



Enterprise Cloud Data Platforms

with Microsoft Azure & Cortana Analytical Suite

Khalid M. Salama, Ph.D.
Business Insights & Analytics
Hitachi Consulting UK

Outline

- Microsoft Azure and Cloud Computing
- Azure Data Architecture and Services
- EDW with Azure SQL Data Warehouse
- Big Data Storage and Analytics on Azure Data Lake and HDInsight
- Data Processing and Movements with Azure Data Factory
- Real-time Processing with Azure Event Hubs and Stream Analytics
- Data Science and Machine Learning on Microsoft Azure
- Data Discoverability with Azure Data Catalog

Microsoft Azure and Cloud Computing

Cloud Computing

Enabling the future of computing

On-demand Self-service

- Web interfaces that deliver capacity in seconds or minutes

Broad Network Access

- Support many devices, such as mobile phones, laptops, tablets and workstations

Resource Pooling

- Thousands to hundred of thousands servers are immediately usable

Rapid Elasticity

- On-demand or automatic (threshold-based) scale-up or down.

Measured Service

- Pay as you use, what you use, Improve operational metric and cost accounting

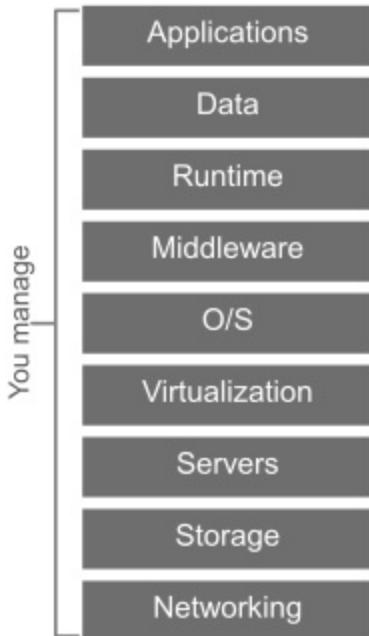
Innovation

- Modern IT capabilities and services provided frequently

Cloud Computing on Microsoft Azure

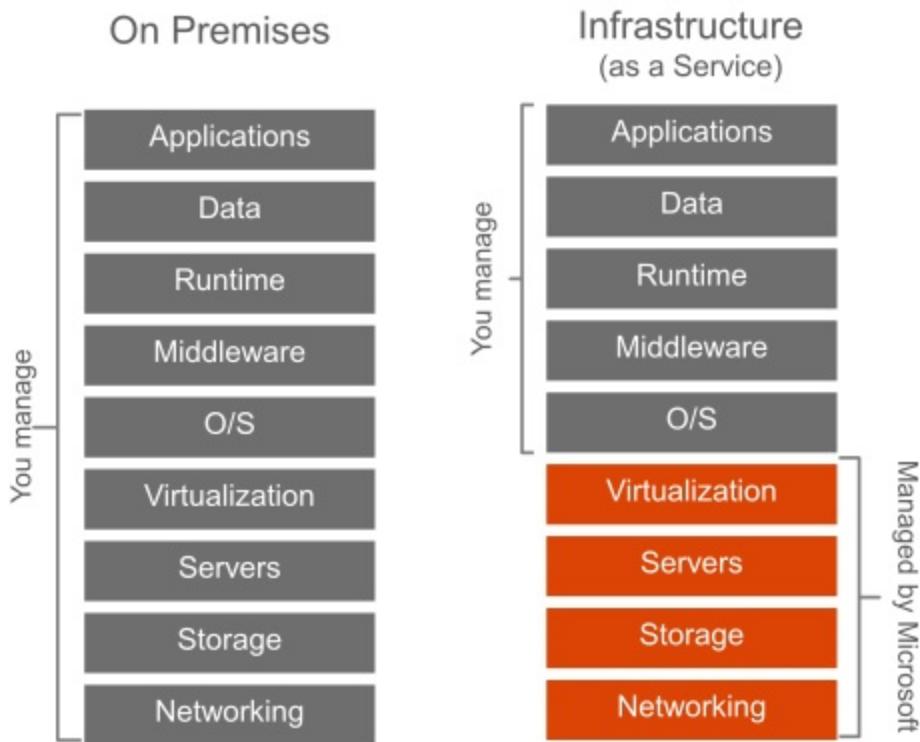
Computing Models

On Premises



Cloud Computing on Microsoft Azure

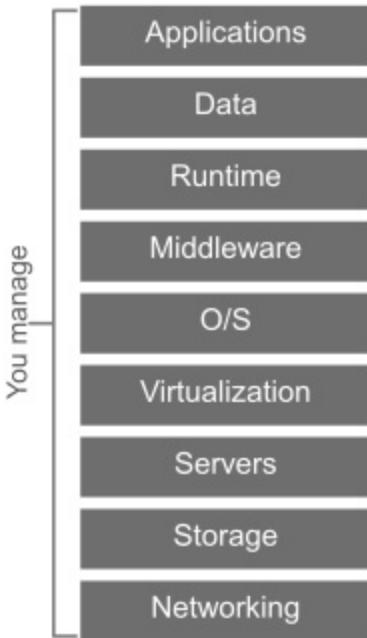
Computing Models



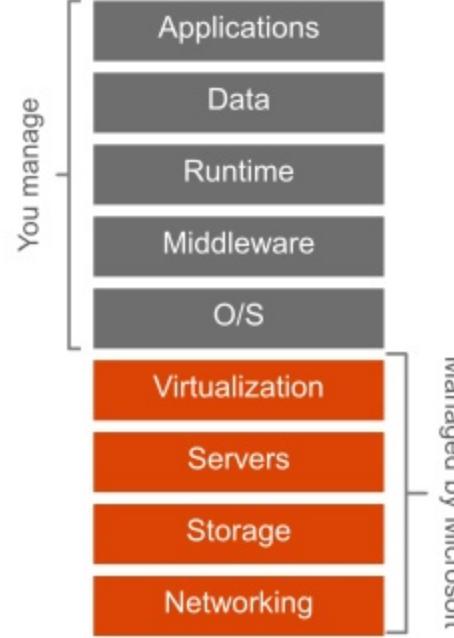
Cloud Computing on Microsoft Azure

Computing Models

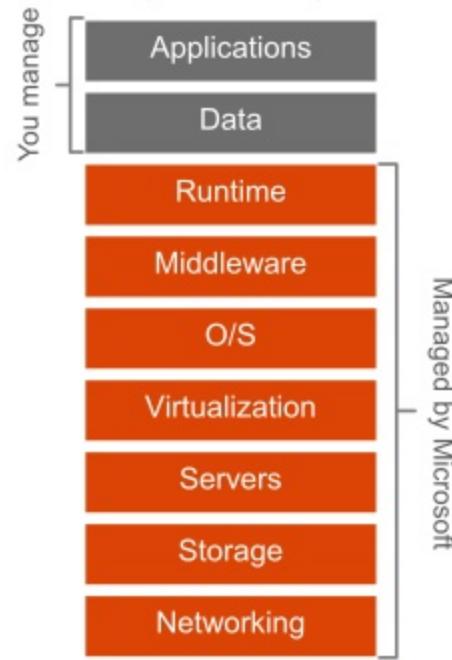
On Premises



Infrastructure (as a Service)

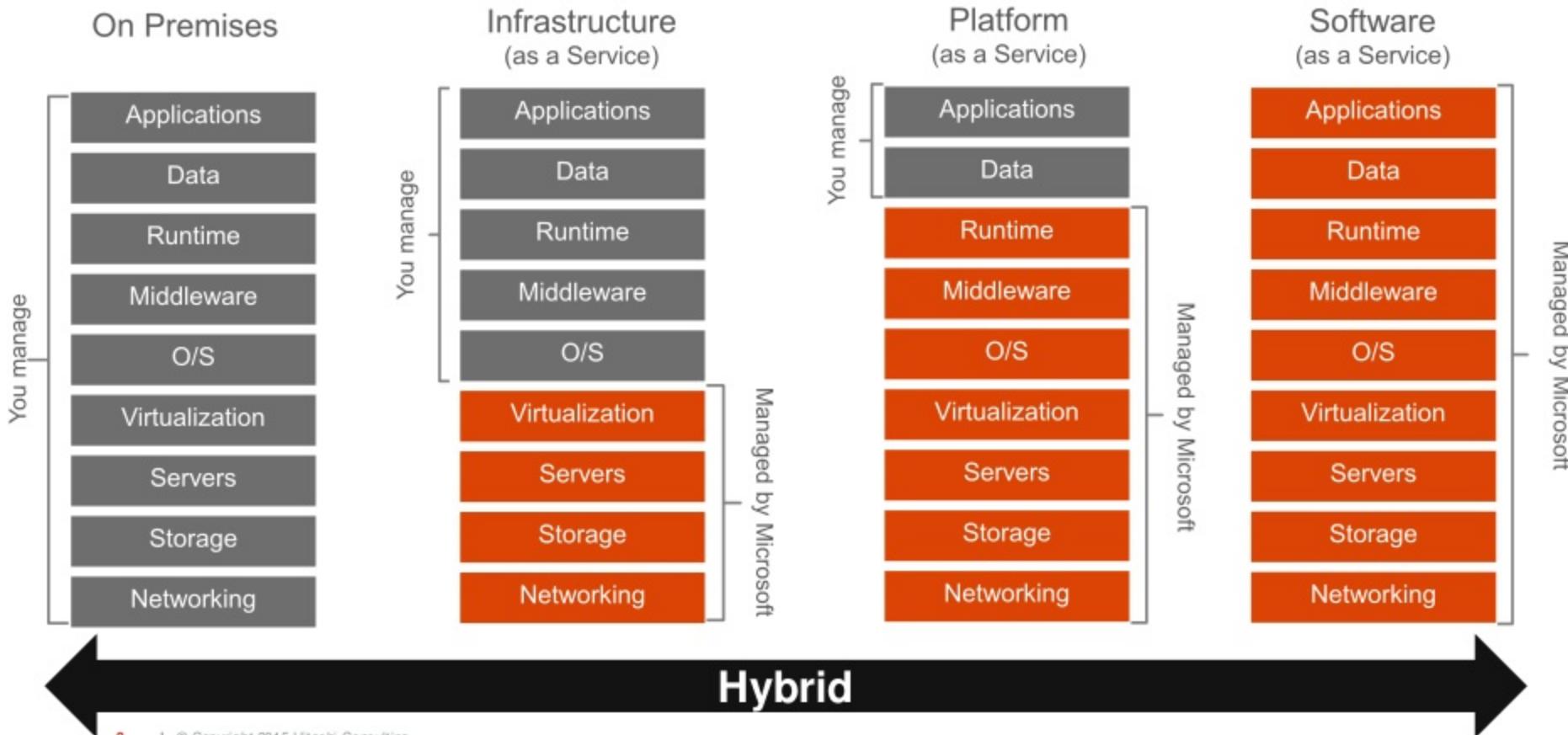


Platform (as a Service)



Cloud Computing on Microsoft Azure

Computing Models



Cloud Computing on Microsoft Azure

Computing Models - Examples

- **Infrastructure as a Service (IaaS)** – VMs and Networking Services
(build Software infrastructure, develop apps, use apps, Software maintained infrastructure)
- **Platform as a Service (PaaS)** – Data and App Services
(develop apps, use apps)
- **Software as a Service (SaaS)** – SharePoint Online, Office 365, Microsoft Dynamics, Visual Studio Online, Azure Cognitive Services, etc.
(use apps)

Cloud Computing on Microsoft Azure

So why the cloud?

Agility

- Configure and provision in minutes
- Automatic/on-demand scaling up/down

Economics

- Pay as you use (100% Utilization)
- Elastic capacity

Innovation

- Connectivity
- Innovative capabilities and services
- New compute models

Azure Data Architecture and Services

Enterprise Data Platforms

Challenges of Enterprise Information Management

Dispersed data sources

Data Explosion
(Big Data/ IoT)

Absence of the single version of truth

Limited automation for data acquisition

No conformed business definitions

Limited data discoverability

More demand on Data Science and Advanced Analytics capabilities

Pragmatic business vs dogmatic IT dilemma in building data products

(quick & dirty added-value vs. enterprise class solution)

Need for both batch and real-time data processing

Lack of integration with operational systems

Self-serviceability, innovation, and enrichment

Source Systems



LOB CRM



ERP OLTP



Web logs



Devices



Sensors



Social



Clickstream

Azure Enterprise Data Platform

Source Systems



LOB CRM



ERP OLTP



Web logs



Devices



Sensors



Social

Clickstream

Azure Enterprise Data Platform

Data Lake

<Blob Storage/
Data Lake>

- Raw Data Files
- Unknown Value



Filter,
Reform,
Transpose
Extract, Load

Source Systems



LOB CRM



ERP

OLTP



Web logs



Devices



Sensors



Social

Clickstream

Azure Enterprise Data Platform

Data Lake

<Blob Storage/
Data Lake>

- Raw Data Files
- Unknown Value



Filter,
Reform,
Transpose
Extract, Load

Enterprise Data Warehouse

<Azure SQL DW>

- Single Version of the truth
- High Value Data
- Integrated Data store
- Unified Business Rules
- Conformed Data Model
- Historical, Non-volatile
- Optimized for Batch Processing

MPP

Processing &
Transformation

Aggregate/
Calculate

Source Systems



LOB CRM



ERP OLTP



Web logs



Devices



Sensors



Social



Clickstream

Azure Enterprise Data Platform

Data Lake

<Blob Storage/
Data Lake>

- Raw Data Files
- Unknown Value

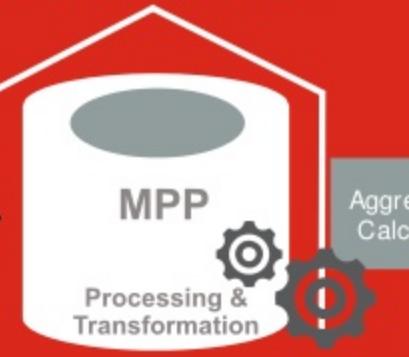


Filter,
Reform,
Transpose
Extract, Load

Enterprise Data Warehouse

<Azure SQL DW>

- Single Version of the truth
- High Value Data
- Integrated Data store
- Unified Business Rules
- Conformed Data Model
- Historical, Non-volatile
- Optimized for Batch Processing



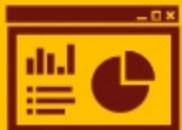
Information Marts

<Azure SQL Database/
SQL Server 2016 VM>

- Subject Oriented
- Aggregated/ Summarized
- Simple and Intuitive Model for Analysts
- Optimized for Interactive Querying
- Based on the Single Version of the truth



Front-End Analytics



Power BI
Dashboards



Power Pivot/
Power View/
Power Query



SQL Server
Analysis Services

Source Systems



LOB CRM

ERP OLTP

Web logs



Devices

Sensors



Clickstream

Azure Enterprise Data Platform

High Performance Computing

- Azure Batch

Data Orchestration & Movement

- Azure Data Factory

Data Lake

<Blob Storage/
Data Lake>

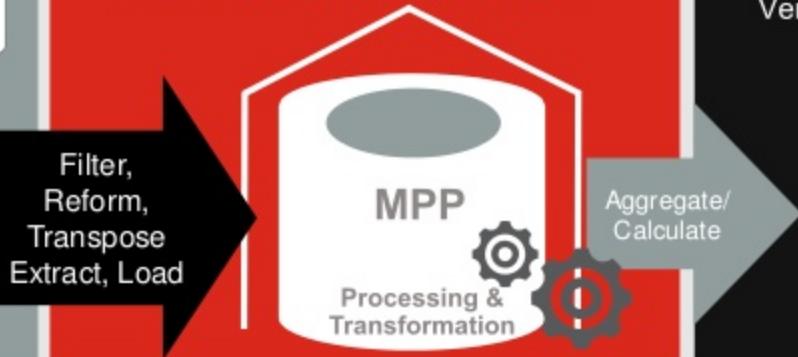
- Raw Data Files
- Unknown Value



Enterprise Data Warehouse

<Azure SQL DW>

- Single Version of the truth
- High Value Data
- Integrated Data store
- Unified Business Rules
- Conformed Data Model
- Historical, Non-volatile
- Optimized for Batch Processing



Information Marts

<Azure SQL Database/
SQL Server 2016 VM>

- Subject Oriented
- Aggregated/ Summarized
- Simple and Intuitive Model for Analysts
- Optimized for Interactive Querying
- Based on the Single Version of the truth



Front-End Analytics



Power BI
Dashboards



Power Pivot/
Power View/
Power Query



SQL Server
Analysis Services

Big Data Processing

- Hive/Pig
- U-SQL
- Spark
- Storm



Data Science & Machine Learning

- Azure Machine Learning
- Microsoft R Server
- Spark ML
- Azure Cognitive Services



Real-time Stream Processing

- Stream Analytics
- Event Hubs
- Storm
- Spark Streaming



NoSQL Data Stores

- Table Storage
- HBase
- DocumentDB



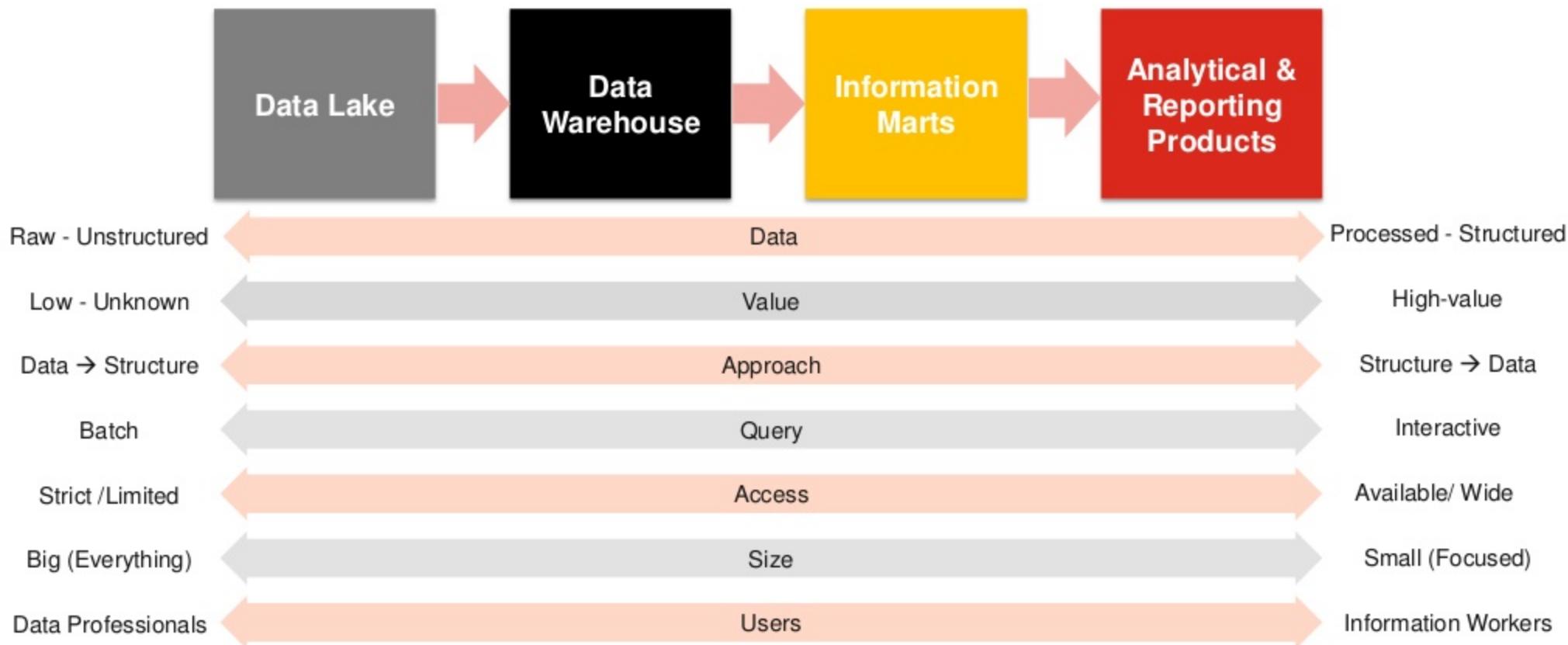
Data Discovery & Search

- Azure Data Catalog
- Azure Search



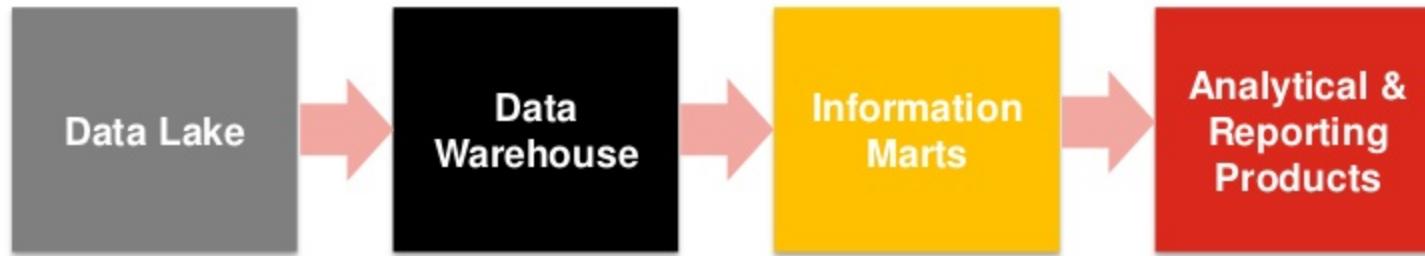
Azure Enterprise Data Platform

Layers and relevance



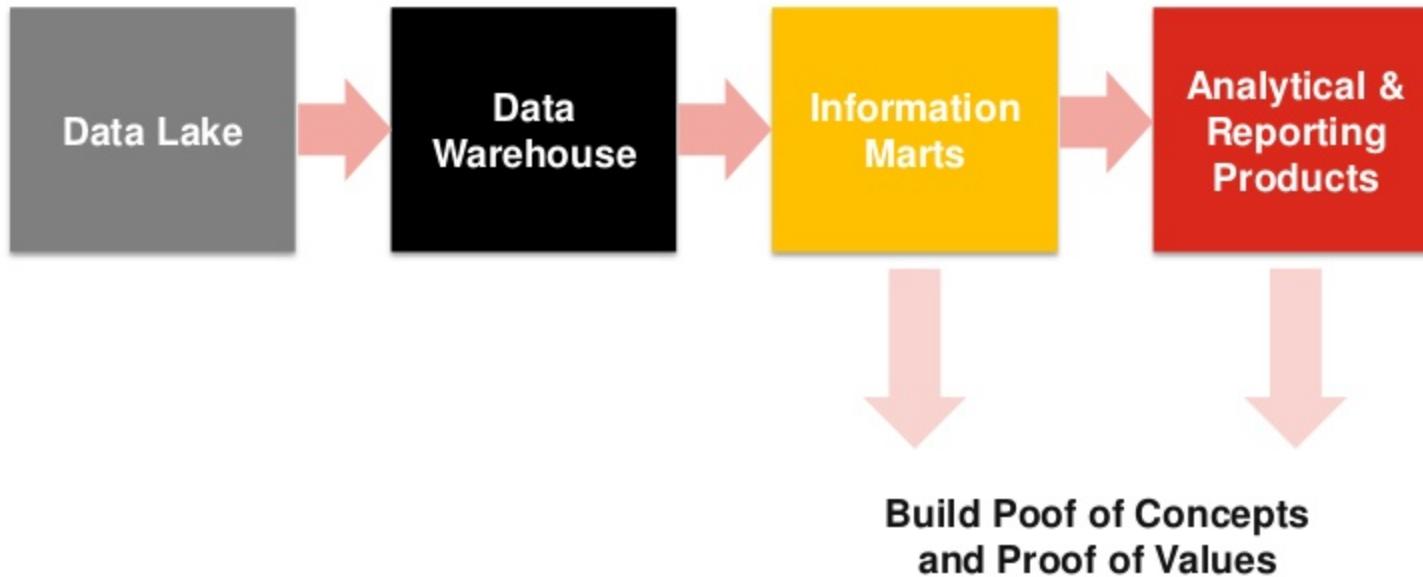
Azure Enterprise Data Platform

Business Agility & IT Governance



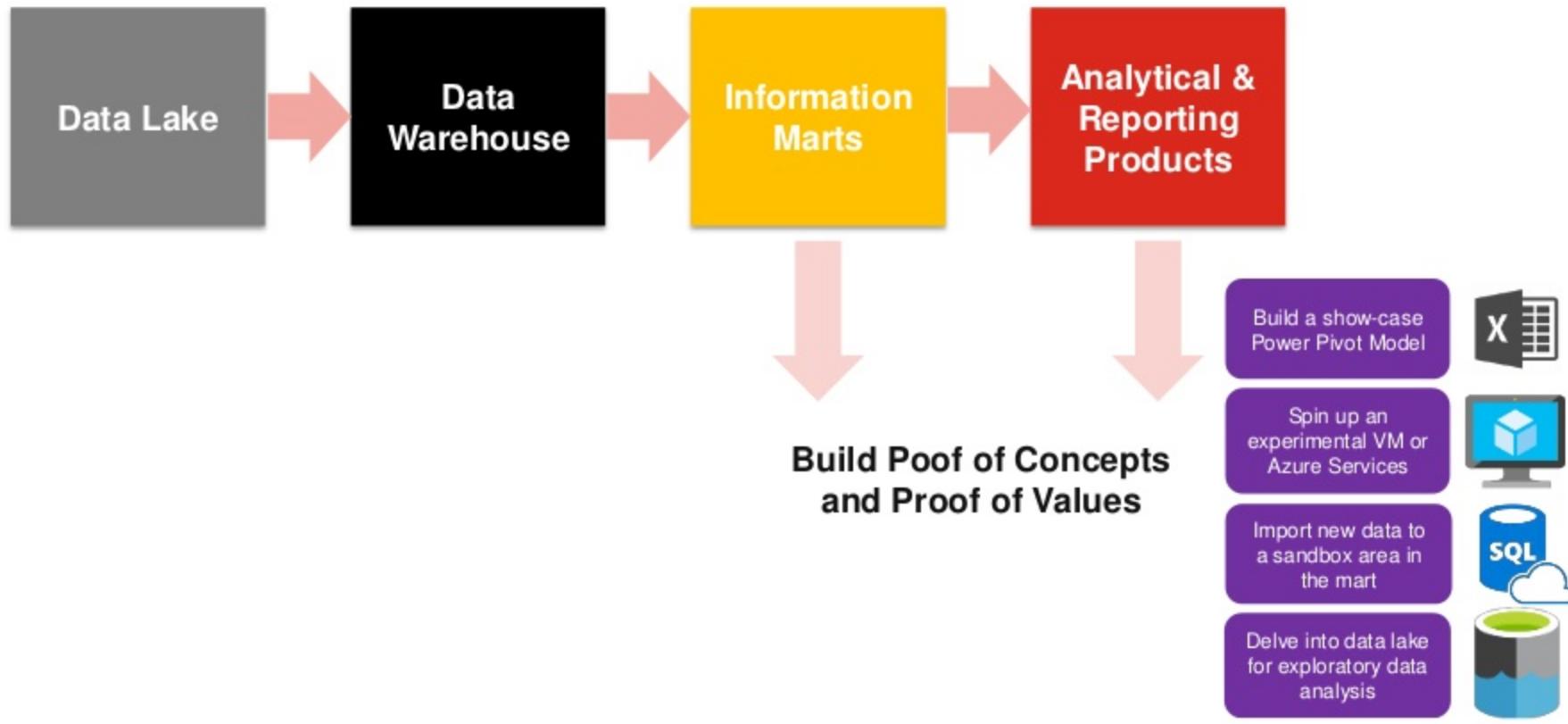
Azure Enterprise Data Platform

Business Agility & IT Governance



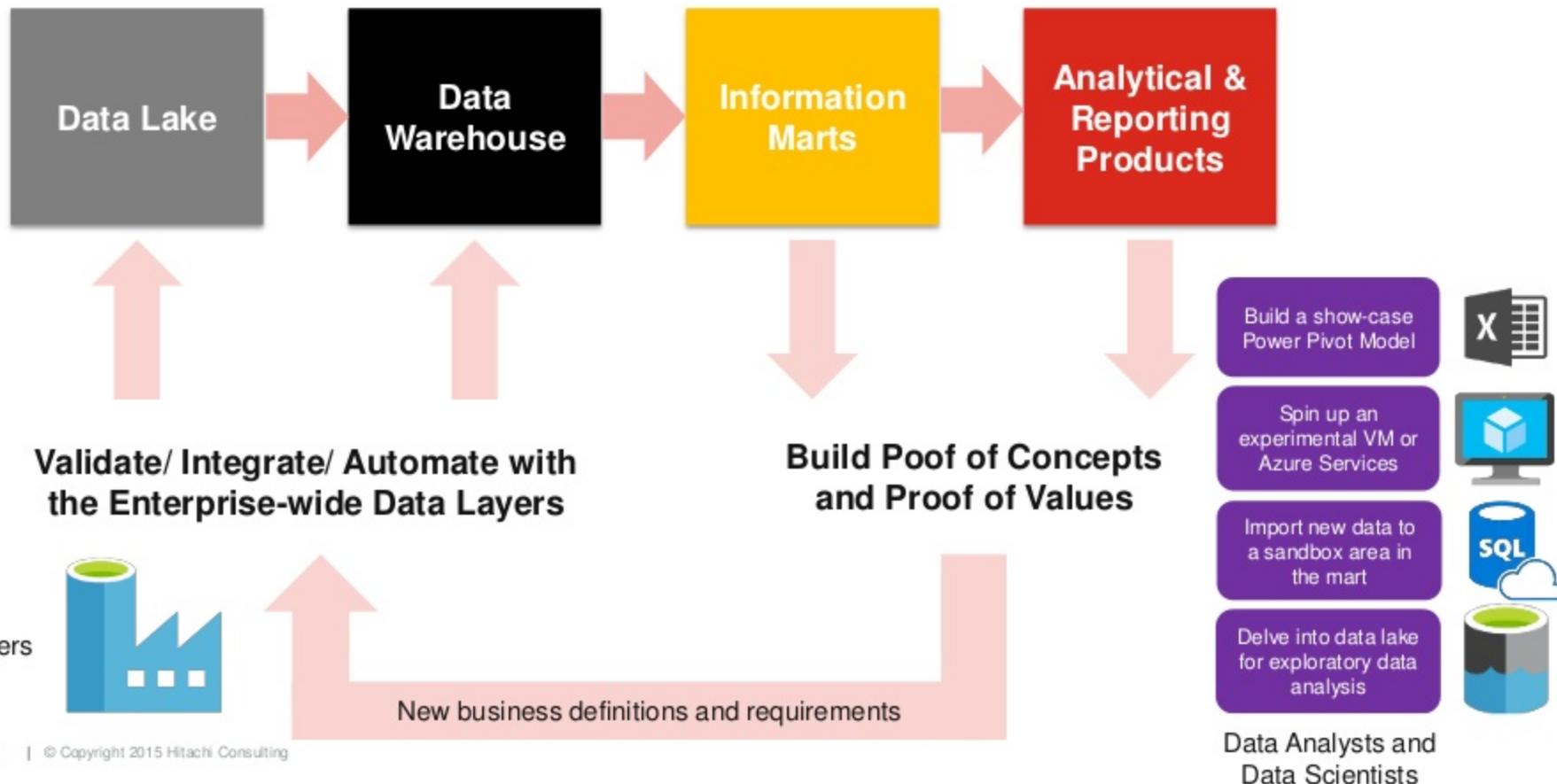
Azure Enterprise Data Platform

Business Agility & IT Governance



Azure Enterprise Data Platform

Business Agility & IT Governance



Azure Data Lake Storage and Analytics

Data Lake – The Concept

What is that? And Why it is needed?

What is a Data lake

A repository that holds raw data file extracts of all the Enterprise source systems

Usually a Distributed File System (HDFS, Blob Storage, Amazon S3, GFS, etc.) to support Big Data Processing

Pentaho (Hitachi Data Systems) co-founder and CTO, James Dixon coined the term ‘Data Lake’

Why a Data Lake

Single and cheap repository of Enterprise-wide (raw) data

Abstracts your DW/ Business Intelligence Solutions from source systems

Supports potential Big Data Processing and Analytics

Data Lake – The Concept

Data Processing

Data stored in its native form (Structured, Semi-structured, and Unstructured):
Text Files, JSON, XML, Logs, Mainframe, etc.

Big Data technologies, such as MapReduce, Pig, Hive, Spark, and U-SQL, are used
to process the data in a procedural fashion:

Transpose/ reshape data
into a tabular form

Extract/filter valuable
and relevant data
elements

Load prepared datasets
into a relational EDW



Azure Data Lake

Azure Data Lake Services

HDInsight

- Hadoop Cluster as a Service
- MapReduce | Pig | Hive | Spark
- HBase | Oozie | Storm
- Unit of Cost: Running nodes



Data Lake Analytics

- Big Data Jobs as a Service
- U-SQL
- Unit of Cost: running Job



Data Lake Storage

- Built for data files, rather than Blobs
- WebHDFS REST APIs for Hadoop integration
- Optimized for analytical workloads
- Unlimited Storage, petabyte files



Azure Data Lake Analytics

U-SQL

- C# meets SQL (.NET data type, C# expressions + SQL set statements)
- Natively Parallel
- Set-based and procedural processing
- Works with structured and unstructured data
- Extensible (User-defined Functions, Operators and Aggregate Functions)
- Reuse of custom build .NET assemblies

New U-SQL Job

PREVIEW

Submit Job Data Explorer Import File

* Job Name: Summarise Products Priority: 1000 Parallelism: 1

```
1 @searchlog =
2     EXTRACT ProductKey Int32,
3         OrderDateKey Int32,
4             PromotionKey Int32,
5                 CurrencyKey Int32,
6                     SalesTerritoryKey Int32,
7                         SalesOrderNumber String,
8                             SalesOrderLineNumber Int16,
9                                 RevisionNumber Int16,
10                                OrderQuantity Int16,
11                                    SalesAmount Decimal,
12                                        ShipDate DateTime,
13                                            filename String
14    FROM "wasb://csvuploads/prspocsw.blob.core.windows.net/{filename:}.csv"
15    USING Extractors.Csv();
16
17 @productsummary =
18     SELECT ProductKey,
19         CurrencyKey,
20             SUM(SalesAmount) AS SalesTotal,
21                 ShipDate < Today()? "Shipped" : "Outstanding"
22     FROM @searchlog
23
24 OUTPUT @productsummary
25     TO "/AdventureWorks/AdventureWorksSummary.csv"
26 USING Outputters.Csv();
```

Azure HDInsight

Hadoop on Microsoft Cloud

Applications



Yet Another Resource Negotiator (YARN)

Windows Azure Blob Storage | Data Lake (DFS)

Azure SQL Data Warehouse

Azure SQL Data Warehouse

Enterprise Data Warehouse – why?

Azure SQL Data Warehouse

Enterprise Data Warehouse – why?

“That is great! Now we have a data lake that has all the data that we want. Let’s spin up a VM, build a database, and produce the reports we want.”

– Said Mr. quick-win from the business.

Azure SQL Data Warehouse

Enterprise Data Warehouse – why?

“That is great! Now we have a data lake that has all the data that we want. Let’s spin up a VM, build a database, and produce the reports we want.”

– Said Mr. quick-win from the business.

“No. The data lake has raw data with no business structure. High-value data needs to be loaded into a canonical relational form, where enterprise-wide business entities , calculations, and rules are defined. Building reporting products on top of the data lake will results in unreliable information silos without single version of the truth.”

– Said Mr. enterprise data architect from the IS.

Azure SQL Data Warehouse

A Parallel DW on the cloud

An elastic, Massively Parallel Processing (MPP), cloud-based data warehouse as a service, with enterprise-class features



Suitable for

- Batch data processing workloads
- Consolidate high-value data into a single enterprise-wide structured store
- Model, transform and aggregate data
- Perform query analysis across large

Not Suitable for

- Operational workloads
- High frequency reads & writes
- Singleton operations → (Information Mart)
- Row by row processing needs → (Data Lake)
- High concurrency → (Information Mart)

Massively Parallel Processing (MPP)

Scale-out

Distributed
Queries

Azure SQL Data Warehouse

MPP - Logical Overview

Storage

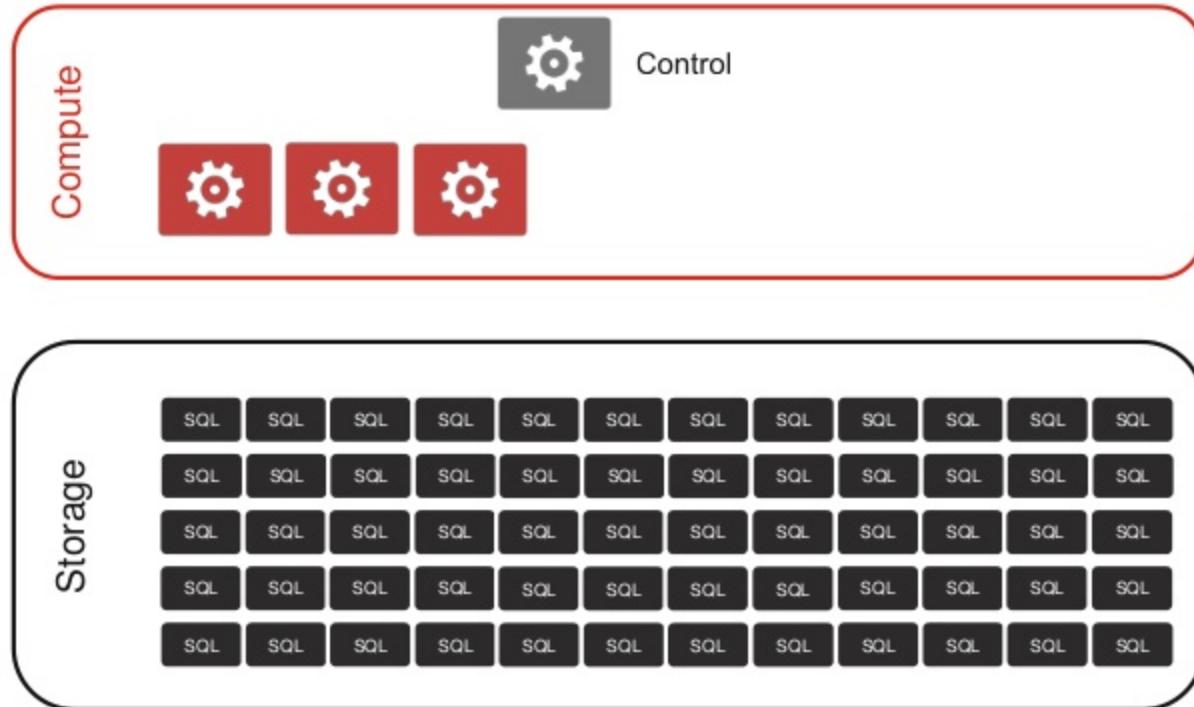
SQL											
SQL											
SQL											
SQL											
SQL											

- 60 Data Distribution Nodes
- Each node contains a portion of data (distribution) for parallel processing
- Storage is de-coupled of compute

Azure SQL Data Warehouse

MPP - Logical Overview

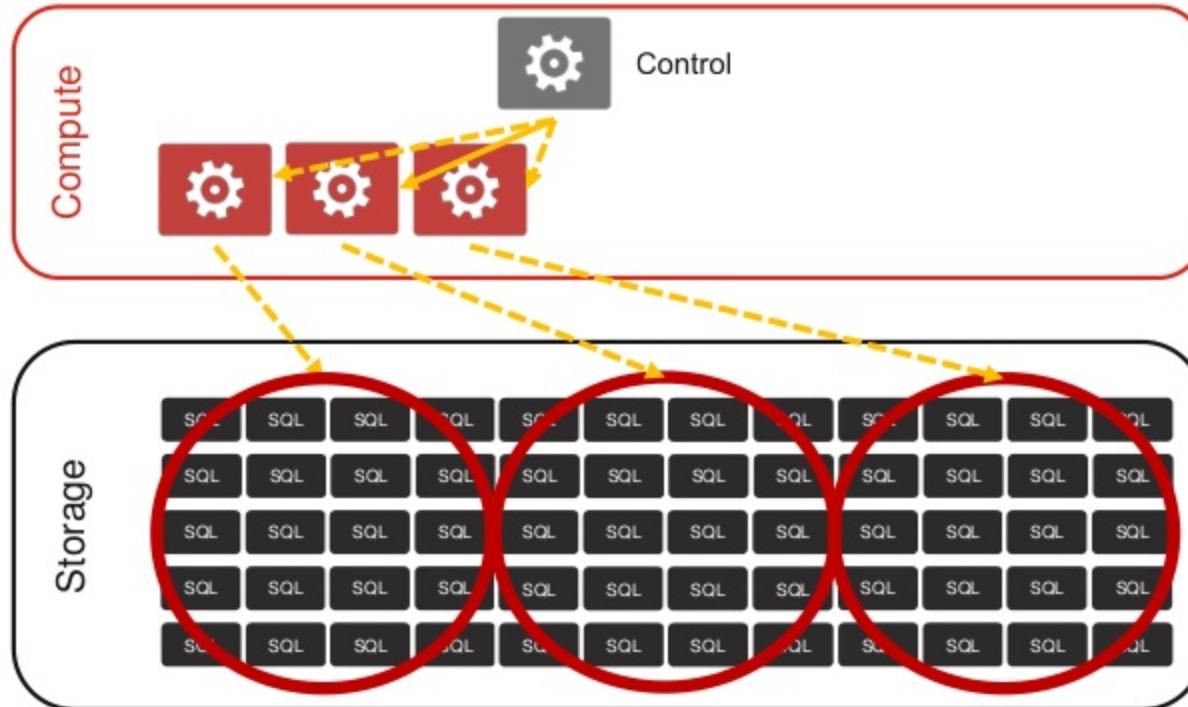
- A scale-out distributed query engine



Azure SQL Data Warehouse

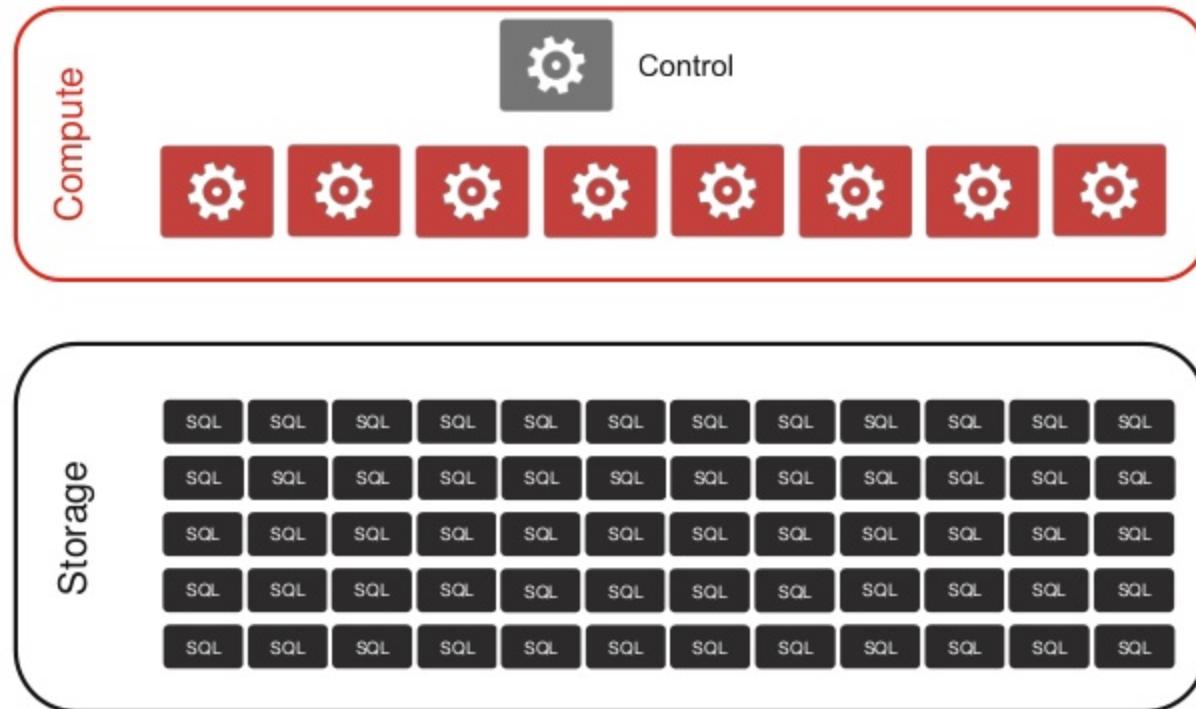
MPP - Logical Overview

- A scale-out distributed query engine



Azure SQL Data Warehouse

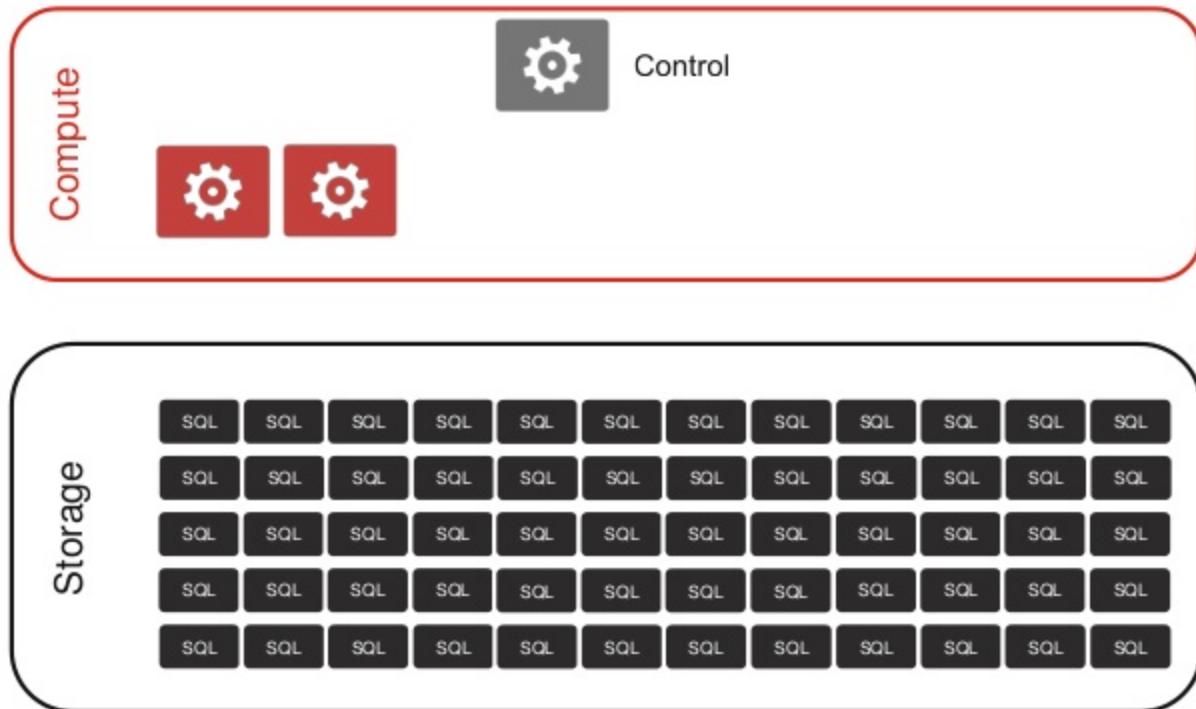
MPP - Logical Overview



- A scale-out distributed query engine
- You get more compute capacity by increasing Data Warehouse Units (DWUs)

Azure SQL Data Warehouse

MPP - Logical Overview

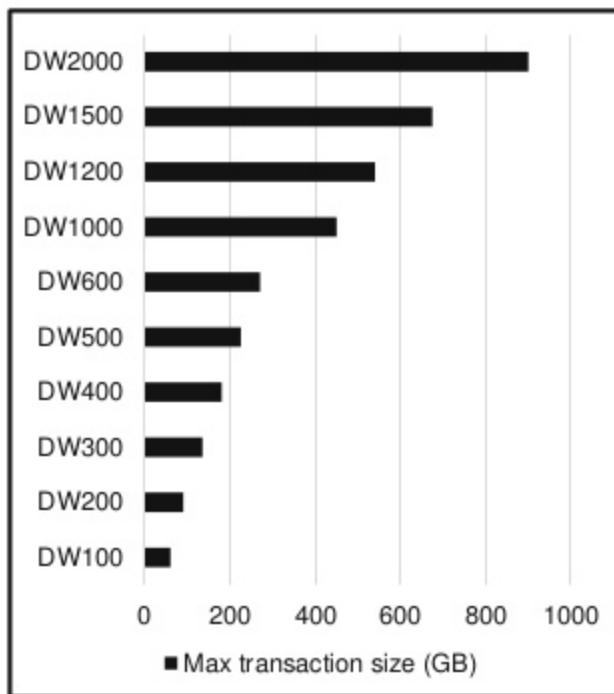


- A scale-out distributed query engine
- You get more compute capacity by increasing Data Warehouse Units (DWUs)
- Elastic – scale up/down as needed, and pay per usage (hourly-level)
- Scaling takes minutes as – compute is independent of data, thus no data re-distribution is needed.

Azure SQL Data Warehouse

Data Warehouse Units

RAM

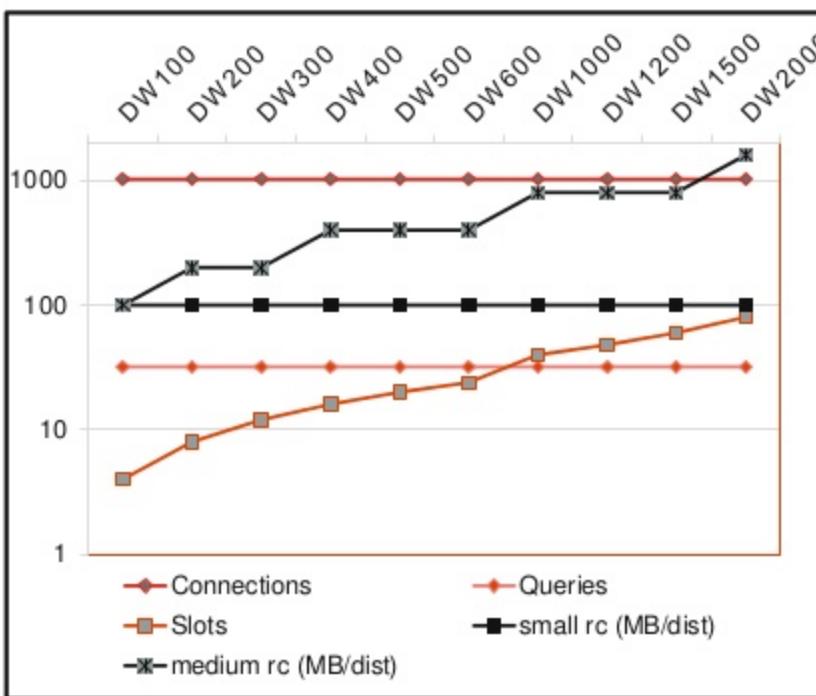
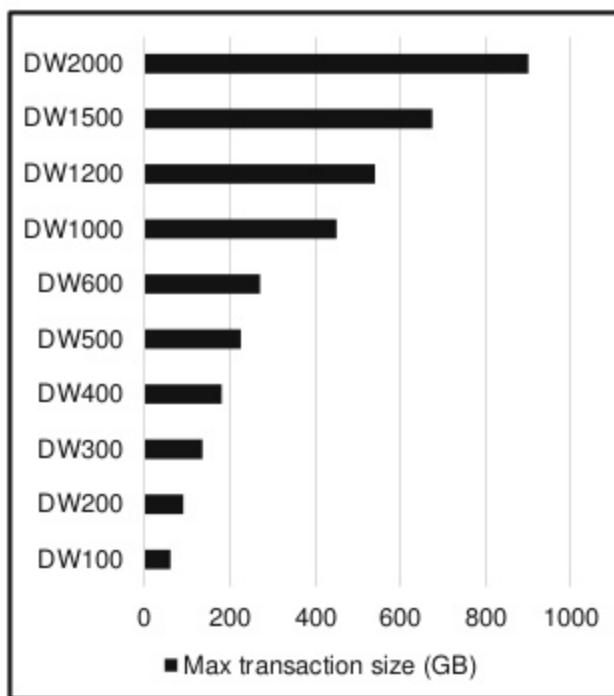


Azure SQL Data Warehouse

Data Warehouse Units

RAM

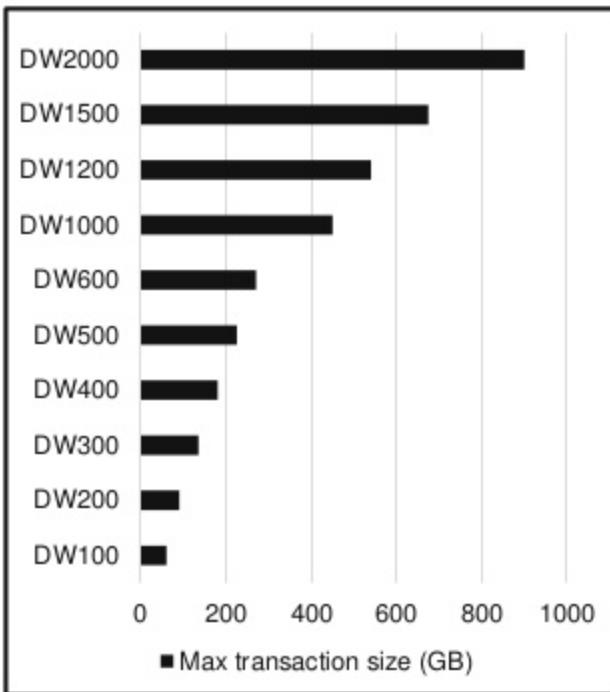
CPU



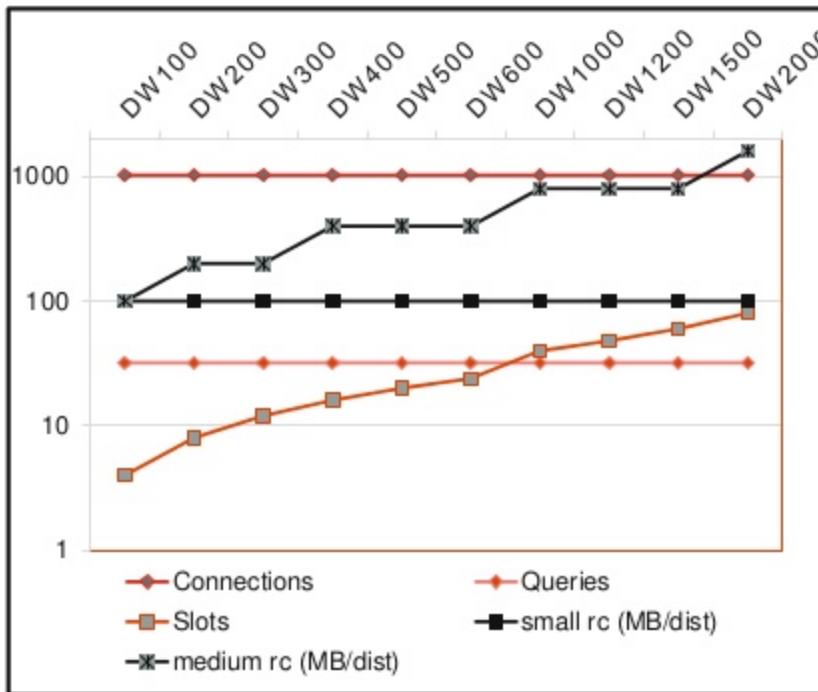
Azure SQL Data Warehouse

Data Warehouse Units

RAM



CPU



I/O

DWU	Readers	Writers
DW100	8	60
DW200	16	60
DW300	24	60
DW400	32	60
DW500	40	60
DW600	48	60
DW1000+	60	60

Azure SQL Data Warehouse

Distribution vs. Partitioning

Distribution

- Aim to avoid **skewing** to maximize **balance parallel processing**.
 - if most of your data in one node, no parallel processing is gained.
- Aim to avoid **expensive data shuffling** across data nodes.
- Consider columns used in **JOIN** or **GROUP BY** for distribution keys to reduce shuffling
- Consider columns with **high cardinality** for distribution keys to avoid skewing
- **Round-Robin** distribution is the default - and a good compromise.
- Use **Replication** for small tables (yet to be available).

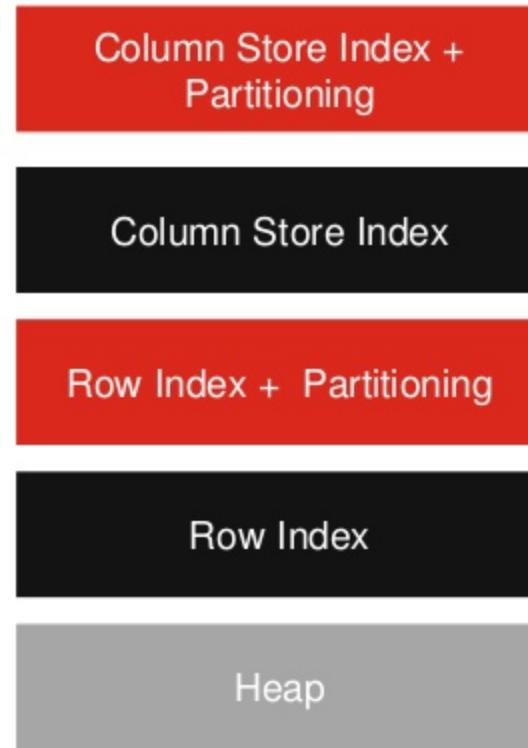
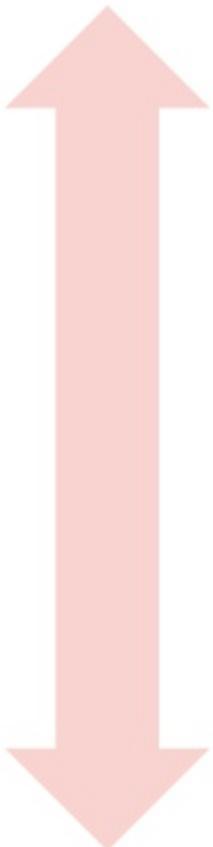
Partitioning

- In a given data node, aim to minimize the amount of data to be considered in the query
 - that is, **partition elimination**.
- Consider columns used in **WHERE** clauses for partition keys.

Azure SQL Data Warehouse

Indexing

Better query performance



Used in tables expecting aggregation queries

Used in staging/data processing tables

Better load performance

Used in data landing tables

Azure SQL Data Warehouse

Data Model

- Capture elementary source data structure, main business entities and relationships
- Atomic, fine-grained, integrated, detail-oriented, historical, bi-temporal
- Provides a flexible, yet unified, data foundation allow agile business rules to construct new data products and information models.

Data Vault
(Hubs, Links, Satellites)

Dan Linstedt

3rd Normal form
(3NF)

Bill Inmon

Dimensional Model
(Facts, Dimension)

Ralph Kimball

Azure SQL Data Warehouse

PolyBase

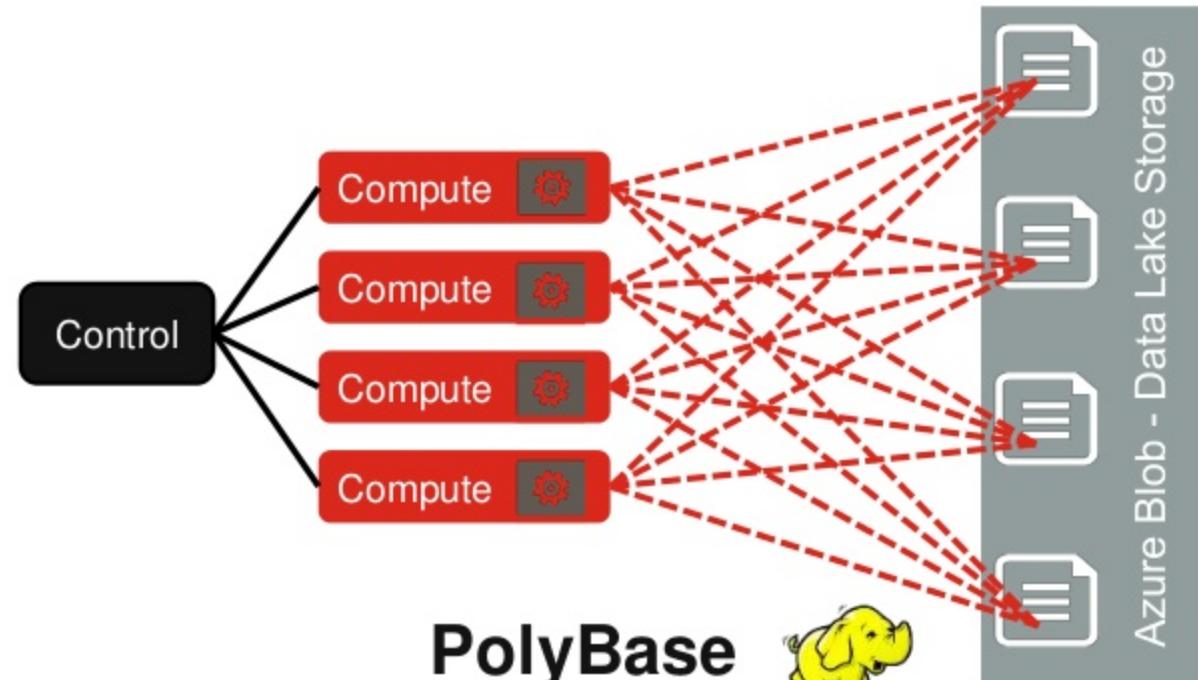
Hadoop-base technology
(very similar to Hive)

Allow bringing data in/out Azure
SQL DW/SQL from/to Azure Blob
Storage/Data Lake

Create External table
(Schema-on-Read)
– Identify file location and format

Parallel Data Movement – Direct
connectivity between compute
and data files (back door)

Best way to get data in and out of
Azure SQL DW



Azure SQL Data Warehouse

PolyBase

```

CREATE EXTERNAL DATA SOURCE MyAzureDataSource
WITH
    (TYPE          = HADOOP
    ,LOCATION      =
        'wasbs://[container]@accountname.blob.core.windows.net/path'
    );
    
```



```

CREATE EXTERNAL FILE FORMAT MyTextFileFormat
WITH
    (FORMAT_TYPE      = DELIMITEDTEXT
    ,FORMAT_OPTIONS
        (
            FIELD_TERMINATOR = '|',
            STRING_DELIMITER = ',',
            DATE_FORMAT= 'yyyy-MM-dd',
            USE_TYPE_DEFAULT= TRUE
        )
    ,DATA_COMPRESSION =
        | 'org.apache.hadoop.io.compress.GzipCodec'
    );
    
```

```

CREATE EXTERNAL TABLE [asb].[FactOnlineSales]
    ([ProductKey]      int      NOT NULL
    ,[StoreKey]        int      NOT NULL
    ,[DateKey]         int      NOT NULL
    ,[CustomerKey]     int      NOT NULL
    ,[PromotionKey]   int      NOT NULL
    ,[SalesQuantity]  int      NOT NULL
    ,[UnitPrice]       money    NOT NULL
    ,[SalesAmount]     money    NOT NULL
    )
    WITH
    (LOCATION='wasbs://filepath_or_directory'
    ,DATA_SOURCE      = MyAzureDataSource
    ,FILE_FORMAT      = MyTextFileFormat
    ,REJECT_TYPE      = VALUE
    ,REJECT_VALUE     = 0
    ,REJECT_SAMPLE_VALUE = 1000
    );
    
```

Azure SQL Data Warehouse

Data Processing

Avoid Singleton
Inserts, Deletes,
and Updates

Used Create Table
As Select (CTAS)
for processing
data

Distribution

Partitioning

Indexing

Create and Update
Statistics

Extract the data from the source, Load the data, then
perform the set-based Transformations as SQL
procedures to benefit from MPP - No SSIS data flow task!

Materialize external tables
before processing



Information Marts

Information Marts

Role in the Enterprise Data Platform

Subject-Oriented	Owned by business departments	Contains data aggregations and calculations to answer pre-defined business questions	Azure SQL Database Azure SQL 2016 VM Analysis Service Model
Dimensional Model or pre-canned datasets	May contain a sandbox area for experimentation		Data is processed on the MPP platform (Azure SQL DW) and then loaded in the mart
Built on top of the single version of the truth in the EDW		Suitable for interactive query and high concurrency	



End-to-end Data Processing

Data Lake
<Azure Data Lake | Blob Storage>



Enterprise DW
<Azure SQL DW>



Information Mart
<Azure SQL DB | SQL Server 2016>



Big Data Processing (TEL)
<Data Lake Analytics | HDInsight>

Massively Parallel Processing (ELT)
<PolyBase + T-SQL>

Populate (ETL)
<Azure Data Factory | SSIS>

Azure Data Factory

Azure Data Factory

Cloud data processing & movement services

A managed Azure cloud service for building & operating
data processing and movements

Scalable

Reliable

Pay-as-you use

Compose, monitor & schedule
data pipelines

Automatic cloud resource
management



Azure Data Factory

Concepts & Components

Data Factory



Linked Service



A connection information to a data store or a compute resource

Dataset



A pointer to the data object to be used as input or an output of an Activity

Pipeline



logical grouping of Activities with certain sequence & schedule

Azure Data Factory

Concepts & Components

Data Factory



Pipeline



Activity



actions to perform on data.

Input 0-N dataset(s) - Output 1-N dataset(s)

Linked Service



A connection information to a data store or a compute resource

Dataset



A pointer to the data object to be used as input or an output of an Activity

Pipeline



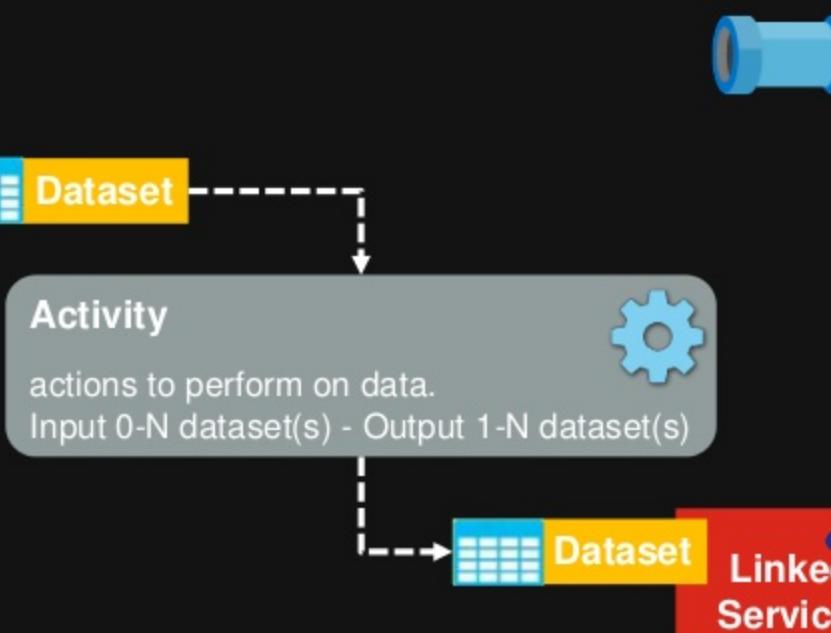
logical grouping of Activities with certain sequence & schedule

Azure Data Factory

Concepts & Components

Data Factory

Pipeline



Linked Service



A connection information to a data store or a compute resource



Dataset

A pointer to the data object to be used as input or an output of an Activity



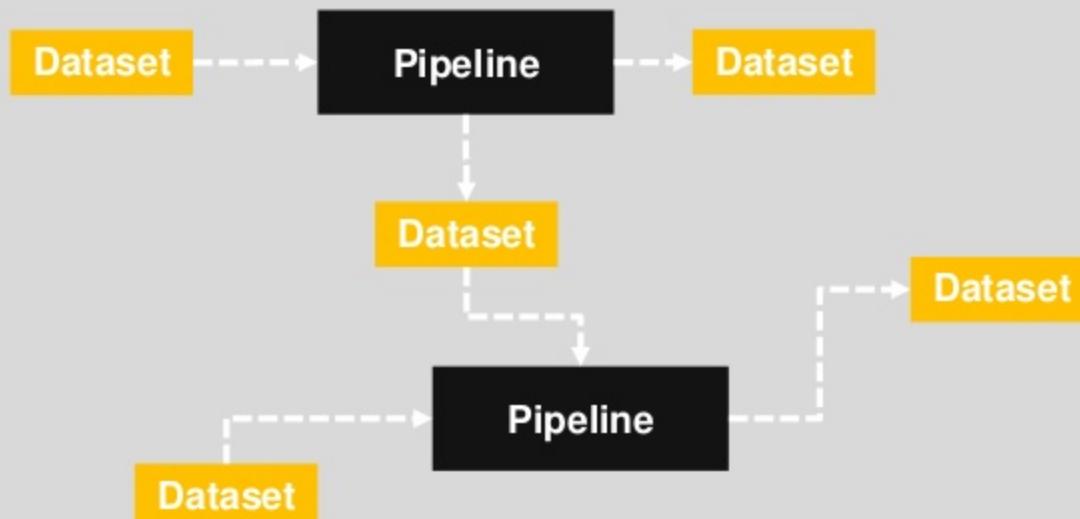
Pipeline

logical grouping of Activities with certain sequence & schedule

Azure Data Factory

Concepts & Components

Data Factory



Linked Service



A connection information to a data store or a compute resource

Dataset



A pointer to the data object to be used as input or an output of an Activity

Pipeline



logical grouping of Activities with certain sequence & schedule

Azure Data Factory

Concepts & Components



Azure Data Factory

Features

Linked Service – Data

- Azure Blob
- Azure Table
- Azure SQL Database
- Azure SQL Data Warehouse
- Azure DocumentDB
- Azure Data Lake Store
- SQL Server
- ODBC data sources
- Hadoop Distributed File System (HDFS)

Linked Service - Compute

- Azure HDInsight On-Demand
- Azure Batch
- Azure Machine Learning
- Azure Data Lake Analytics
- Azure SQL DW (StoredProcedure)
- Azure SQL DB (StoredProcedure)
- SQL Server (StoredProcedure)

Activities

- Data Movement
(copy Data from Source to sink)
- Data Transformation
 - Hive
 - Pig
 - MapReduce
 - Machine Learning
 - Stored Procedure
 - Data Lake U-SQL
 - Custom .NET

Data Gateway - Connect to an on-prem data source

Azure Data Factory

Example

Linked Service

```
{
  "name": "HDInsightOnDemandLinkedService",
  "properties": {
    "type": "HDInsightOnDemand",
    "typeProperties": {
      "version": "3.2",
      "clusterSize": 1,
      "timeToLive": "00:30:00",
      "linkedServiceName": "AzureStorageLinkedService1"
    }
  }
}
```

Input Dataset

```
{
  "name": "AzureBlobInput",
  "properties": {
    "type": "AzureBlob",
    "linkedServiceName": "AzureStorageLinkedService1",
    "typeProperties": {
      "fileName": "input.log",
      "folderPath": "adfgtstarted/inputdata",
      "format": {
        "type": "TextFormat",
        "columnDelimiter": ","
      },
      "availability": {
        "frequency": "Month",
        "interval": 1
      },
      "external": true,
      "policy": {}
    }
  }
}

{
  "name": "Azure Blob Output",
  "properties": {
    "type": "AzureBlob",
    "linkedServiceName": "Azure Storage Linked Service 1",
    "typeProperties": {
      "folderPath": "adfgtstarted/partitioneddata",
      "format": {
        "type": "TextFormat",
        "columnDelimiter": ","
      },
      "availability": {
        "frequency": "Month",
        "interval": 1
      }
    }
  }
}
```

Output Dataset

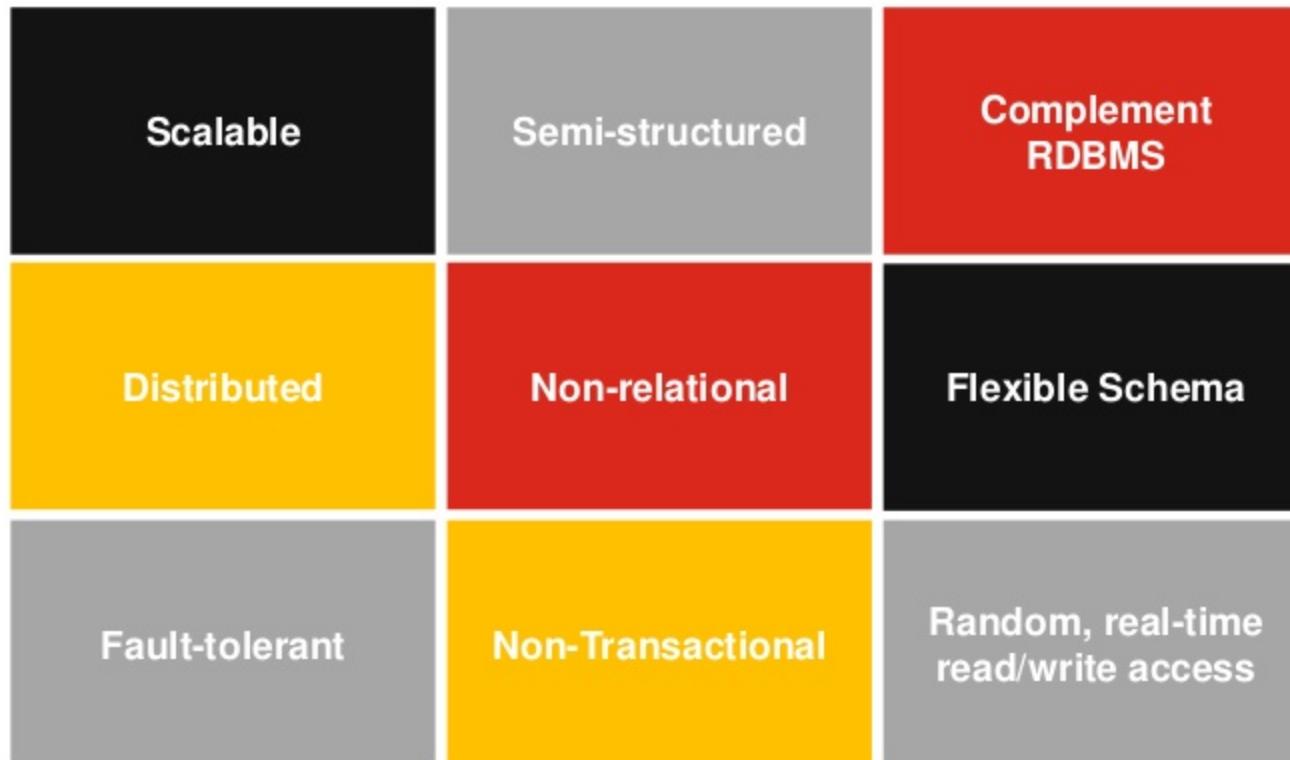
Pipeline

```
{
  "name": "MyFirstPipeline",
  "properties": {
    "description": "My first Azure Data Factory pipeline",
    "activities": [
      {
        "type": "HDInsightHive",
        "typeProperties": {
          "scriptPath": "adfgtstarted/script/partitionweblogs.hql",
          "scriptLinkedService": "AzureStorageLinkedService1",
          "defines": {
            "inputtable": "<storageaccounturl>",
            "partitionedtable": "<folderurl>"
          }
        },
        "inputs": [
          {
            "name": "AzureBlobInput"
          }
        ],
        "outputs": [
          {
            "name": "AzureBlobOutput"
          }
        ],
        "policy": {
          "concurrency": 1,
          "retry": 3
        },
        "scheduler": {
          "frequency": "Month",
          "interval": 1
        },
        "name": "Run Sample Hive Activity",
        "linkedServiceName": "HDInsightOnDemandLinkedService"
      }
    ],
    "start": "2016-04-01T00:00:00Z",
    "end": "2016-04-02T00:00:00Z",
    "isPaused": false
  }
}
```

Azure NoSQL Data Stores

What is NoSQL?

key attributes..



Why NoSQL?

Suitability

Suitable for

Random, real-time read/write access

Reference Data

Variable Data Structures

Singleton Select/ Insert/ update

Not Suitable for

Batch Processing

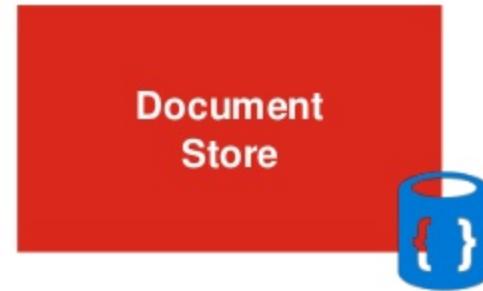
Complex Analytical Queries

Joins

Complex Transactions

NoSQL Data Stores

Categories and breads



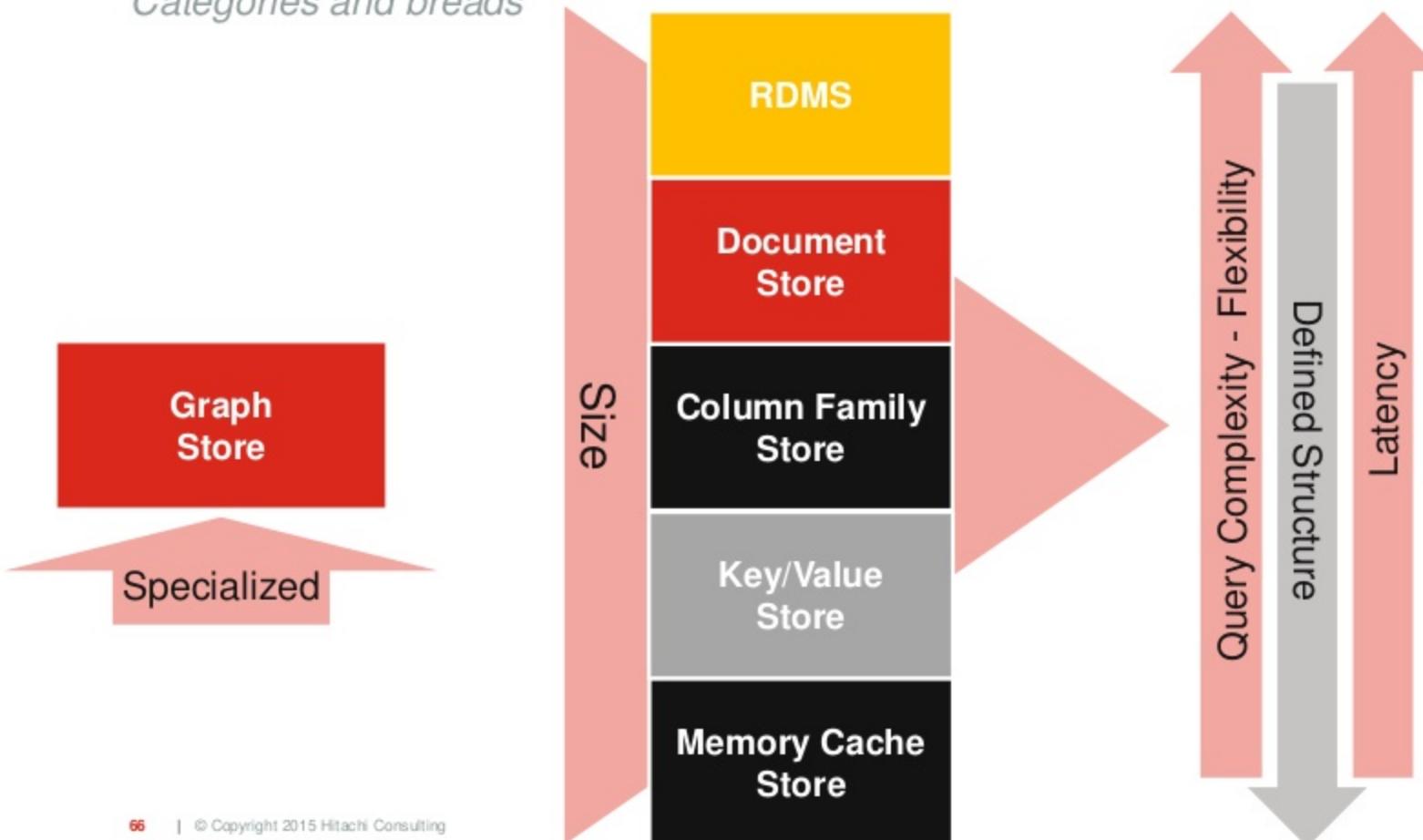
Others

Azure NoSQL Technologies

Key/Value Store	<ul style="list-style-type: none"> ▪ Flexible data structure ▪ Dictionary/ lookup ▪ Value can be anything 	Azure Table Storage
Column Family Store	<ul style="list-style-type: none"> ▪ Column-oriented access ▪ Big Data with real-time read/write random access ▪ Extensible 	 Azure Hbase on HDInsight
Document Store	<ul style="list-style-type: none"> ▪ Query-able data ▪ Objects (complex structure) in JSON, XML, etc. ▪ CRUD apps 	 Azure DocumentDB
Graph Store	<ul style="list-style-type: none"> ▪ Social networks ▪ Fraud detection ▪ Relationship-heavy data 	 Microsoft Graph Engine (Trinity)
Memory Cache Store	<ul style="list-style-type: none"> ▪ Non-durable data ▪ Fast access ▪ Hot data 	 Azure Redis Cache

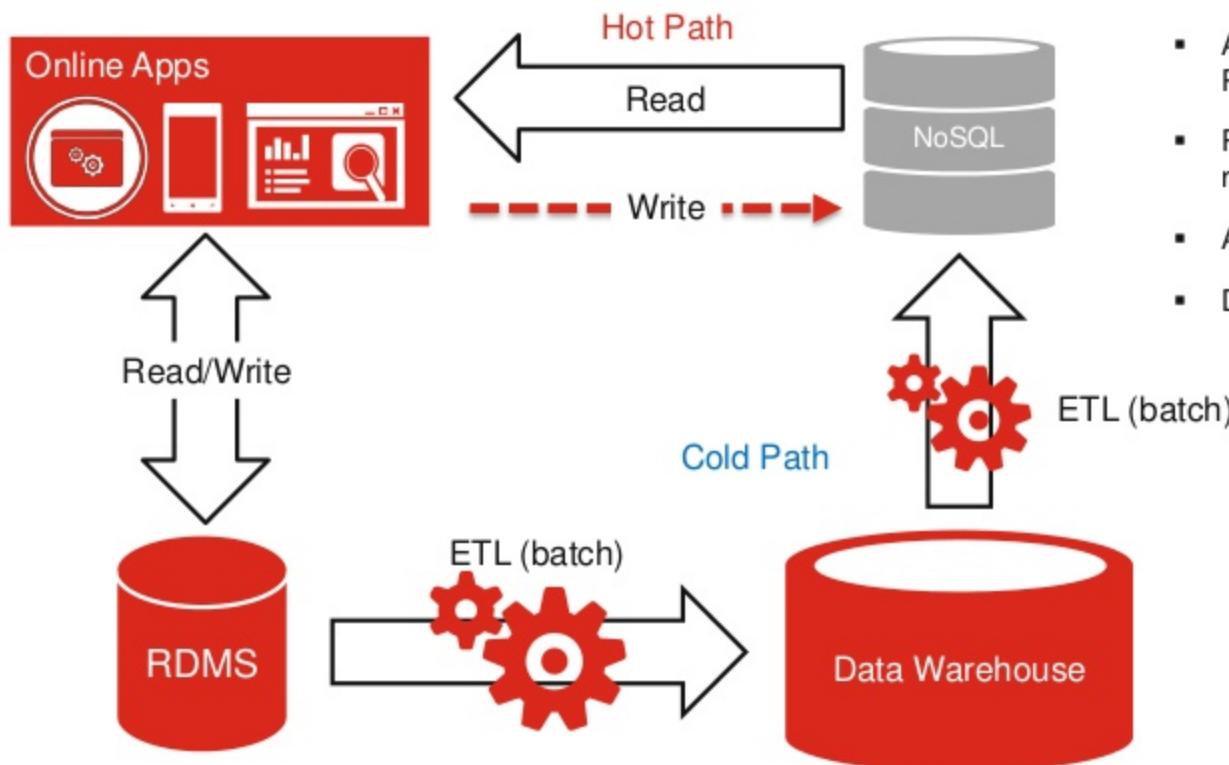
NoSQL Data Stores

Categories and breads



NoSQL Usage Patterns

A common scenario...



- Data is Extracted, Transformed, loaded from OLTP to EDW
- Aggregations, KPIs, and scores are computed via Batch Processing
- Results are populated to a NoSQL data store for reference use in apps
- App **hot read** - ETL **cold write**
- Document & Graph stores

E.g., Single Customer View:

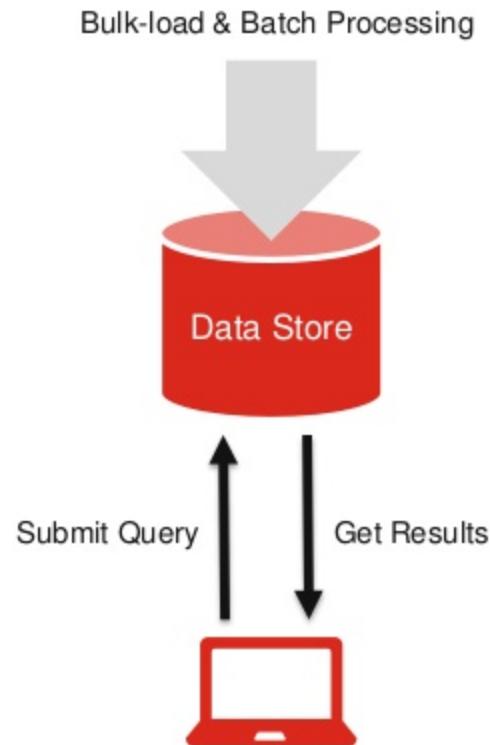
- Customer Matching, customer KPIs, segment assignment, and propensity scoring are performed as batch processing in DW
- The output goes to NoSQL to be used for real-time recommendation, campaigning, targeted advertising, etc.

Real-time stream processing on Microsoft Azure

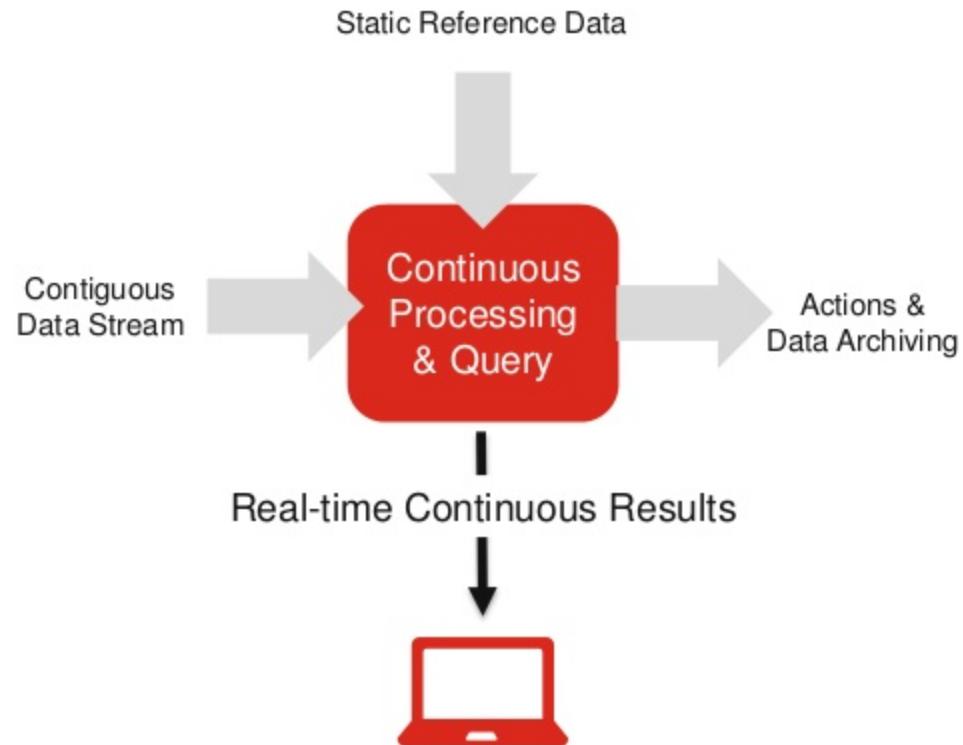
What is Event & Stream Processing?

Data at rest vs. Data in motion

Traditional – Working with data at rest

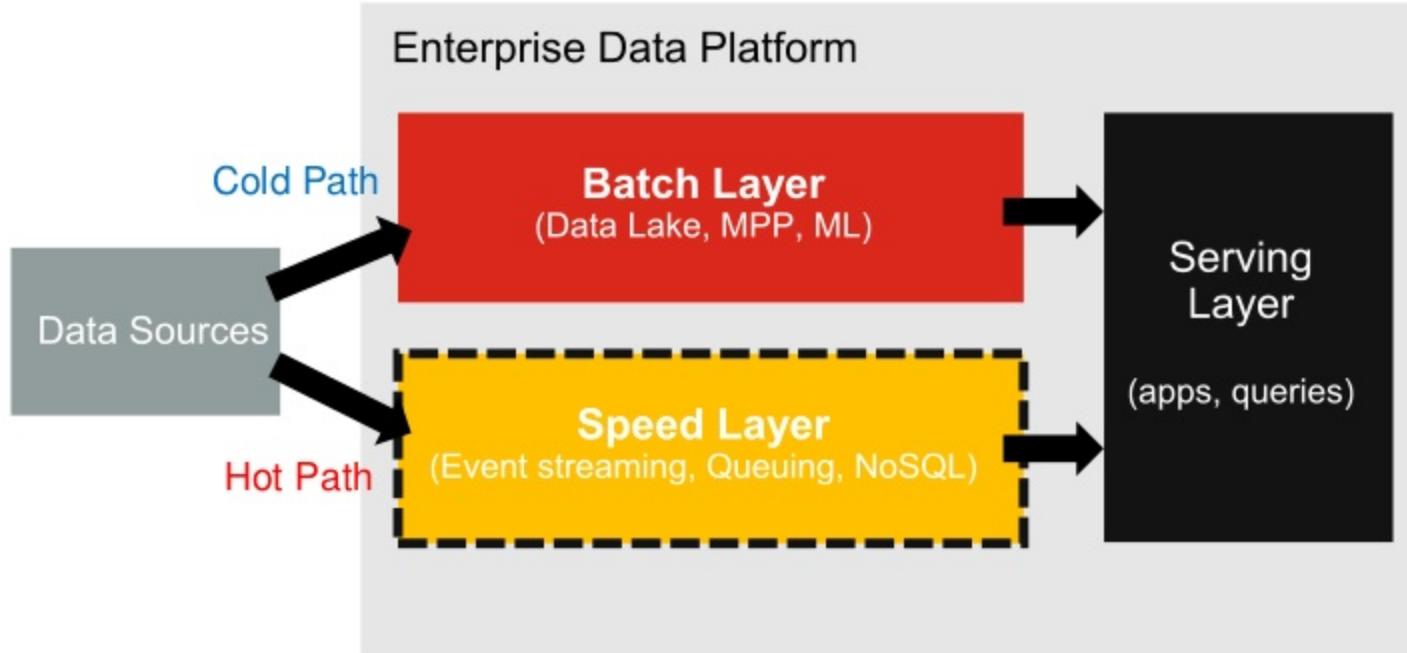


Real-time – Working with data in motion



Lambda Architecture

The speed layer and stream processing



Events & Stream Processing Architecture

Azure Tools & Technologies

Event
Triggers



Applications



Devices



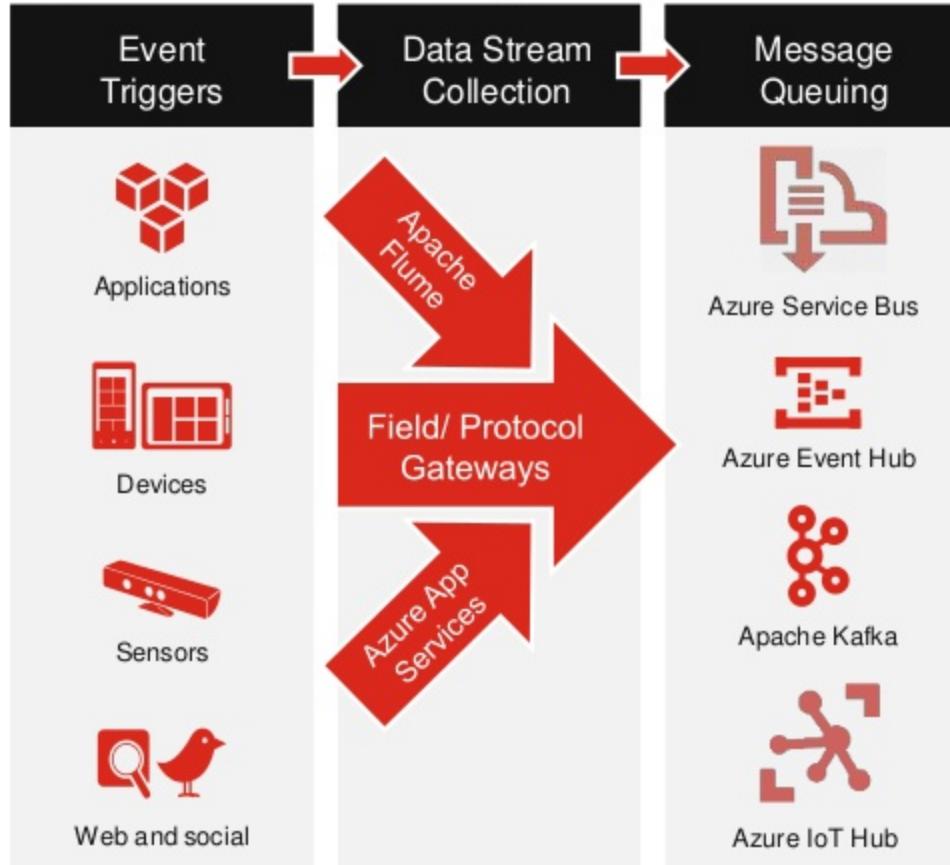
Sensors



Web and social

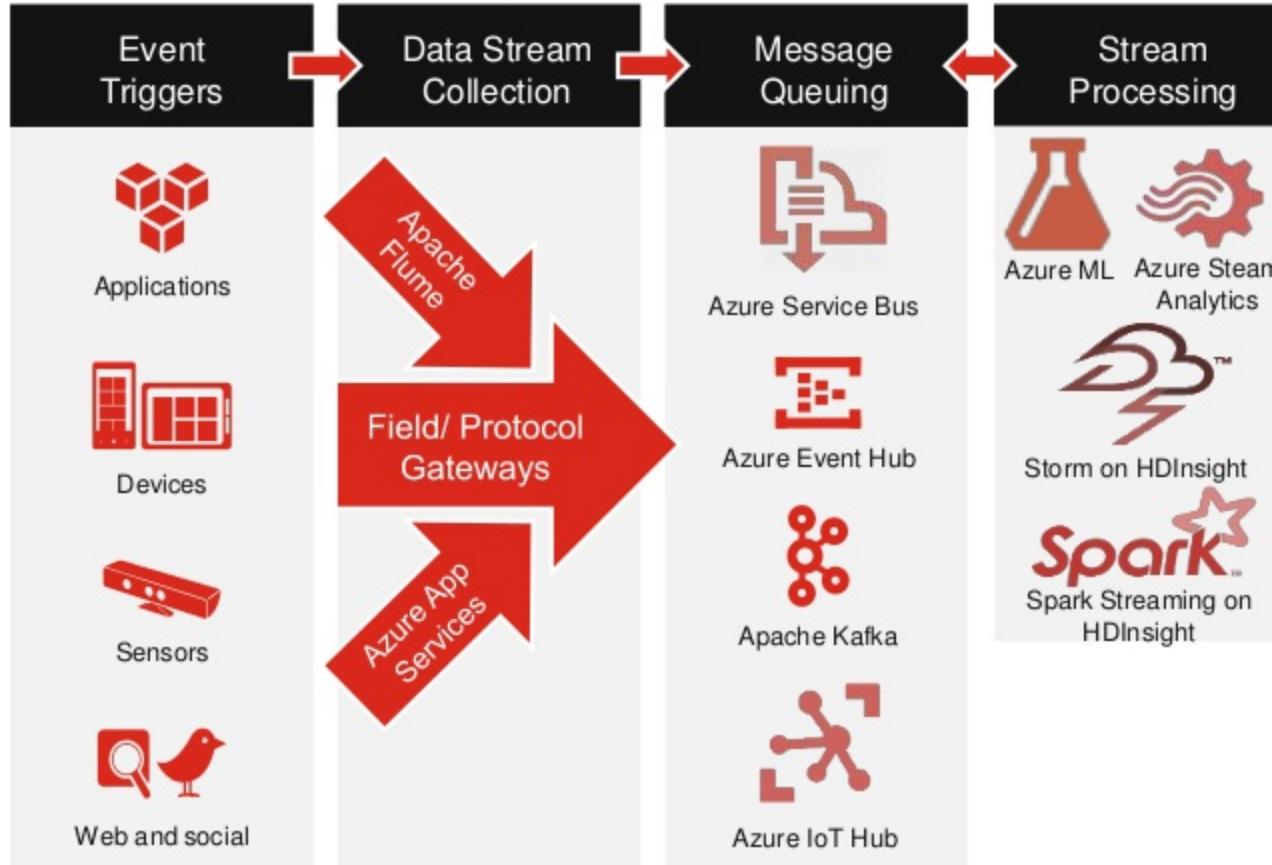
Events & Stream Processing Architecture

Azure Tools & Technologies



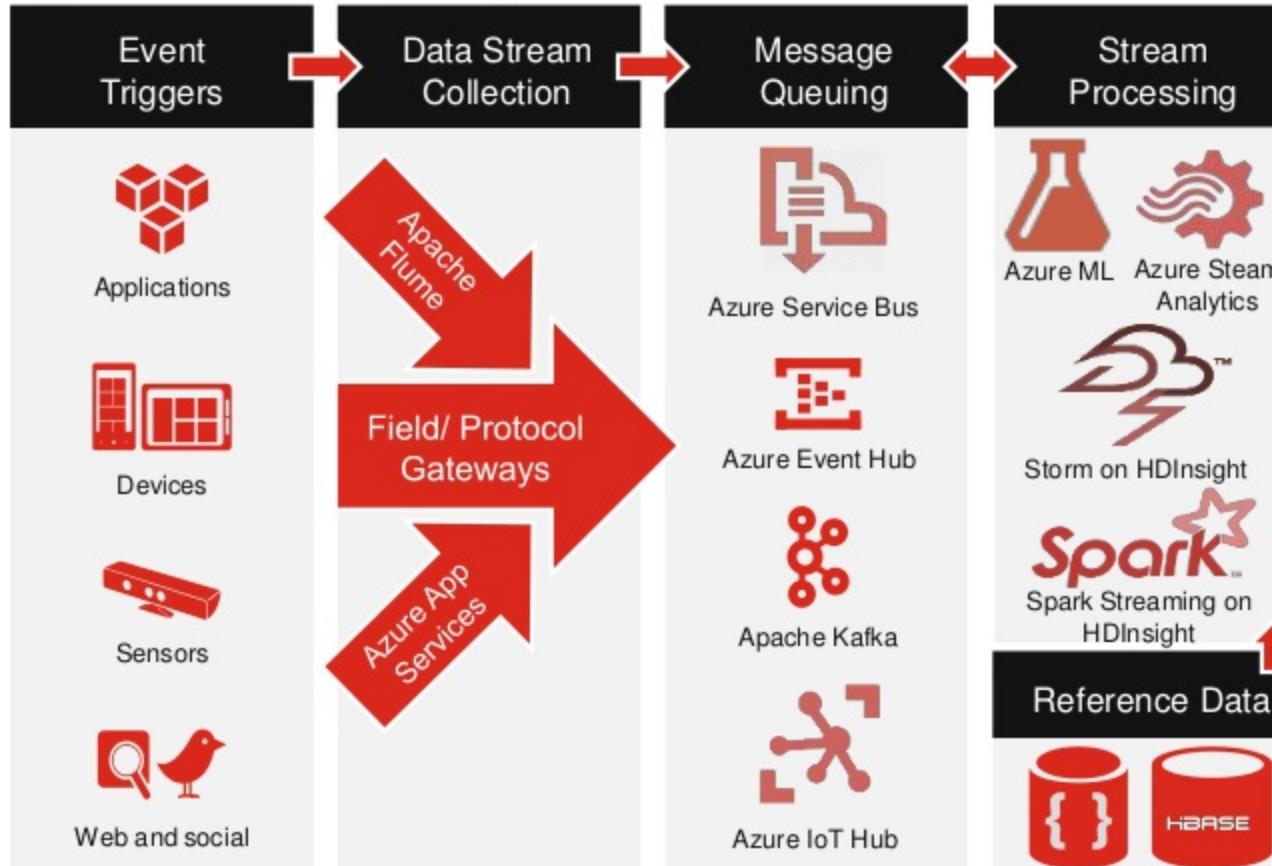
Events & Stream Processing Architecture

Azure Tools & Technologies



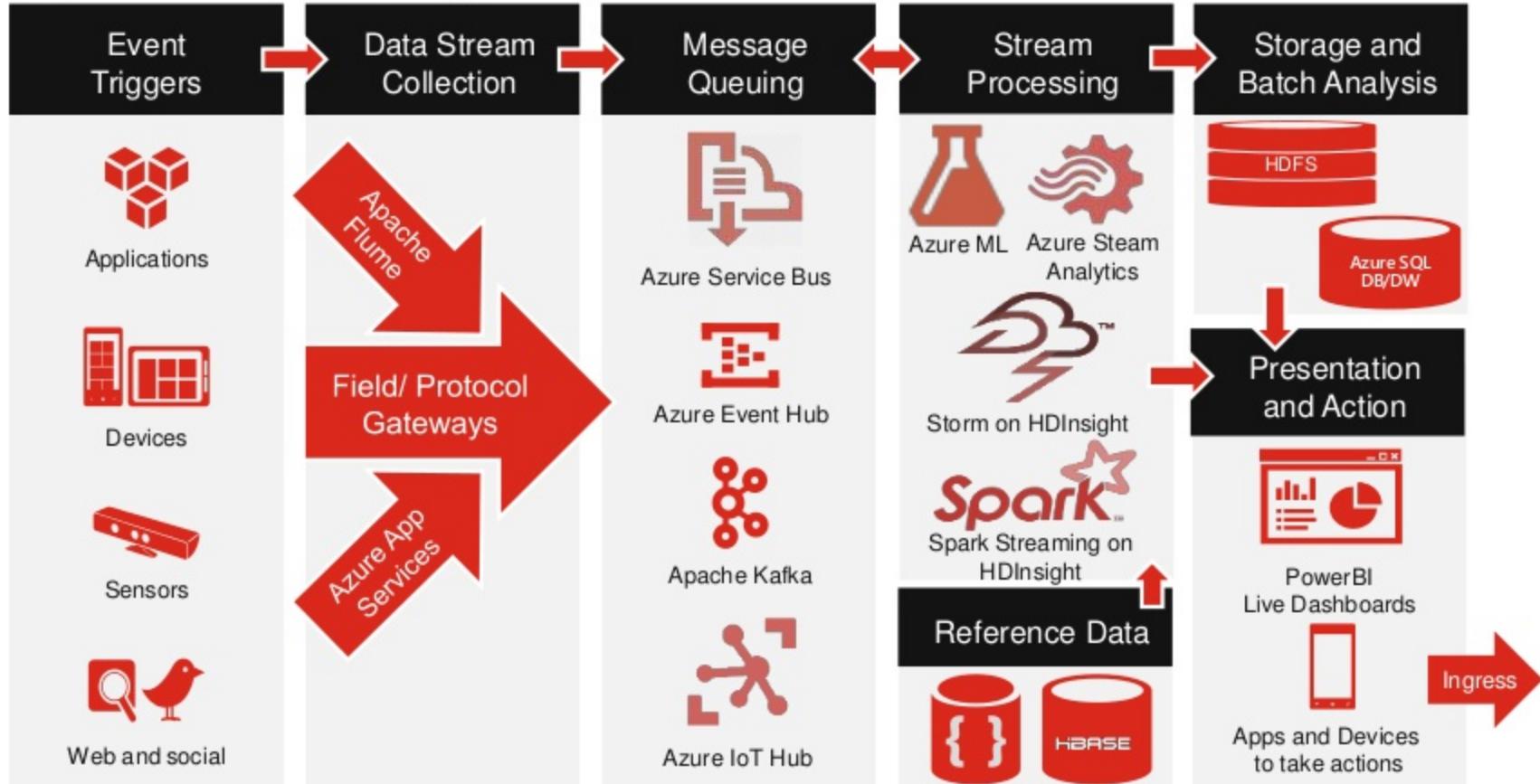
Events & Stream Processing Architecture

Azure Tools & Technologies



Events & Stream Processing Architecture

Azure Tools & Technologies



Events & Stream Processing Architecture

Data Sources

- Devices, Websites, and Apps that continuously produce data streams

Data Collection

- Listen to, collection, and transfer in-bound events

Message Queuing

- De-couples data consumers from data producers
- Reliable, distributed fault-tolerant, high-throughputs short-term storage

Stream Processing

- Aggregate / filter / join incoming event streams
- Temporal engine for analysing data across time-series windows

Reference Data

- High throughputs, random access data store to support processing
- Usually NoSQL data stores

Storage

- Store processed/ aggregated/ filtered data (SQL/NoSQL)
- Consolidate and store raw data into files for batch analysis (DFS)

Presentation

- Rich interactive visualizations for real-time data analysis
- Application integration for process automation

Azure Stream Analytics

A PaaS real-time **complex event processing** (CEP) on Microsoft Azure

Fully-managed real-time processing

- Intake millions of events per second
- Processing on continuous streams of data
- Reference data lookup
- Output to live dashboards and data stores

Mission Critical Reliability

- Guaranteed events delivery
- Preserves event order pre-device basis
- Guaranteed business continuity
- Auto-recovery from failures

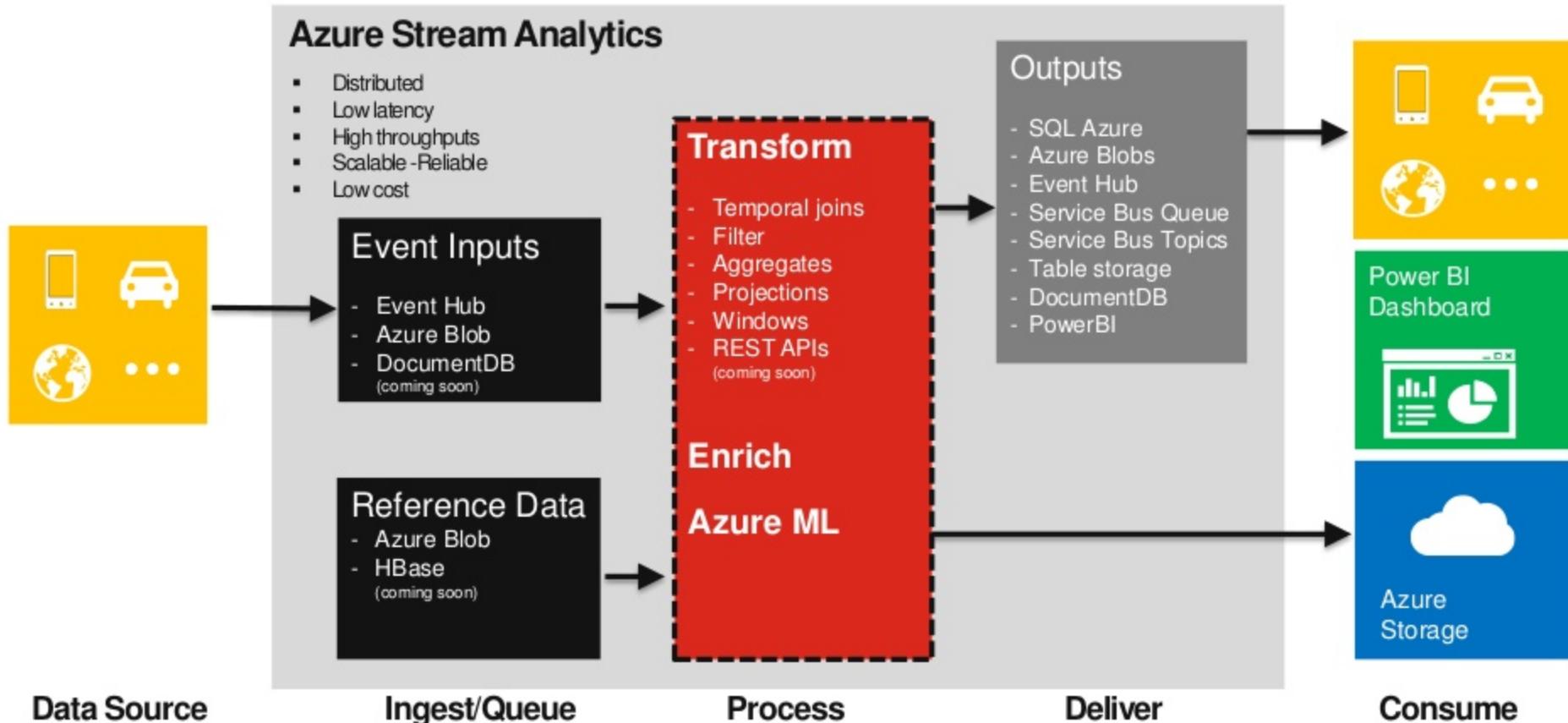
No challenges with Scale

- Elasticity for scale up or scale down
- Distributed, scale-out architecture
- Pay only for the resources you use

Rapid Development & Deployment

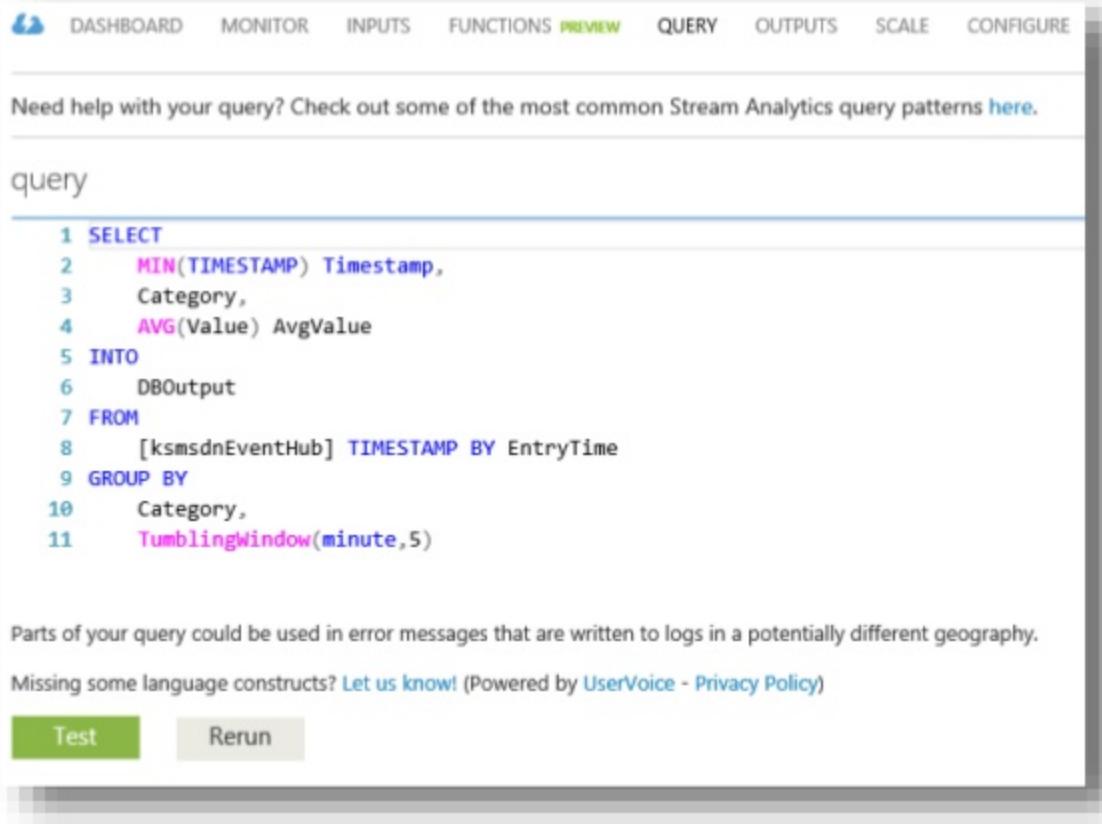
- SQL-like Language
- Built-in temporal semantics
- Up and running in a few clicks
- Scheduling and Monitoring

Azure Stream Analytics



Azure Stream Analytics

Stream Analytics Query



DASHBOARD MONITOR INPUTS FUNCTIONS **PREVIEW** QUERY OUTPUTS SCALE CONFIGURE

Need help with your query? Check out some of the most common Stream Analytics query patterns [here](#).

query

```
1 SELECT
2     MIN(TIMESTAMP) Timestamp,
3     Category,
4     AVG(Value) AvgValue
5 INTO
6     DBOutput
7 FROM
8     [ksmsdnEventHub] TIMESTAMP BY EntryTime
9 GROUP BY
10    Category,
11    TumblingWindow(minute,5)
```

Parts of your query could be used in error messages that are written to logs in a potentially different geography.

Missing some language constructs? [Let us know!](#) (Powered by [UserVoice](#) - [Privacy Policy](#))

Test **Rerun**

Stream Analytics Query Language

Windowing Functions

- In **data streams**, a common requirement is to perform aggregation (max, min, sum, count, etc.) over messages that **arrive within a specified period of time** (window) - **to detect events**.
- Each **Group By** requires a **windowing function**
- Each window operation **outputs a single event** at the end of the window
- All windows have a **fixed length**

Tumbling window

Aggregate per time interval

Hopping window

Schedule overlapping windows

Sliding window

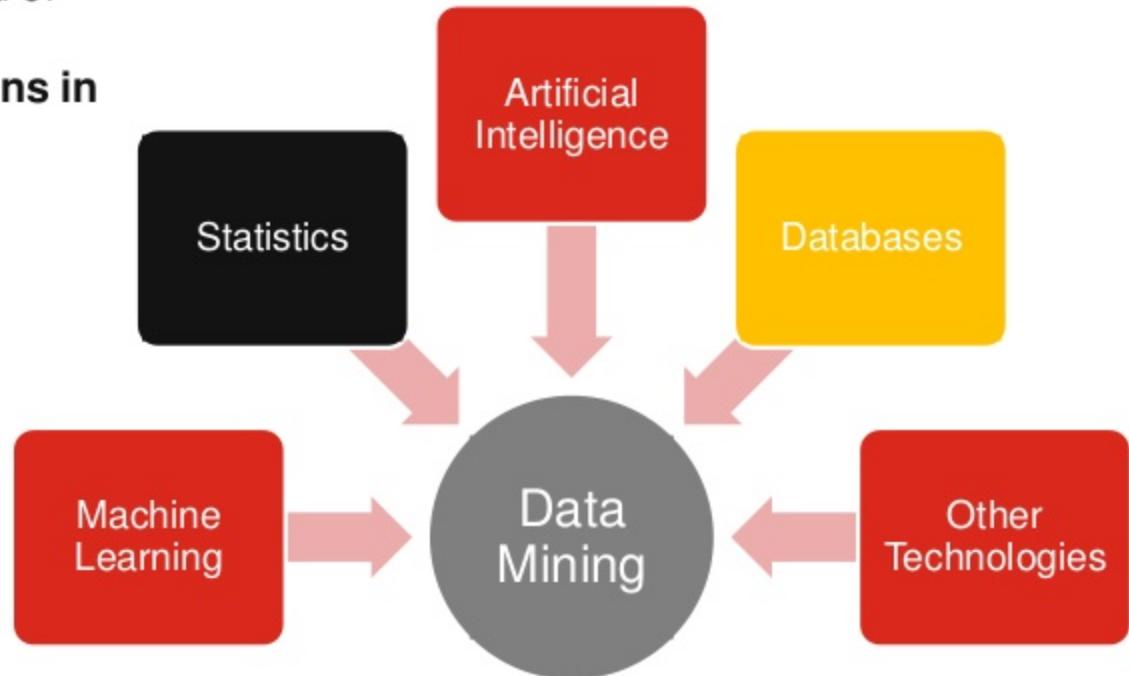
Windows constant re-evaluated

Data Science and Machine Learning on Microsoft Azure

Data Science and Machine Learning

What?

“**Data mining**, an interdisciplinary subfield of computer science, is the computational process of **automatic discovering patterns in large data sets**”



Other Related Technologies:

- Visualization
- Big Data
- High Performance Computing
- Cloud Computing
- Others..

Data Science and Machine Learning

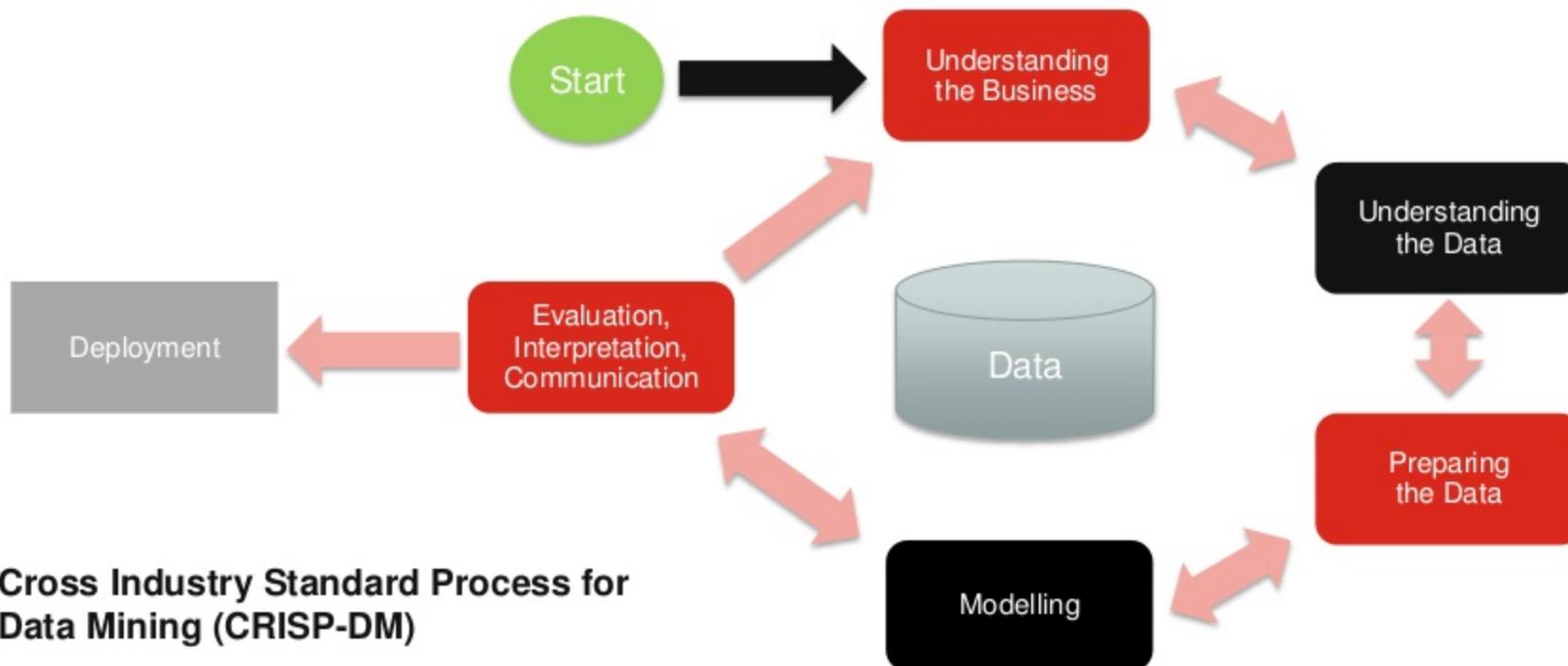
Why?

Machine learning & predictive analytics are core capabilities that are needed throughout your business



Data Science and Machine Learning

How?



**Cross Industry Standard Process for
Data Mining (CRISP-DM)**

Data Science and Machine Learning

How?

Classification Learning

Build a model that can predict the target class of an input case



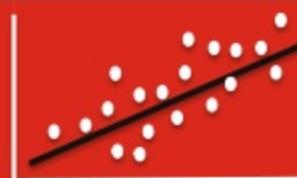
Time Series Analysis

Analysis of temporal data to forecast future values



Regression Modeling

Build a model that can estimate the response value given an input case



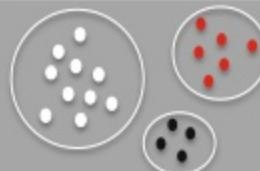
Probabilistic Modeling & Estimation

Compute the probability of an event to occur



Cluster Analysis

Discover natural groupings within the data points



Similarity Analysis

Identify similar cases to a given input case



Association Rule Discovery

Extract frequent patterns present in the data

```
IF ... AND ... AND ...
THEN A
ELSE IF ... AND ...
THEN C
ELSE IF ... AND ...
THEN B
...
```

Collaborative Filtering

Filtering of information using techniques involving collaboration viewpoints



Data Science and Machine Learning

Microsoft tools and technologies

**Data Mining - SQL Server
Analysis Services**



Azure Machine Learning



Spark ML – Azure HDInsight



Microsoft R Server



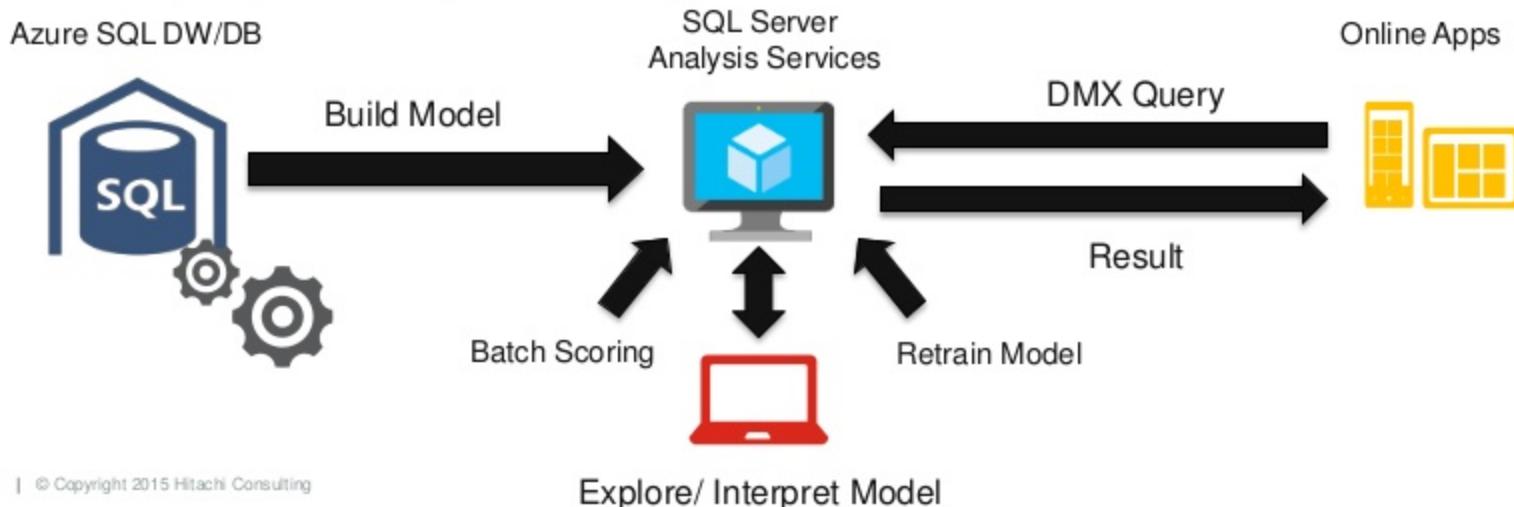
Azure Cognitive Services



Data Science and Machine Learning

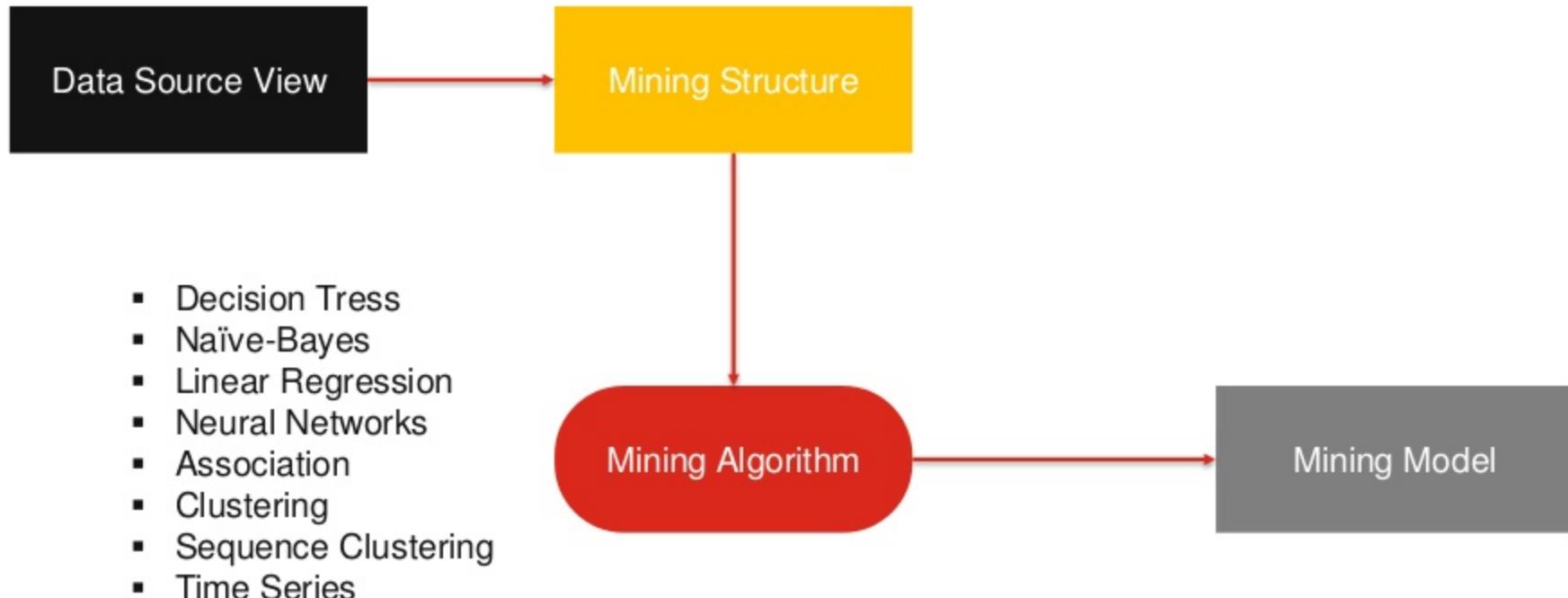
Microsoft Analysis Services

- SQL Server Analysis Services (SQL Server on a VM, no PaaS)
- Process data from many OLEDB and ODBC data sources
- Easy to **build, interpret, deploy, and productionize**
- Limited extensibility
- Custom APIs can be built to use constructed models (DMX)
- Limited pre-processing functionality



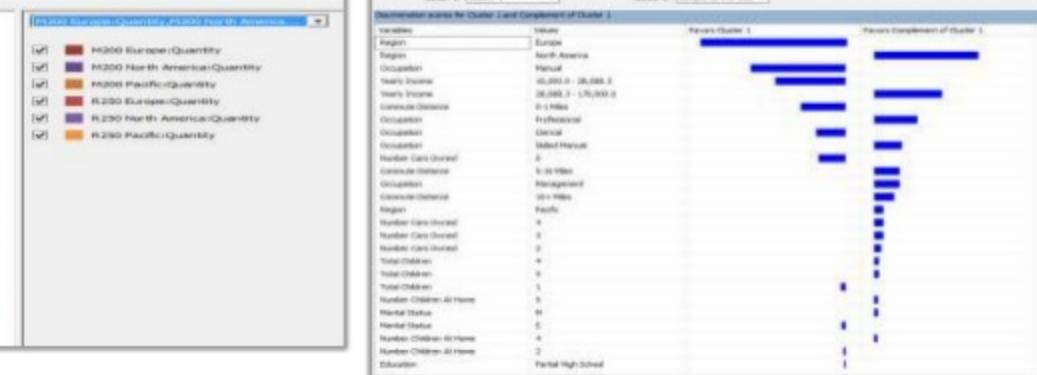
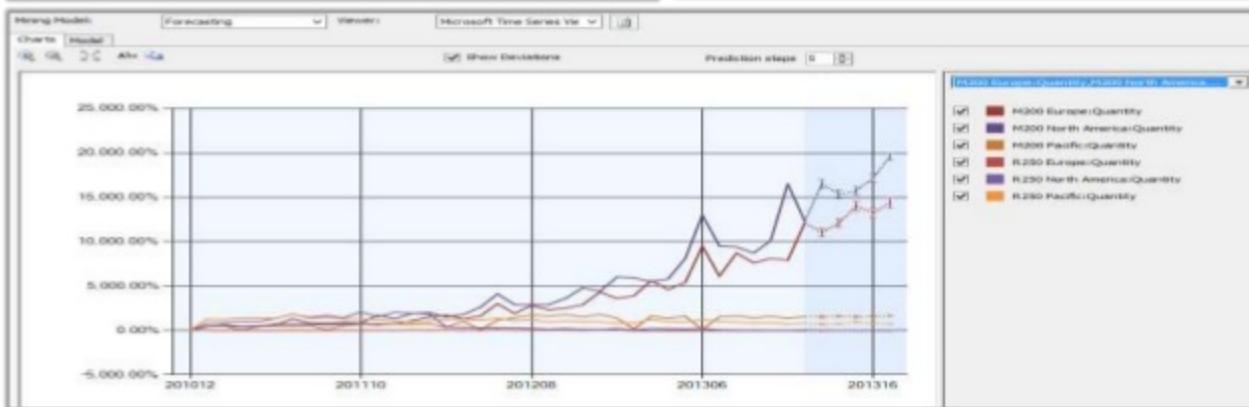
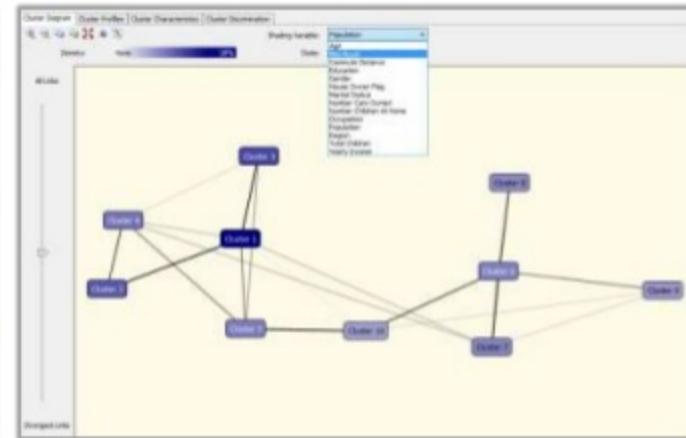
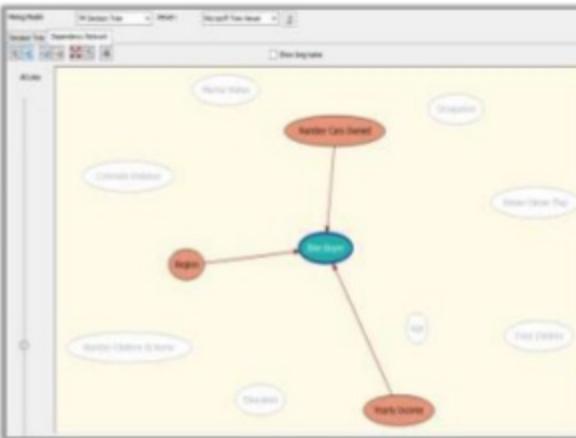
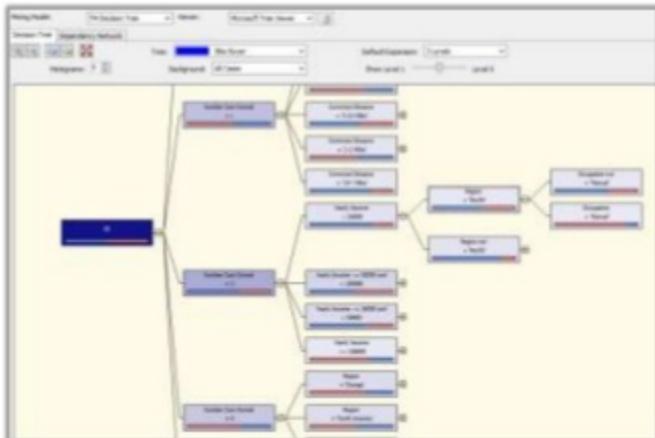
Data Science and Machine Learning

Microsoft Analysis Services



Data Science and Machine Learning

Microsoft Analysis Services



Data Science and Machine Learning

Spark ML on HDInsight

- HDInsight-based ML Services provided by Spark
- Spark MLlib and Spark ML functionality
- Stream mining functionality
- Imports data from HDFS (Blob Storage/Data Lake)
- Scalable – Highly Distributed
- Extensible – Python and R
- Suitable for Big Data - **Batch Model Training and Data Scoring**

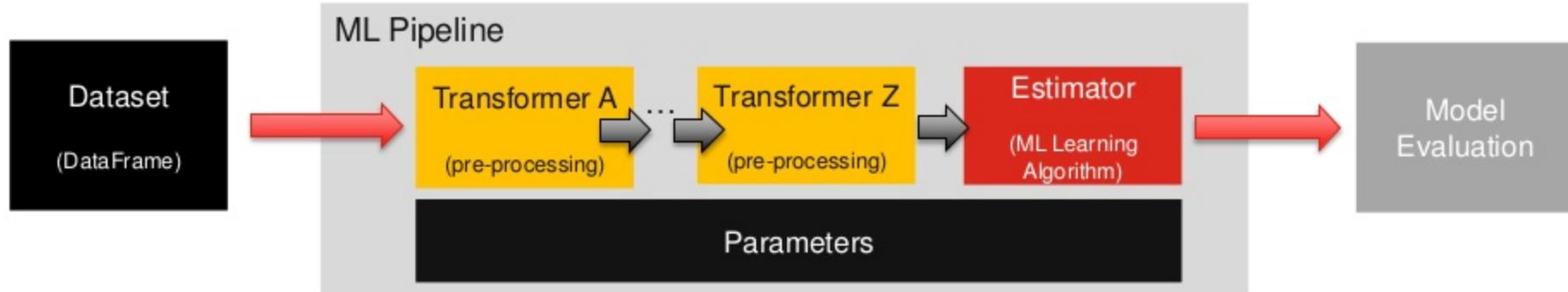


Data Science and Machine Learning

Spark ML on HDInsight – Pipelines

Spark ML standardizes APIs for machine learning algorithms to make it easier to combine multiple algorithms into a single pipeline, or workflow.

- **Transformers** – used for data pre-processing. Input: DataFrame - Output:DataFrame
- **Estimators** – ML algorithm used to build a predictive model. Input: DataFrame - Output: Model.
- **Parameters** – Configurations for Transformers and Estimators
- **Pipeline** – Chains Transformers and Estimators



Data Science and Machine Learning

Spark ML on HDInsight - functionality

Transformers

Text Feature Extraction

- TF-IDF (HashingTF and IDF)
- Word2Vec
- CountVectorizer
- Tokenizer
- StopWordsRemover
- n-gram

Feature Selection

- VectorSlicer
- RFormula
- ChiSqSelector

Dimensionality Reduction

- PCA

Features Vector Preparation

- VectorAssembler
- VectorIndexer
- StringIndexer
- IndexToString

Feature Type Conversion

- Binarizer
- Discrete Cosine Transform (DCT)
- OneHotEncoder
- Bucketizer
- QuantileDiscretizer

Feature Scaling

- Normalizer
- StandardScaler
- MinMaxScaler

Feature Construction

- SQLTransformer
- ElementwiseProduct
- PolynomialExpansion

Estimators (supervised)

Classification

- Decision Trees – Ensembles
- Naïve-Bayes
- SVM

Regression

- Linear Regression
- SVM

Other (Unsupervised)

Clustering

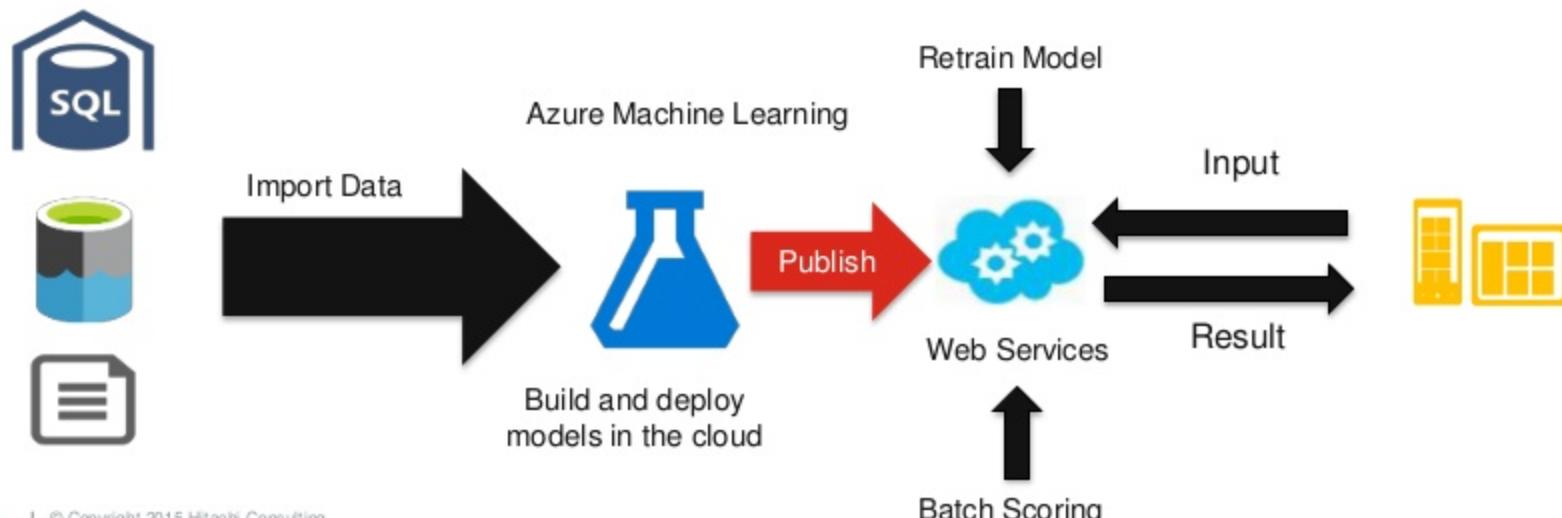
Collaborative Filtering

Frequent Pattern Mining

Data Science and Machine Learning

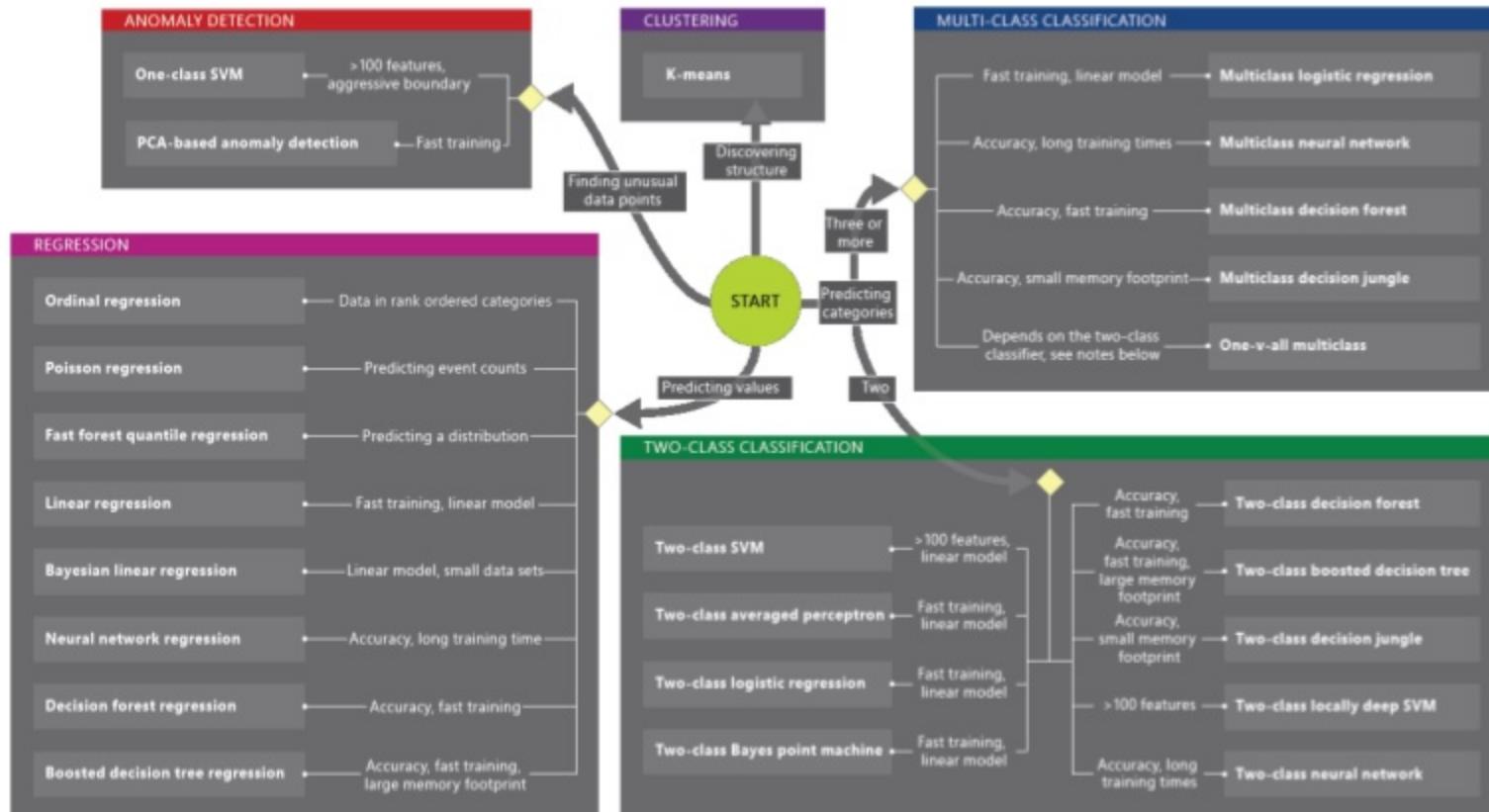
Azure Machine Learning

- Cloud-based Machine Learning Services
- Imports data from everywhere
- Easy to **develop and productionize** – Web Services
- Extensible via R and Python scripts
- Limited in exploration and model interpretation
- Rich pre-processing functionality



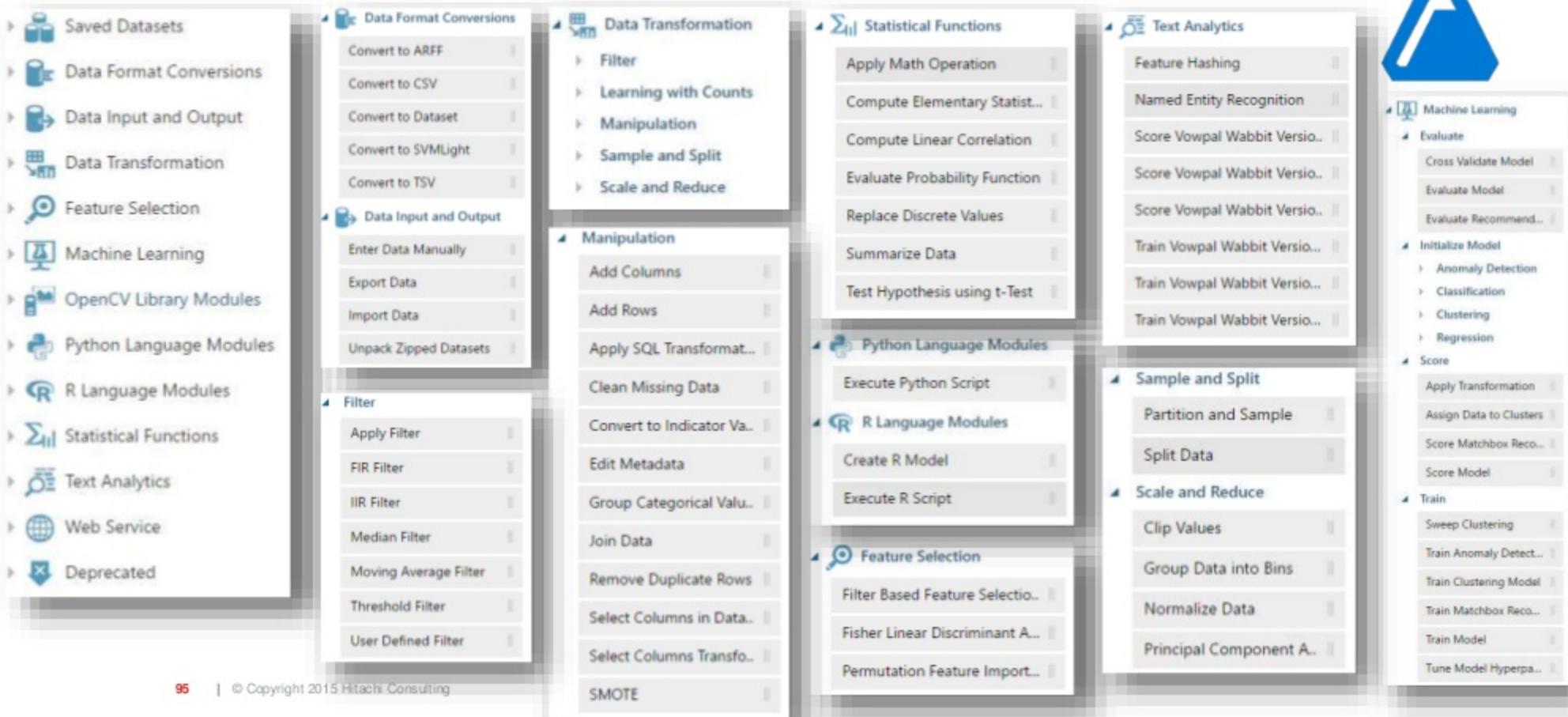
Data Science and Machine Learning

Azure Machine Learning – algorithms cheat sheet



Data Science and Machine Learning

Azure Machine Learning – ML Studio

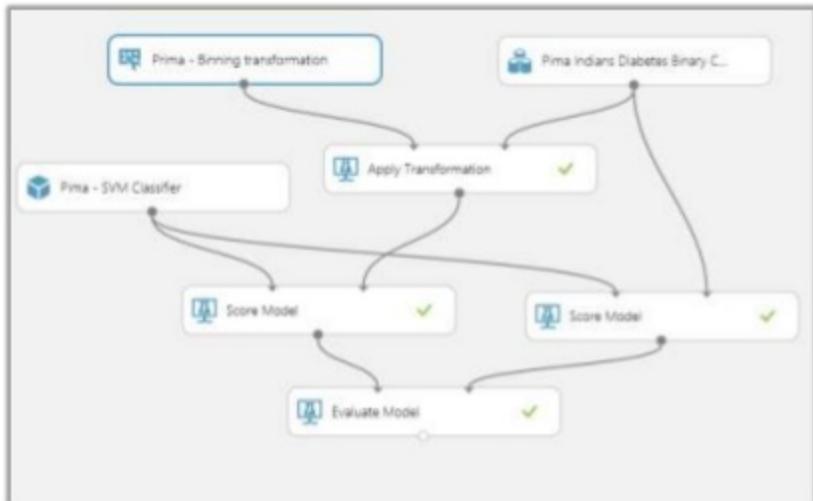
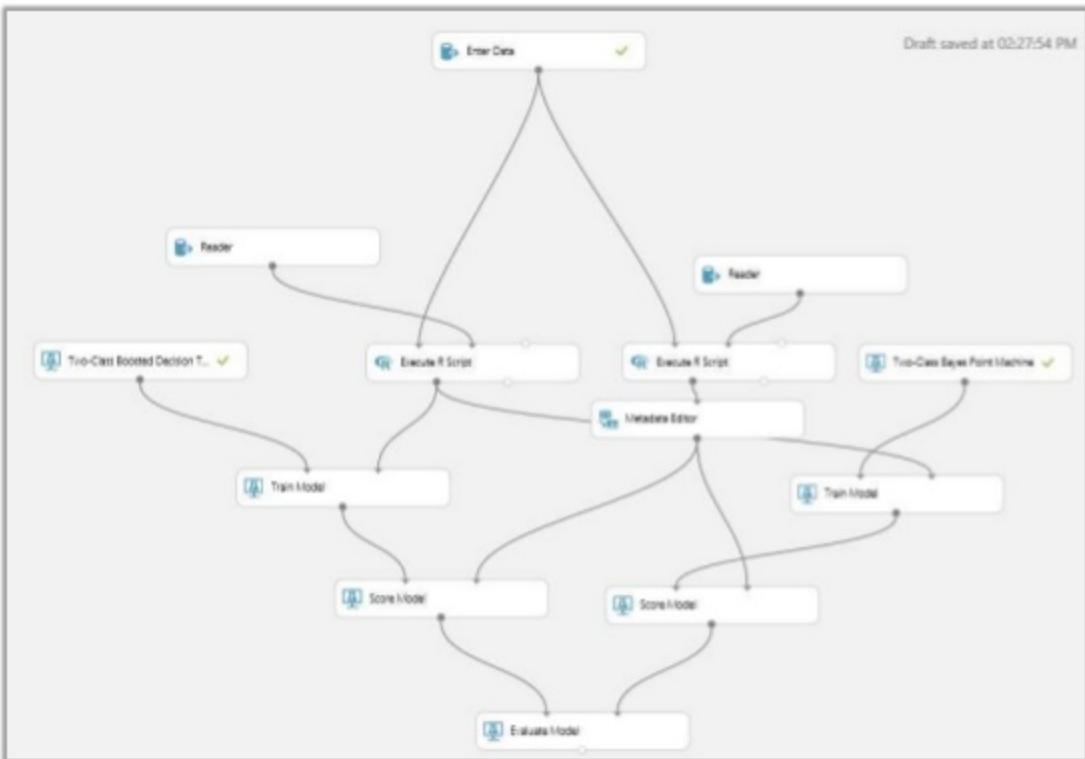


The screenshot displays the Azure Machine Learning - ML Studio interface, showing a hierarchical list of available modules:

- Saved Datasets**
- Data Format Conversions**
 - Convert to ARFF
 - Convert to CSV
 - Convert to Dataset
 - Convert to SVMLight
 - Convert to TSV
- Data Input and Output**
 - Enter Data Manually
 - Export Data
 - Import Data
 - Unpack Zipped Datasets
- Data Transformation**
 - Filter**
 - Apply Filter
 - FIR Filter
 - IIR Filter
 - Median Filter
 - Moving Average Filter
 - Threshold Filter
 - User Defined Filter
 - Learning with Counts**
 - Manipulation**
 - Add Columns
 - Add Rows
 - Apply SQL Transformation...
 - Sample and Split**
 - Scale and Reduce**
- Statistical Functions**
 - Apply Math Operation
 - Compute Elementary Statist...
 - Compute Linear Correlation
 - Evaluate Probability Function
 - Replace Discrete Values
 - Summarize Data
 - Test Hypothesis using t-Test
- Text Analytics**
 - Feature Hashing
 - Named Entity Recognition
 - Score Vowpal Wabbit Versio...
 - Score Vowpal Wabbit Versio...
 - Score Vowpal Wabbit Versio...
 - Train Vowpal Wabbit Versio...
 - Train Vowpal Wabbit Versio...
 - Train Vowpal Wabbit Versio...
- Machine Learning**
 - Evaluate**
 - Cross Validate Model
 - Evaluate Model
 - Evaluate Recommend...
 - Initialize Model**
 - Anomaly Detection
 - Classification
 - Clustering
 - Regression
 - Score**
 - Apply Transformation
 - Assign Data to Clusters
 - Score Matchbox Reco...
 - Score Model
 - Train**
 - Sweep Clustering
 - Train Anomaly Detect...
 - Train Clustering Model
 - Train Matchbox Reco...
 - Train Model
 - Tune Model Hyperpa...
- Python Language Modules**
 - Execute Python Script
- R Language Modules**
 - Create R Model
 - Execute R Script
- Feature Selection**
 - Filter Based Feature Selectio...
 - Fisher Linear Discriminant A...
 - Permutation Feature Import...
- Sample and Split**
 - Partition and Sample
 - Split Data
- Scale and Reduce**
 - Clip Values
 - Group Data into Bins
 - Normalize Data
 - Principal Component A...

Data Science and Machine Learning

Azure Machine Learning – ML Studio



Data Science and Machine Learning

Microsoft R Server

Microsoft R Open (MRO)

- Based on latest Open Source R (3.2.2.) - Built, tested, and distributed by Microsoft
- Enhanced by Intel Math Kernel Library (MKL) to speed up linear algebra functions
- More efficient and multi-threaded math computation
- Revolutions Open-Source R packages (**ParallelR**, **Rhadoop**, **AzureML**, etc.)
- Compatible with all R-related software



Data Science and Machine Learning

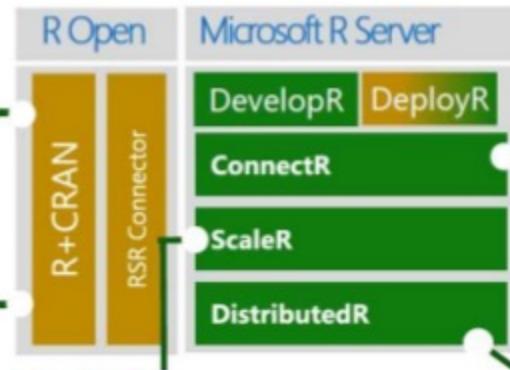
Microsoft R Server

R+CRAN

- Open source R interpreter
 - R 3.1.2
- Freely-available huge range of R algorithms
- Algorithms callable by RevoR
- Embeddable in R scripts
- 100% Compatible with existing R scripts, functions and packages

MRO

- Performance enhanced R interpreter
- Based on open source R
- Adds high-performance math library to speed up linear algebra functions



ConnectR

- High-speed & direct connectors

Available for:

- High-performance XDF
- SAS, SPSS, delimited & fixed format text data files
- Hadoop HDFS (text & XDF)
- Teradata Database & Aster
- EDWs and ADWs
- ODBC

ScaleR

- Ready-to-Use high-performance big data big analytics
- Fully-parallelized analytics
- Data prep & data distillation
- Descriptive statistics & statistical tests
- Range of predictive functions
- User tools for distributing customized R algorithms across nodes
- Wide data sets supported – thousands of variables

DistributedR

- Distributed computing framework
- Delivers cross-platform portability

Data Science and Machine Learning

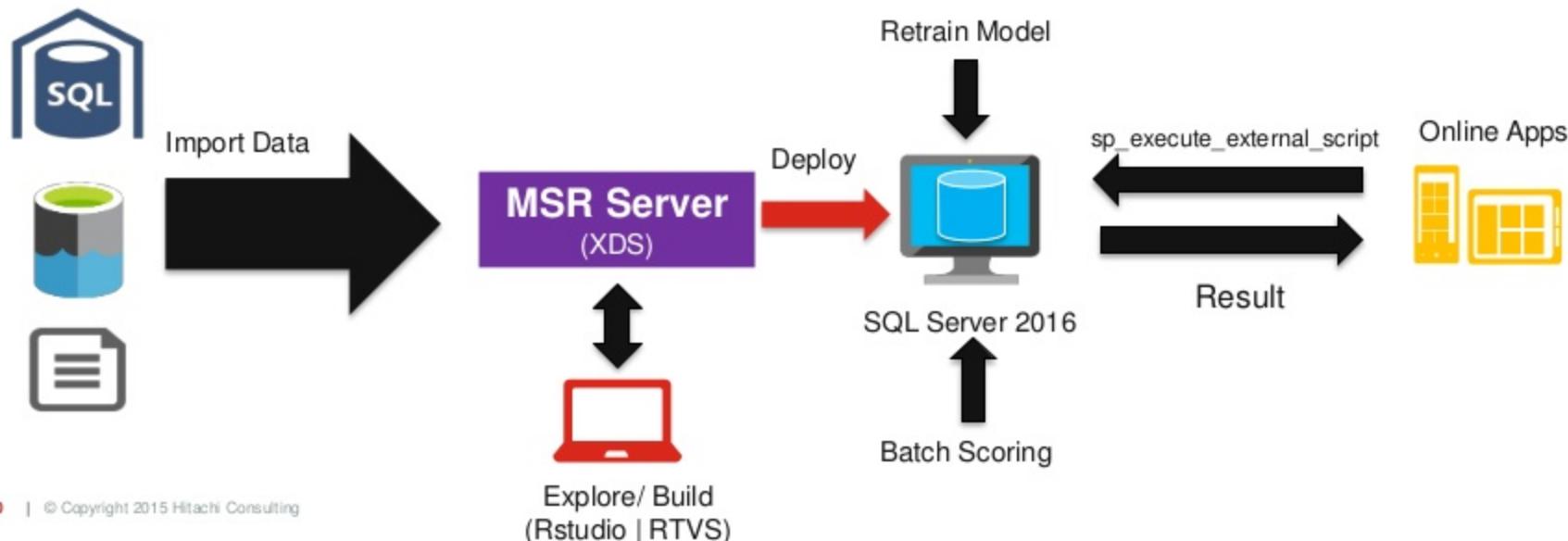
Microsoft R Server

	CRAN	MRO	MRS
Data size	In-memory	In-memory	In-memory & disk
Efficiency	Single threaded	Multi-threaded	Multi-threaded, parallel processing 1:N servers
Support	Community	Community	Community + Commercial
Functionality	7500+ innovative analytic packages	7500+ innovative analytic packages	7500+ innovative packages + commercial parallel high-speed functions
Licence	Open Source	Open Source	Commercial license.

Data Science and Machine Learning

Microsoft R Server

- Powerful and scalable
- Installed on **Windows, Linux, or as R Services** in SQL Server 2016
- Work in different compute contexts – **local, Hadoop, Spark, or SQL Server 2016**
- Very suitable for interactive Data Science activities
- Models are deployed to **SQL Server 2016, AzureML, or Web Service on R Server (DeployR)**



Data Science and Machine Learning

Microsoft R Server – ScaleR functionality

Data Preparation

- Data import – Delimited, Fixed, SAS, SPSS, OBDC
- Variable creation & transformation
- Recode variables
- Factor variables
- Missing value handling
- Sort, Merge, Split
- Aggregate by category (means, sums)

Descriptive Statistics

- Min / Max, Mean, Median (approx.)
- Quantiles (approx.)
- Standard Deviation
- Variance
- Correlation
- Covariance
- Sum of Squares (cross product matrix for set variables)
- Pairwise Cross tabs
- Risk Ratio & Odds Ratio
- Cross-Tabulation of Data (standard tables & long form)
- Marginal Summaries of Cross Tabulations

Statistical Tests

- Chi Square Test
- Kendall Rank Correlation
- Fisher's Exact Test
- Student's t-Test

Sampling

- Subsample (observations & variables)
- Random Sampling

Predictive Models

- Sum of Squares (cross product matrix for set variables)
- Multiple Linear Regression
- Generalized Linear Models (GLM) exponential family distributions: binomial, Gaussian, inverse Gaussian, Poisson, Tweedie. Standard link functions: cauchit, identity, log, logit, probit. User defined distributions & link functions.
- Covariance & Correlation Matrices
- Logistic Regression
- Classification & Regression Trees
- Predictions/scoring for models
- Residuals for all models

Variable Selection

- Stepwise Regression

Simulation

- Simulation (e.g. Monte Carlo)
- Parallel Random Number Generation

Cluster Analysis

- K-Means

Classification

- Decision Trees
- Decision Forests
- Gradient Boosted Decision Trees
- Naïve Bayes



Combination

- rxDataStep
- rxExec
- PEMA-R API Custom Algorithms

Data Science and Machine Learning

Azure Cognitive Services - SaaS

Language

Allow your apps to process natural language, evaluate sentiment and topics, and learn how to recognise what users want.



Language Understanding Intelligent Service

Teach your apps to understand commands from your users



Text Analytics API

Easily evaluate sentiment and topics to understand what users want



Web Language Model API

Use the power of predictive language models trained on web-scale data



Bing Spell Check API

Detecting and correcting spelling mistakes in your app

Speech

Processing spoken language in your applications



Bing Speech API

Convert speech to text and back again to understand user intent



Speaker Recognition API

Use speech to identify and authenticate individual speakers

Search

Make your apps, web pages and other experiences smarter and more engaging with the Bing Search APIs.



Bing Search APIs

Search, image, video and news APIs for your apps



Bing Autosuggest API

Give your app intelligent autosuggest options for searches



Vision

State-of-the-art image processing algorithms to build more personalised apps by returning smart insights such as faces, images and emotion recognition.



Face API

Detect, analyse, organise and tag faces in photos



Emotion API

Personalise user experiences with emotion recognition

Knowledge

Map complex information and data in order to solve tasks such as intelligent recommendations and semantic search.



Recommendations API

Predict and recommend items that your customers want

Azure Data Catalog

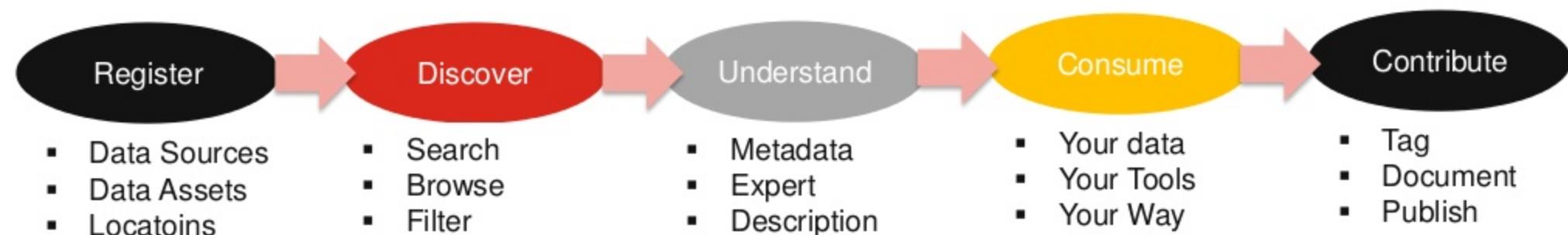
Azure Data Catalog

Enterprise-wide metadata catalogue for data assets discovery

A cloud-based service into which data source can be registered, described and discovered by enterprise-wide users



- Registration of central data sources
- Self-service business intelligence
- Capturing tribal knowledge



Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Data Asset

container that manages data assets.
(Database table/view, Blob Storage file, etc.)

Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Data Asset

container that manages data assets.
(Database table/view, Blob Storage file, etc.)

Data Asset Location

location of a data source or data asset, which can be used to connect to the source using a client application

Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Descriptive Metadata

describe the purpose and the relevance of the data asset (free-form description, tags, documentations, etc.)

Structural Metadata

metadata extracted from a data source that describes the structure of a data asset (names, data types, etc.)

Data Asset

container that manages data assets. (Database table/view, Blob Storage file, etc.)

Data Asset Location

location of a data source or data asset, which can be used to connect to the source using a client application

Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Descriptive Metadata

describe the purpose and the relevance of the data asset (free-form description, tags, documentations, etc.)

Structural Metadata

metadata extracted from a data source that describes the structure of a data asset (names, data types, etc.)

Data Asset

container that manages data assets.
(Database table/view, Blob Storage file, etc.)

Data Asset Location

location of a data source or data asset, which can be used to connect to the source using a client application

Owner

has additional privileges for managing a data asset

Expert

identified as having an informed perspective for a data asset

Azure Data Catalog

Terminology

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Data Asset

container that manages data assets.
(Database table/view, Blob Storage file, etc.)

Data Asset Location

location of a data source or data asset, which can be used to connect to the source using a client application

Descriptive Metadata

describe the purpose and the relevance of the data asset (free-form description, tags, documentations, etc.)

Structural Metadata

metadata extracted from a data source that describes the structure of a data asset (names, data types, etc.)

Preview

a snapshot of up to 20 records that can be extracted from the data source during registration, and stored in the catalog

Owner

has additional privileges for managing a data asset

Expert

identified as having an informed perspective for a data asset

Profile

a snapshot of table-level and column-level metadata about the content of a registered data asset

Azure Data Catalog

Terminology

Glossary

Terms, parent terms, and definitions

Catalog

metadata repository in which data sources and data assets can be registered.

Data Source

container that manages data assets. (Database Server, Blob Storage Account, etc.)

Data Asset

container that manages data assets.
(Database table/view, Blob Storage file, etc.)

Data Asset Location

location of a data source or data asset, which can be used to connect to the source using a client application

Descriptive Metadata

describe the purpose and the relevance of the data asset (free-form description, tags, documentations, etc.)

Profile

a snapshot of table-level and column-level metadata about the content of a registered data asset

Structural Metadata

metadata extracted from a data source that describes the structure of a data asset (names, data types, etc.)

Owner

has additional privileges for managing a data asset

Preview

a snapshot of up to 20 records that can be extracted from the data source during registration, and stored in the catalog

Expert

identified as having an informed perspective for a data asset

Azure Data Catalog

Preview

Microsoft Azure Data Catalog

2005

Filter Current Filters: Search Term: sales Select All

Tags: GMD (12) Sort Drinks (11) Carbonated Drinks (2) See more

Object Type: View (164) Table (54) Measure (3) KPI (1) See more

Source Type: SQL Server (147) SQL Server Analysis Services Multidimensional (2) SQL Server Analysis Services Tabular (2) SQL Server Reporting Services (1) Experts: HQ Sales Team (12) See more

253 search results, 1 selected | Select All

Soft Drink Sales Amounts... Channel Revenue vw_Fact_MSS_Sales Internet Average Sales A...

The tabular model representing sales amounts of all our core Soft Drink... test Adventure works... click Use to add a description... Experts: HQ Sales Team (4) GMD [] Sales Amounts: 2015 YTD

Sort Drinks Tabular Soft Drinks Sales Amounts: 2015 YTD

Contained in Model: Softdrinks ANALYSIS SERVICES MEASURE Open In... []

test Adventure works... Experts: HQ Sales Team (4) Revenue: Soft Drinks

Contained in Model: Adventure WorksDW2012M... ANALYSIS SERVICES KPI Open In... []

vw_Fact_MSS_Sales click Use to add a description... Experts: GMD []

Contained in Database: [] SQL SERVER VIEW Open In... []

Internet Average Sales A... click Use to add a description... Experts: []

Contained in Model: Adventure WorksDW2012M... ANALYSIS SERVICES MEASURE Open In... []

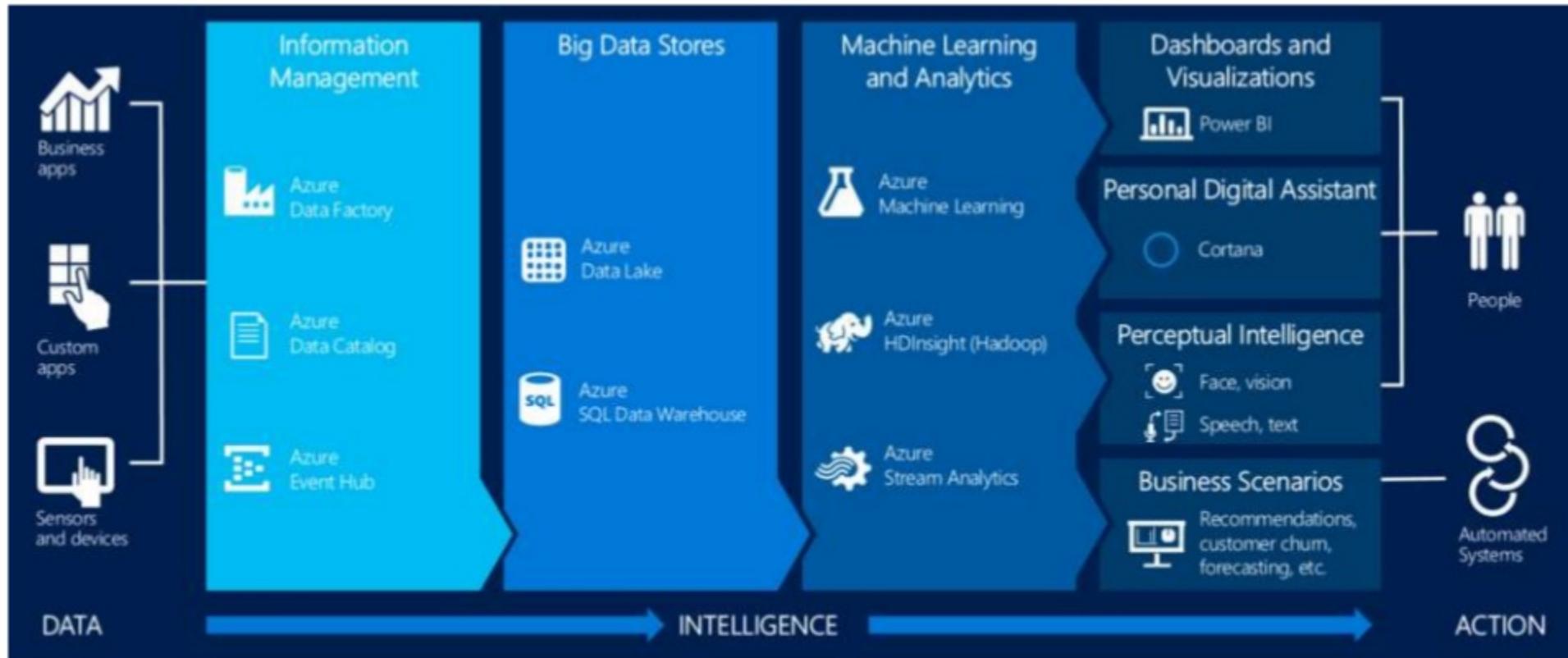
Properties

Name: vw_Fact_MSS_Sales Friendly Name: add friendly name... Description: add your description... Experts: Add... Tag: GMD Add... Connection Info: Server Name: Database Name: Schema Name: sales Object Name: sales

Preview Columns

Column Name	Data Type	Tags	Description
CalculatedName	varchar	Add...	
SubsidiaryKey	tinyint	Add...	add your description...
SubsidiaryCode	varchar	Add...	add your description...
BusinessKey	int	Add...	add your description...

Cortana Analytical Suite



Acknowledgement

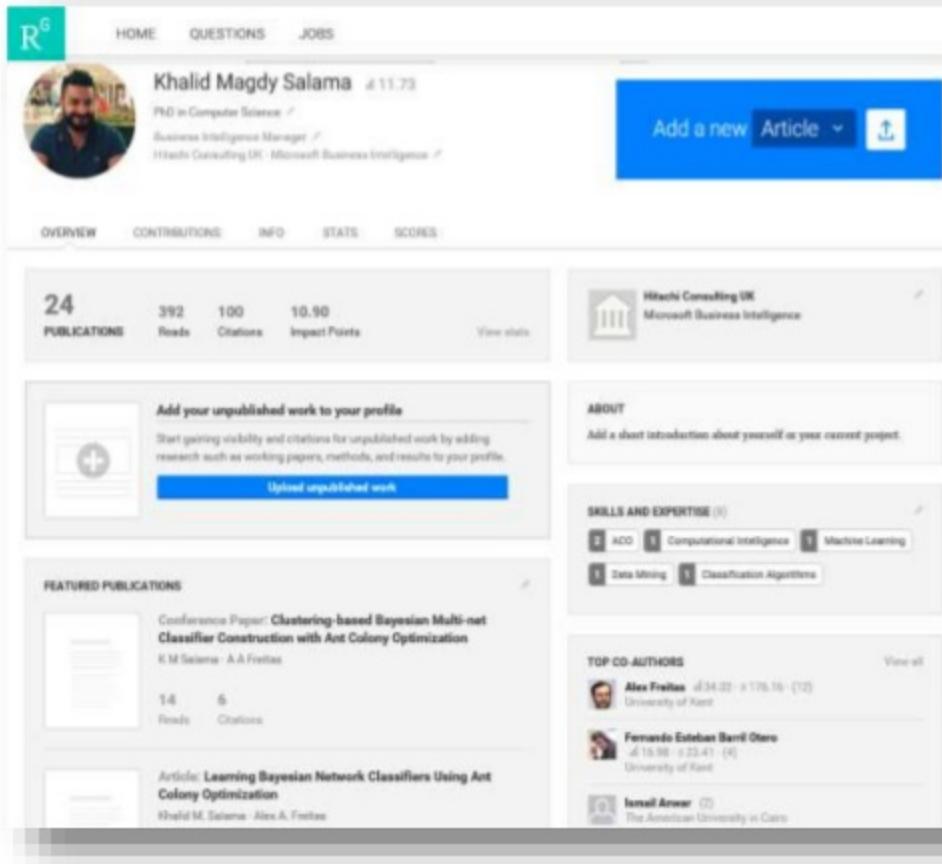
This guy is awesome

Thanks to Paul Lineham, the best big data architect I worked with, for supplying me with practical ideas and knowledge on architecting and delivering data platforms.

My Background

Applying Computational Intelligence in Data Mining

- Honorary Research Fellow, School of Computing , University of Kent.
- Ph.D. Computer Science, University of Kent, Canterbury, UK.
- M.Sc. Computer Science , The American University in Cairo, Egypt.
- 25+ published journal and conference papers, focusing on:
 - *classification rules induction,*
 - *decision trees construction,*
 - *Bayesian classification modelling,*
 - *data reduction,*
 - *instance-based learning,*
 - *evolving neural networks,* and
 - *data clustering*
- **Journals:** *Swarm Intelligence, Swarm & Evolutionary Computation, Applied Soft Computing, and Memetic Computing.*
- **Conferences:** ANTS, IEEE CEC, IEEE SIS, EvoBio, ECTA, IEEE WCCI and INNS-BigData.



The screenshot shows Khalid Magdy Salama's profile on ResearchGate. At the top, there is a green header with the RG logo, followed by 'HOME', 'QUESTIONS', and 'JOBS'. Below the header is a circular profile picture of Khalid Magdy Salama, a man with a beard, wearing a blue shirt. To his right, the text reads 'Khalid Magdy Salama' with a small 'd. 11.73' icon, 'PhD in Computer Science', 'Business Intelligence Manager', and 'Hitachi Consulting UK - Microsoft Business Intelligence'. On the far right of the header is a blue button with white text: 'Add a new Article' and a small upload icon.

Below the header, there are tabs for 'OVERVIEW', 'CONTRIBUTIONS', 'INFO', 'STATS', and 'SCORES'. Under 'OVERVIEW', it shows '24 PUBLICATIONS', '392 Reads', '100 Citations', and '10.90 Impact Points'. There is also a 'View stats' link. A large central box contains a placeholder for 'Add your unpublished work to your profile' with a 'Upload unpublished work' button. To the right of this box are sections for 'ABOUT' (with a placeholder 'Add a short introduction about yourself or your current project.'), 'SKILLS AND EXPERTISE' (listing 'ACO', 'Computational Intelligence', 'Machine Learning', 'Data Mining', and 'Classification Algorithms'), and 'TOP CO-AUTHORS' (listing Alex Freitas, Fernando Esteban Barril Otero, and Ismail Awad).

Under 'FEATURED PUBLICATIONS', two publications are listed:

- Conference Paper: 'Clustering-based Bayesian Multi-net Classifier Construction with Ant Colony Optimization' by K M Salama & A A Freitas. It has 14 Reads and 6 Citations.
- Article: 'Learning Bayesian Network Classifiers Using Ant Colony Optimization' by Khalid M. Salama & Alex A. Freitas. It has 14 Reads and 6 Citations.

ResearchGate.org



Thank you!