime (IR)



- ADF compute environment with multiple capabilities:
 - Activity dispatch & monitoring
 - Data movement
 - SSIS package execution
- To integrate data flow and control flow across the enterprises' hybrid cloud, customer can instantiate multiple IR instances for different network environments:
 - On premises (similar to DMG in ADF V1)
 - In public cloud
 - Inside VNet
- Bring a consistent provision and monitoring experience across the network environments

←→ Command & Control

↔ Data Flow



UX & SDK

Authoring | Monitoring/Mgmt

Azure Data Factory Service

Scheduling | Orchestration | Monitoring

On Premises Apps & Data



TERADATA



cloudera



ORACLE

Cloud Apps, Svcs & Data



Microsoft Azure



Adobe

workday

↔ Command & Control

↔ Data Flow



UX & SDK

Authoring | Monitoring/Mgmt

Azure Data Factory Service

Scheduling | Orchestration | Monitoring

Azure Cloud

PaaS Cloud Host

Integration Runtime

On Premises Apps & Data



TERADATA



cloudera



ORACLE

Cloud Apps, Svcs & Data



Adobe



↔ Command & Control

↔ Data Flow



UX & SDK

Authoring | Monitoring/Mgmt

Azure Data Factory Service

Scheduling | Orchestration | Monitoring

Azure Cloud

PaaS Cloud Host

Integration Runtime

Installable Agent

Integration Runtime



On Premises Apps & Data



TERADATA



cloudera

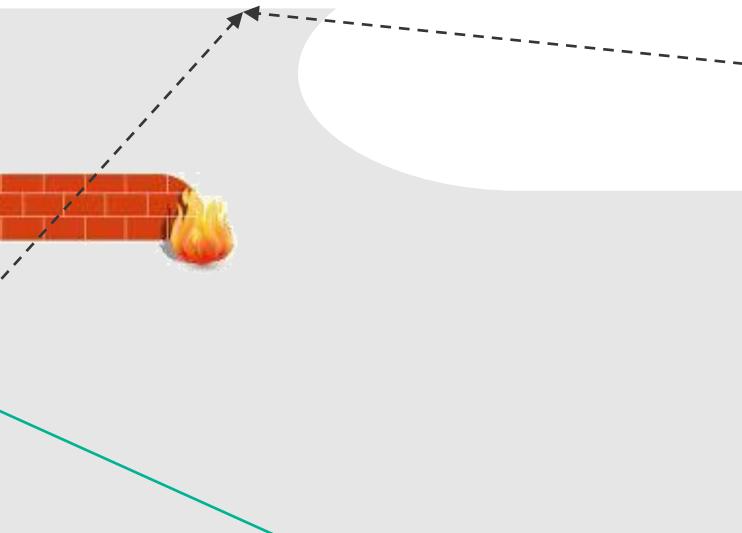


ORACLE

Cloud Apps, Svcs & Data



Adobe



←→ Command & Control

↔ Data Flow



UX & SDK

Authoring | Monitoring/Mgmt

Azure Data Factory Service

Scheduling | Orchestration | Monitoring

Azure Cloud

PaaS Cloud Host

Runtime

Integration Runtime

- Activity Dispatch/Monitor (spark, copy, ml, etc)
- Data Movement
- SSIS Package Execution

On Premises Apps & Data



TERADATA



cloudera



ORACLE

Cloud Apps, Svcs & Data

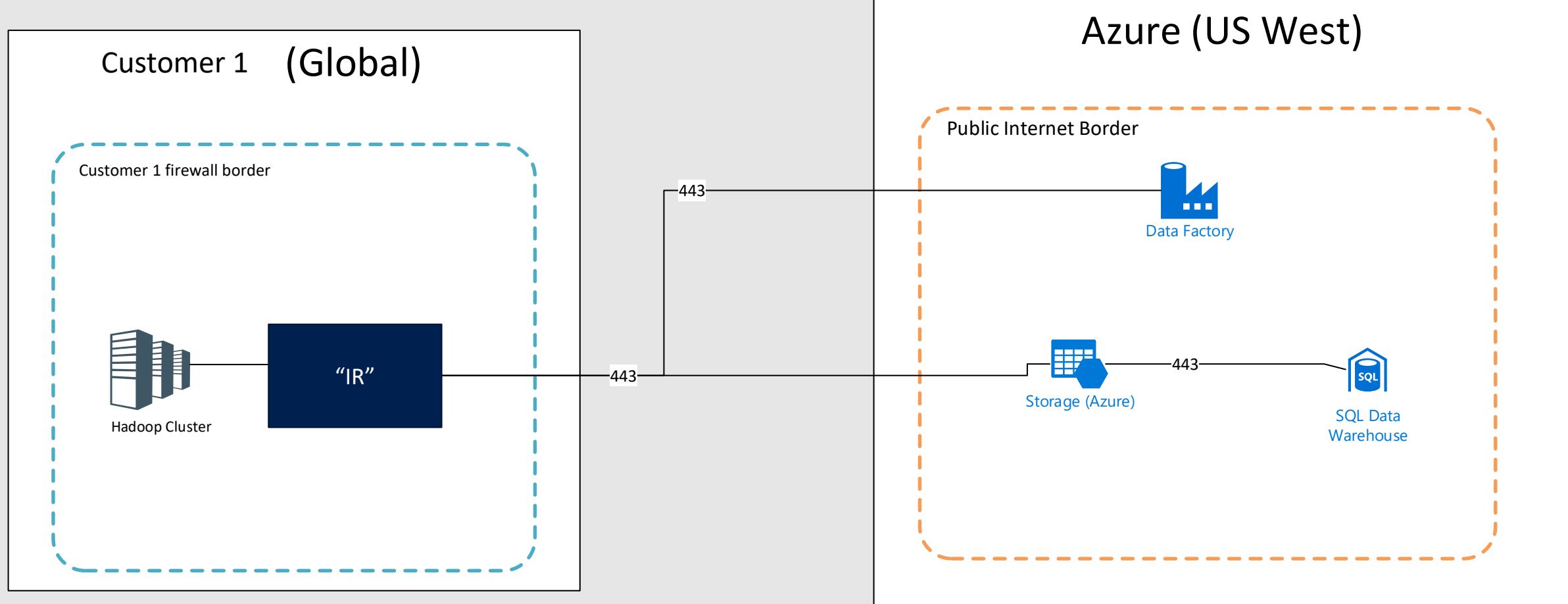


Adobe

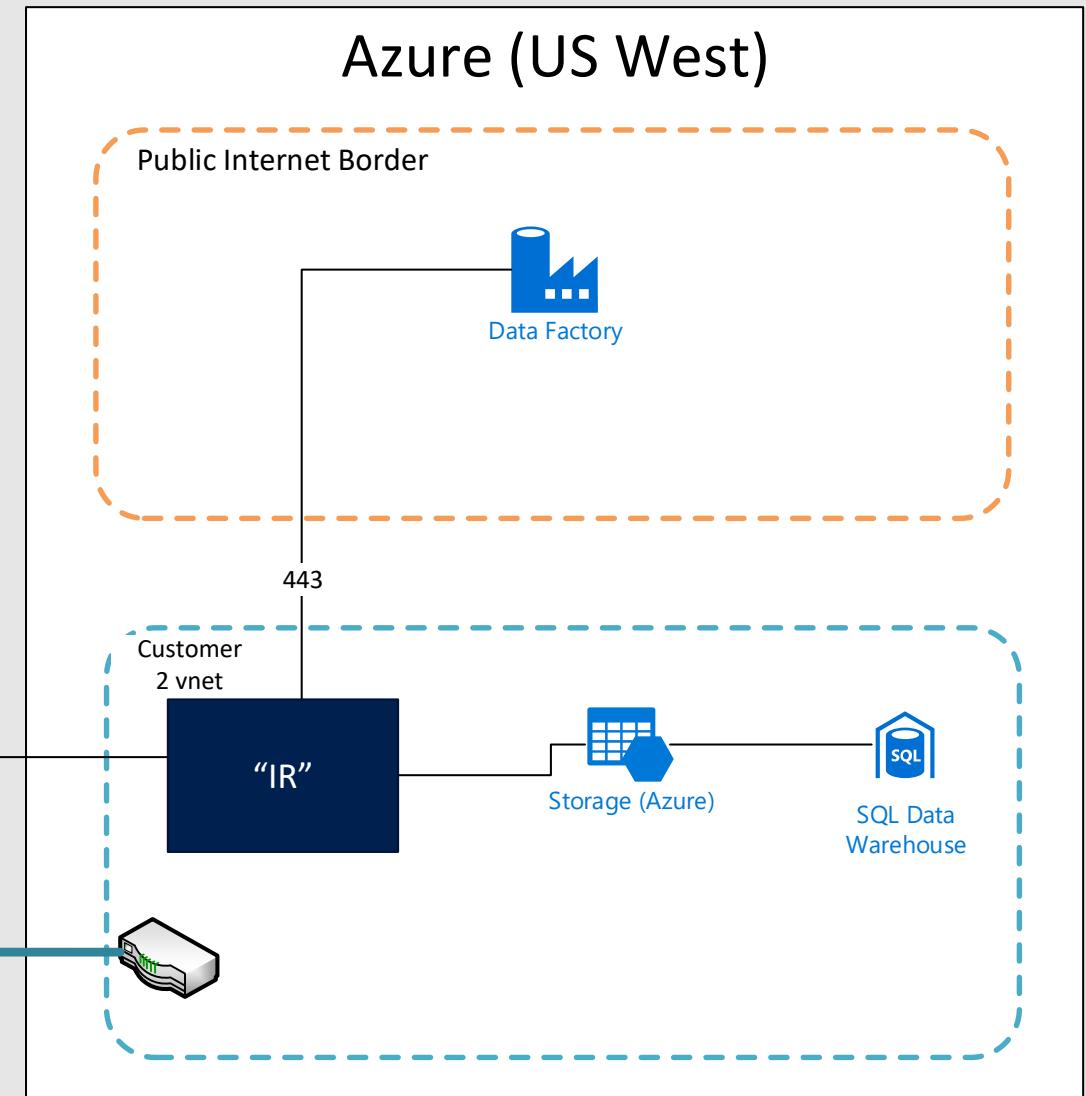
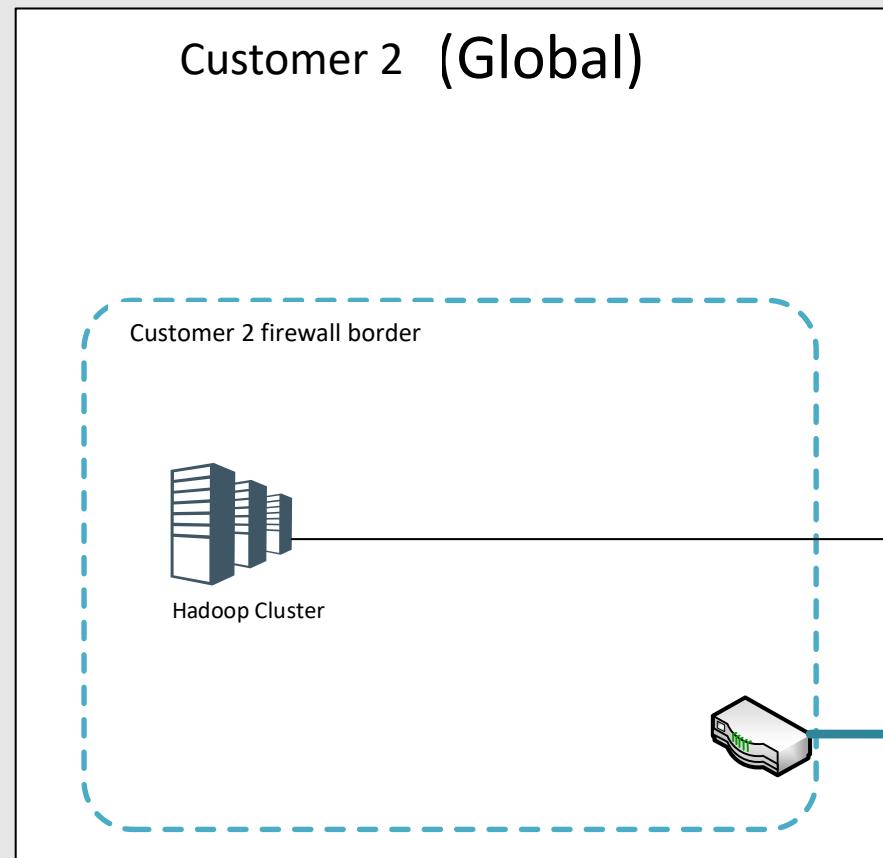
workday



Azure Data Factory “Integration Runtime” deployed on premises for transformation and then moved to cloud

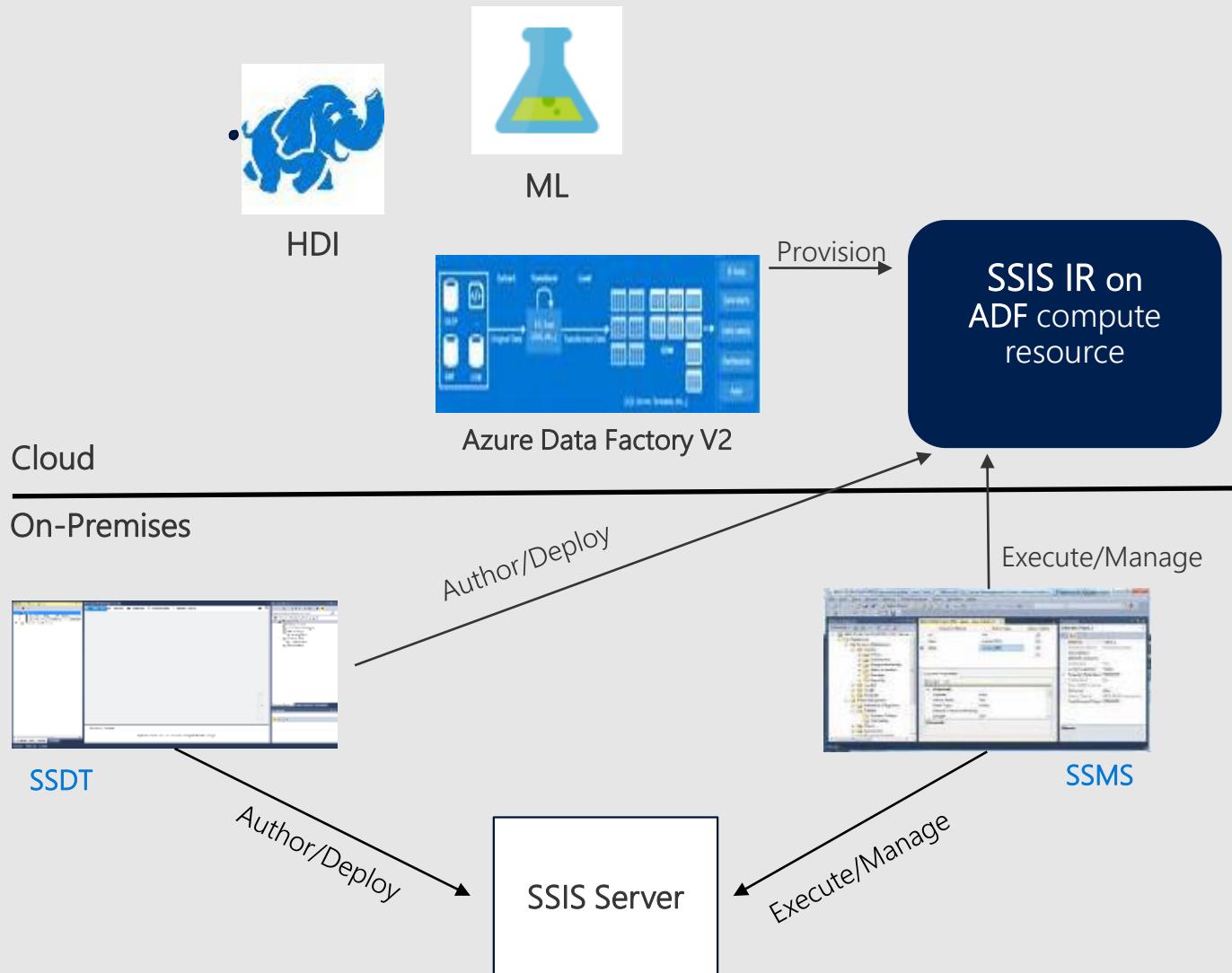


Azure Data Factory “Integration Runtime” deployed inside VNet



SSIS in ADF

SSIS as an Integration Runtime



- Provision SSIS IR via Data Factory
- Use SQL Server Data Tools (SSDT) to author/deploy SSIS packages
- Use SQL Server Management Studio (SSMS) to execute/manage SSIS packages
- Target SSIS customers who want to move all/part of their on-premises workloads and just "lift & shift" many existing packages to Azure
- Independent Software Vendors (ISVs) can build extensions/Software as a Service (SaaS) on SSIS Everest



Traditional ETL



1. SSIS on Prem (to SQL Svr)



- *Lift & Shift w/ compatibility*

2. SSIS on IaaS (to SQL on IaaS | Az DB)



- *Want PaaS benefits (scale, no VM mgmt, etc)*
- *Mix reuse & modernization*

3. SSIS Runtime in ADF

Trigger



SSIS Package A

SSIS Package “As-is”



Traditional ETL



1. SSIS on Prem (to SQL Svr)



- *Lift & Shift w/ compatibility*

2. SSIS on IaaS (to SQL on IaaS | Az DB)



- *Want PaaS benefits (scale, no VM mgmt, etc)*
- *Mix reuse & modernization*

3. SSIS Runtime in ADF

Trigger →

SSIS Package A

SSIS Package “As-is”

Trigger →

Copy

→ SSIS Pkg A

→ AML

→ ...

SSIS Package “Reuse”

Trigger →

Copy

→ Spark

→ AML

→ ...

Transform “@Scale”



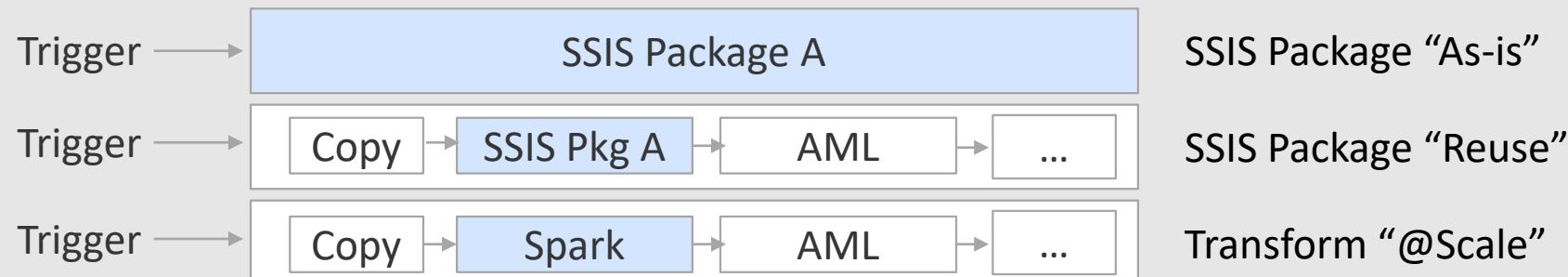
Traditional ETL

- ↓
1. SSIS on Prem (to SQL Svr)
 - *Lift & Shift w/ compatibility*
 2. SSIS on IaaS (to SQL on IaaS | Az DB)
 - *Want PaaS benefits (scale, no VM mgmt, etc)*
 - *Mix reuse & modernization*
 3. SSIS Runtime in ADF



Modern DW & Data Driven SaaS

- ↓
1. ADF V1 (Time-series, Tumbling Window)
(scenarios: log processing, etc.)
 2. ADF v2 (PaaS: Control Flow, Data Flow)
 - *Migrate v1 to v2:*
 - *Much more flexible model*
 - *Similar or cheaper in many scenarios*



ADF V2 Pricing (Preview Prices)

Orchestration

Units: Activity Runs

Activity Run: Single run or re-run of an activity or trigger

Where the activity is orchestrated	Price
Azure Integration Runtime (Public Network environment)	0 – 50K runs: \$.55/1k runs 50K+ runs: \$.50/1k runs
Self-Hosted Integration Runtime (Private Network environment)	\$.75/1K runs

Data Movement

Units: DMU

DMU: Unit of measure used to copy data from source to sink

	Price
Azure Integration Runtime (Public Network environment)	\$.125/hr
Self-Hosted Integration Runtime (Private Network environment)	\$.05/hr

Inactive Pipelines (Pipelines not associated with a trigger with zero runs for a week) : \$.20 / week

Azure Integration Runtime for SSIS

VM	Price for Standard Edition \$/hr	Price for Enterprise Edition \$/hr
A4 v2 (4core)	\$0.420	\$0.956
A8 v2 (8core)	\$0.862	\$1.935
D1 v2 (1core)	\$0.296	\$0.832
D2 v2 (2core)	\$0.397	\$0.933
D3 v2 (4core)	\$0.599	\$1.136
D4 v2 (8core)	\$1.199	\$2.271

The ADF Azure-SSIS IR is a fully-managed service to execute your SSIS packages. You can specify the VM type and the number of nodes for your dedicated IR environment. Note that you have to pay for a SQL Azure DB instance separately in order to host the SSIS catalog in the cloud. If you are using the Azure Data Factory to move data outside Azure network, you are required to pay for the amount of the data egress.

Note that these VM prices are discounted rates from the list prices of SQL Server VMs.

Enterprise SKUs for SSIS will not be available during public preview.

Scenario 1: Data movement across different data sources

Suppose you have a pipeline that has:

1. Copy Activity 1 moves data from On-Premises SQL Server to Azure Blob to be consumed later by a customized application
2. Copy Activity 2 moves data from Azure Blob to Azure SQL Database

Assume that the copy from On-Premises SQL Server to Azure Blob takes 2 hours. And the copy from Azure Blob to Azure SQL Database takes 1 hour with 2 DMUs. The pipeline is executed 30 times in a month.

Copy from On-Premises SQL Server to Azure Blob:

Data Movement: 30 activity runs X 2 hours duration for every run X \$0.05 = \$3

Activity Runs on Self-Hosted Integration Runtime (30 runs) + \$.75

Subtotal = \$3.75

Copy from Azure Blob to Azure SQL Database:

30 activity runs X 1 hour duration for every run X 2 DMUs x \$.125 = \$7.50

Activity Runs on Azure Integration Runtime (30 runs) + \$0.55

Subtotal = \$8.05

Total Price (per month) \$3.75 + \$8.05 = \$11.80

Scenario 2: Orchestrating data transformation activities on compute services

Suppose you have a pipeline that has:

1. Spark Activity to run Spark application on an Azure HDInsight cluster in Azure Virtual Network
2. Azure Data Lake Analytics U-SQL Activity to execute U-SQL on Azure Data Lake Analytics

Assume the pipeline was executed for 30 times in one month.

Spark Activity Data Movement \$0

Activity Runs on Self-Hosted Integration Runtime (30 runs) + \$.75

Subtotal = \$.75

Azure Data Lake Analytics U-SQL Activity Data Movement \$0

Activity Runs on Azure Integration Runtime (30 runs) + \$.55

Subtotal = \$.55

Total Price (per month) \$1.30

Scenario 3: Lift & Shift SSIS Packages

Suppose you have set up a SQL Server Standard and SSIS on a local 4-cores physical machine to run your daily ETL workload. Your ETL workload includes 100 SSIS packages and each of these packages takes about 5 mins to run. These packages are executed once daily and it requires about 8 hours each day to complete all the packages.

In order to lift and shift your SSIS package execution to Azure Data Factory, you need to provision the SSIS dedicated Integration Runtime via the Azure Data Factory portal

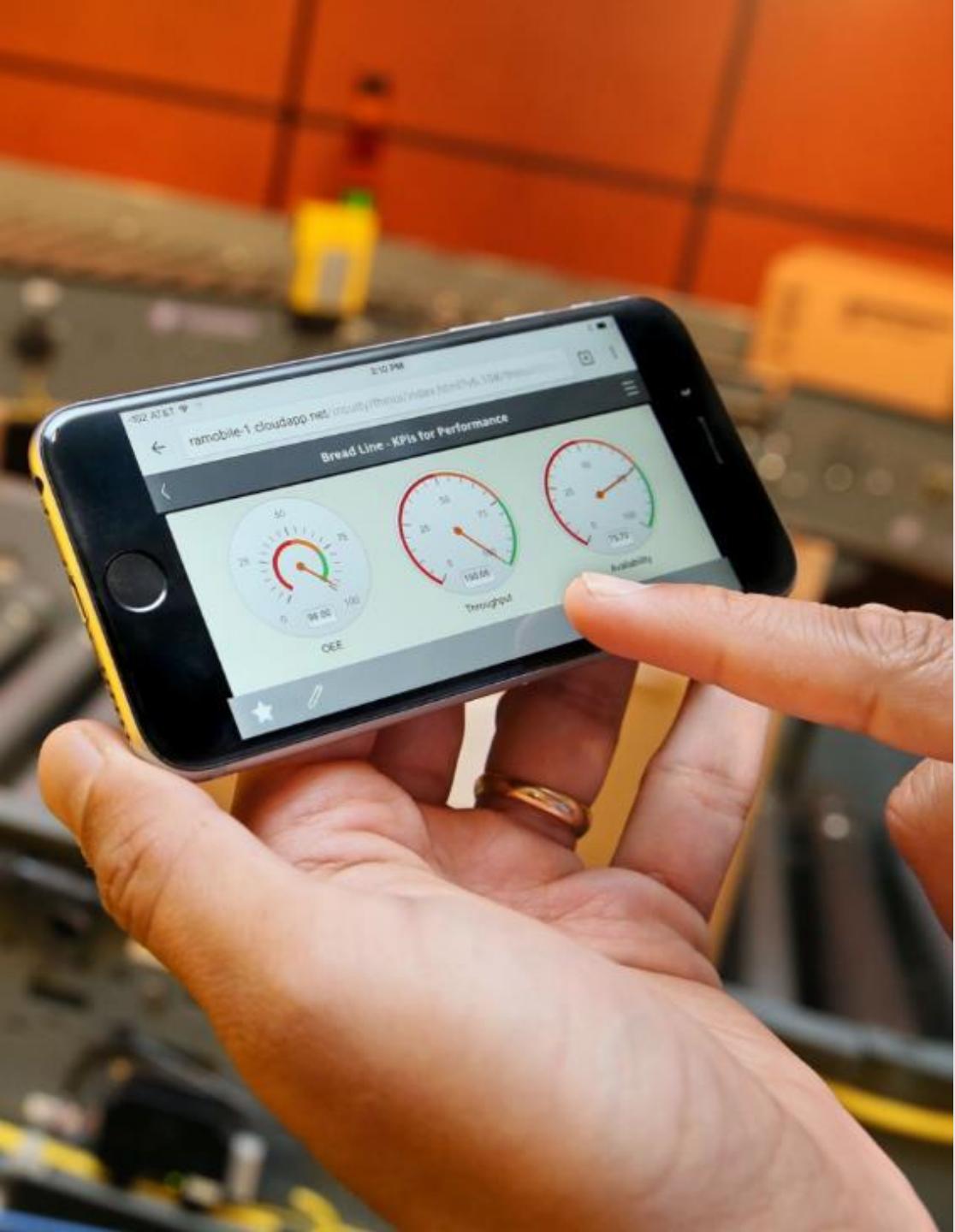
During the provision, you need to select the VM Type that match closely to the compute power you use on-premises. In this example, we will choose the A4 v2 type and 1 node to match the compute power on-premises. You are also required to create a SQL Azure Database instance for hosting the SSIS catalog on Azure.

We choose the S0 Standard for this example. SSIS Integration Runtime is a dedicated pool so the compute resource is dedicated only for you (not shared with other user). You can choose to keep the dedicated integration runtime pool 24/7 available or you can choose to spin it up and down

Integration Runtime Dedicated for 24/7	Cost
SSIS Integration Runtime – A4 v2	\$0.42 / hr
X 1 node	
X 24 hours / day (dedicated pool)	
X 31 days / month	
SSIS Integration Runtime 24/7 cost =	\$312.48 / month
SQL Azure for SSIS Catalog – S0	\$0.02 / hr
X 24 hours / day	
X 31 days / month	
SQL Azure for SSIS Catalog Cost =	\$15.03 / month
Total monthly cost to run SSIS ETL workload	\$327.51 / month

Integration Runtime Dedicated for the execution time only	Cost
SSIS Integration Runtime – A4 v2	\$0.42 / hr
X 1 node	
X 8 hours / day (dedicated pool)	
X 31 days / month	
SSIS Integration Runtime 24/7 cost =	\$104.16 / month
SQL Azure for SSIS Catalog – S0	\$0.0202 / hr
X 24 hours / day	
X 31 days / month	
SQL Azure for SSIS Catalog Cost =	\$15.03 / month
ADF to orchestrate the start / stop SSIS Integration runtime	
2 custom activities – start / stop IR	2 activity runs / day
X 31 days / month	
Total ADF Activity Cost	\$0.55 / month for 1 – 1000 runs
Total monthly cost to run SSIS ETL workload	\$119.74 / month

Customer Evidence



Industrial automation firm cuts up to 90 percent of costs

Challenge

Manage data in excess of the existing data warehouse and systems capacity

Integrate data on-premises with cloud and batch processing with live streaming

Monitor a complex supply chain – undersea wells, refining facilities, gas stations

Solution

Integrated complex data systems with Azure Data Factory

Composed, orchestrated, and managed highly available data pipelines

Reduced up to 90 percent of costs



Actuarial firm unlocks the potential of human capital to drive business

Challenge

Increasing complexity of insurance products and regulations

Diversion from firm's core purpose

Rise in on-premises hardware investments to handle high-performance grid computing

Solution

Built software-as-a-service on Azure to address modeling and reporting of large, compute-intensive workloads and used Azure Data Factory to automate extract, transform, and load (ETL) processes

Gained real-time insight into financial data



Improved Health Outcomes by integrating healthcare data

Challenge

Integration of Healthcare data from multiple hospitals and heterogenous sources

Transformation and orchestration of dataflow at scale

Increasing costs of scaling on-premise data pipelines

Solution

Built end-to-end data integration with Azure Data Factory, leveraging big data and compute offerings in Azure

Enabled data ingestion from hybrid sources of data

Orchestrated complex data workflows at scale

