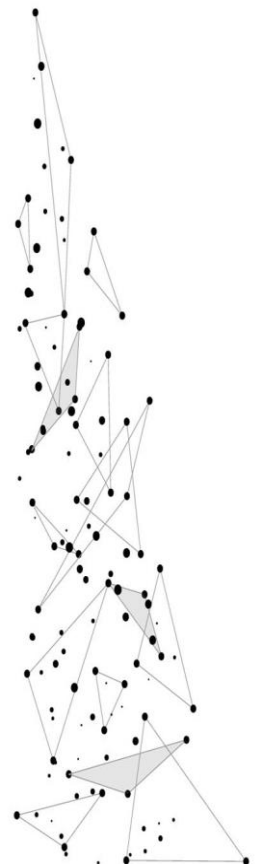


# Analiza **Big Data** z **Azure Data Lake**

Tomasz Krawczyk

# Agenda

- Big Data
- Azure Data Lake
- Azure **Data Lake Store** and **Analytics**
  - **U-SQL (Extensions)**
    - Azure Data Lake Analytics Runtime
    - Costs and Optimization
      - Azure Data Lake and Azure Data Factory
- Demo and Q&A



# Big Data - introduction



# 163 Zetta bytes by 2025

163 Zetta bytes =

163 000 000 000 000 000 000 000 000 bytes

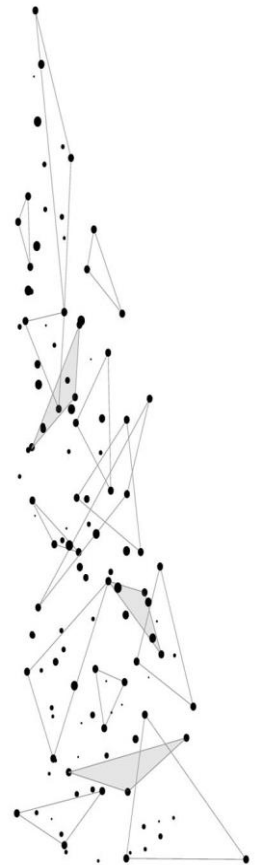
163 \* 2<sup>70</sup> bytes      163 \* 1024<sup>7</sup> bytes

If 1 bytes = One grain of rice

then **Zetta byte** = Fills the Pacific Ocean

Watching the Netflix catalog

**489 million times**



# 3Vs of Big Data

## • Data Volume

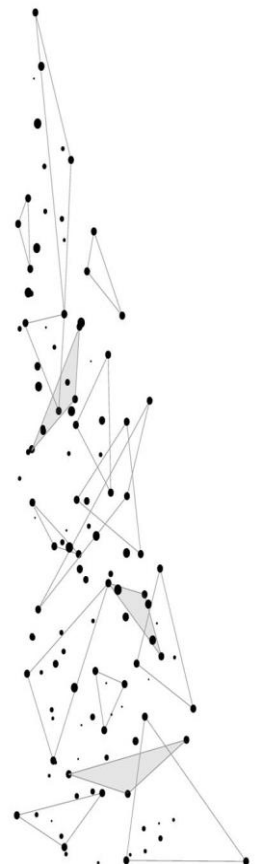
- **Byte** One grain of rice
- **Kilobyte** Cup of rice
- **Megabyte** 8 bags of rice
- **Gigabyte** 3 semi trucks
- **Terabyte** 2 container ships
- **Petabyte** Blankets Manhattan
- **Exabyte** Blankets west coast states
- **Zettabyte** Fills the Pacific Ocean
- **Yottabyte** As earth-sized rice ball

## • Data Variety

- Structured
- Unstructured
- Semi-structured
- All the above

## • Data Velocity

- Near to Real Time
- Batch



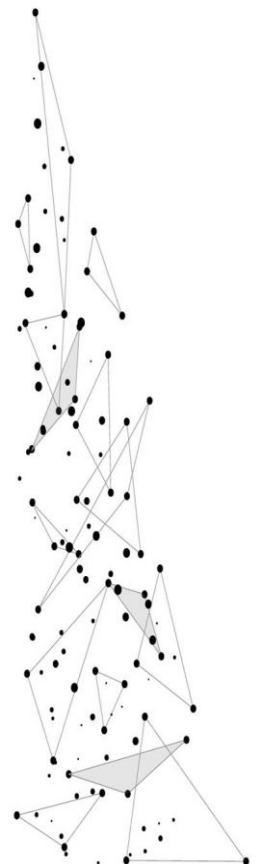
# Schema-on-Read vs Schema-on-Write

## SCHEMA-ON-READ (HADOOP OR ADLS):

- Copy data in its native format
- Create schema + parser
- Query Data in its native format  
(does ETL on the fly)

## SCHEMA-ON-WRITE (RDBMS):

- Create static DB schema
- Transform data into RDBMS
- Query data in RDBMS  
format

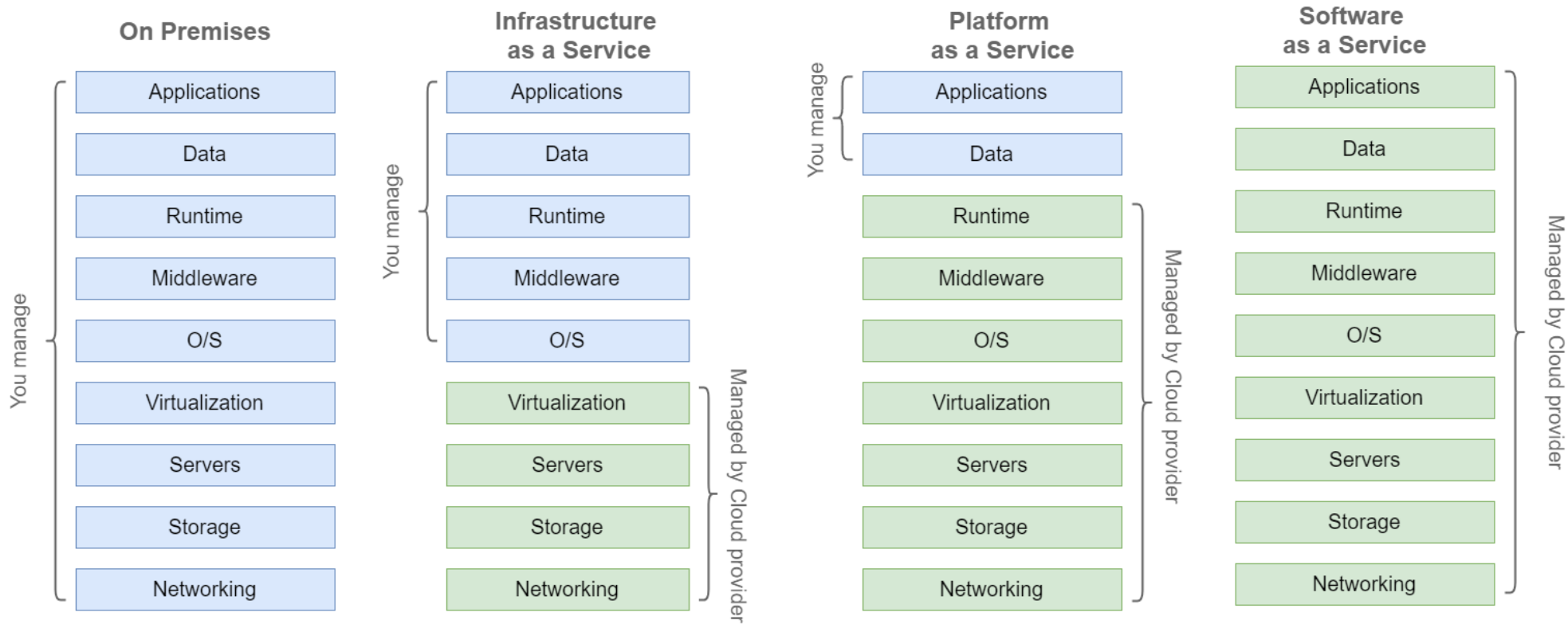




**How do you process all the data ?  
Why Cloud?**

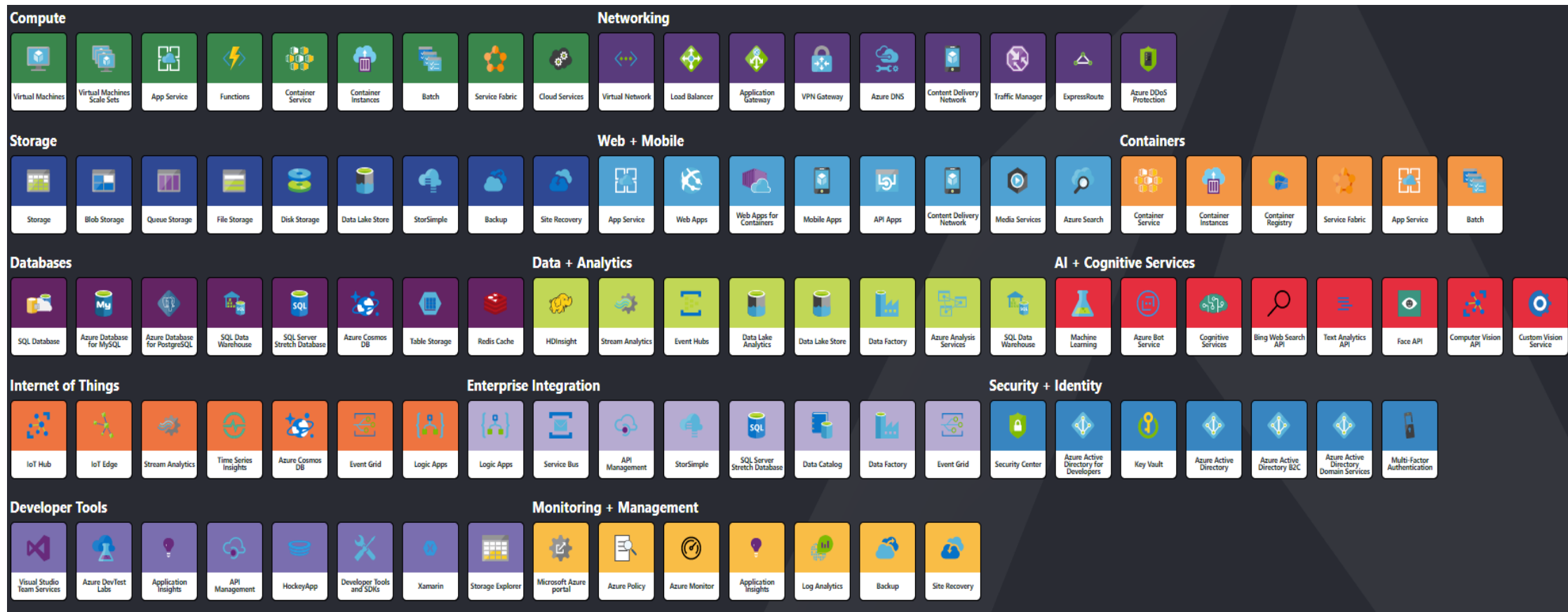


# Cloud Models





# Azure – Services



# Data Lake Approach

## What is Data Lake ?

“If you think of a **datamart** (a subset of a data warehouse) as a store of bottled water – cleansed and packaged and structured for easy consumption – the **data lake** is a large body of water in a more **natural state**,”

Pentaho CTO James Dixon



Source: <https://premiumwaters.com>



Source :<https://snowbrains.com>

I(ngest) S(tore) A(nalyse) S(urface) A(ct)

**Make Me More Money**

# Data Lake ?

**Ingest all data**  
regardless of  
requirements

**Store all data**  
in native format  
without schema  
definition

**Do analysis**  
Using analytic  
engines like Hadoop



# Big Data on Azure

Analytics

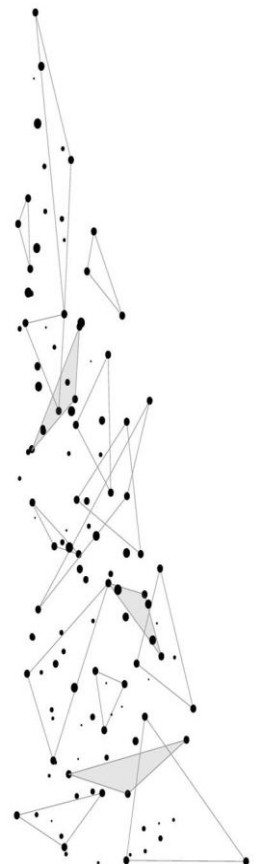


Storage



# Azure Data Lake Store (Gen1 & Gen2)

- Built for Hadoop
- **WebHDFS**-compatible **REST** interface
- Unlimited storage, petabyte files
- Highly-available and secure
- Supports files and folders objects
- Files are split apart into **Extents** (250 MB)
- For availability and reliability, extents are replicated (3 copies).
- Enables: **Parallel read** and **Parallel write**



# Azure Data Lake Store vs Azure Blob Storage

	Azure Data Lake Store	Azure Blob Storage
Purpose	<b>Optimized storage for big data analytics workloads</b>	General purpose object store for a wide variety of storage scenarios
Use Cases	Batch, interactive, streaming analytics and machine learning data such as log files, IoT data, click streams, large datasets	Any type of text or binary data, such as application back end, backup data, media storage for streaming and general purpose data
Key Concepts	Data Lake Store account contains folders, which in turn contains data stored as files	Storage account has containers, which in turn has data in the form of blobs
Structure	<b>Hierarchical file system</b>	Object store with flat namespace
API	REST API over HTTPS	REST API over HTTP/HTTPS
Hadoop File System Client	Yes	Yes
Data Operations - Authentication	Based on Azure Active Directory Identities	Based on shared secrets - Account Access Keys and Shared Access Signature Keys.
Data Operations - Authorization	<b>POSIX Access Control Lists (ACLs). ACLs based on Azure Active Directory Identities can be set file and folder level.</b>	For account-level authorization – Use Account Access Keys For account, container, or blob authorization - Use Shared Access Signature Keys



# Azure – Big Data Storage



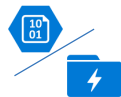
## Azure Blob Storage

- General purpose object store
- Object store with flat namespace
- Hot/cold/archive tiers
- Data replication and redundancy options



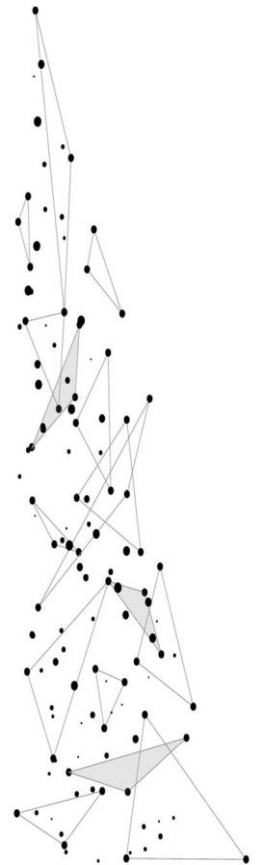
## Azure Data Lake Storage (Gen1)

- Unlimited storage, petabyte files
- **WebHDFS**-compatible REST interface
- Hadoop and big data optimizations
- Supports files and folders objects



## Azure Data Lake Storage (Gen2)

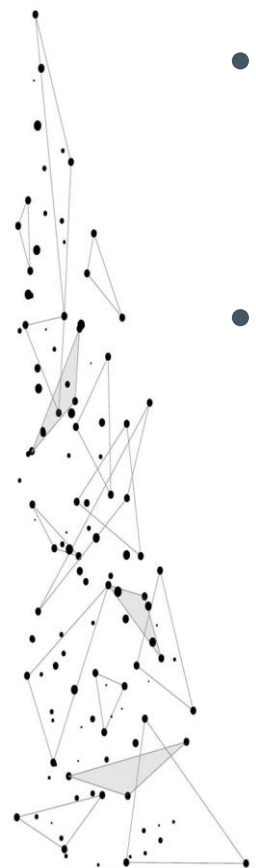
- Multi-modal combining features from both of the above
- Not a separate service: Azure Storage with new features



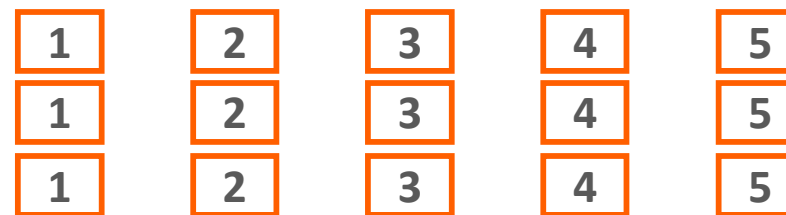


# Azure Data Lake Storage

- Files are split apart into **Extents** (250 MB)
- For availability and reliability, extents are replicated (3 copies).
- Enables:
  - Parallel read
  - Parallel write

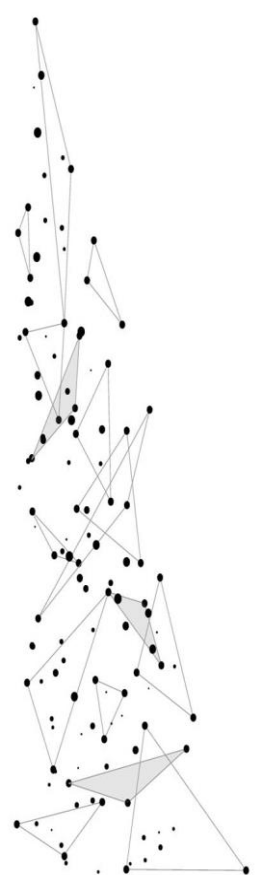


A VERY BIG FILE



# Security in Azure Data Lake Storage

- Authentication (Azure Active Directory)
- Authorization (RoleBaseAccessControl)
- Network isolation (Firewall and virtual networks )
- Data protection (Encryption +Azure Key Vault)
- Auditing
- Set file expiry



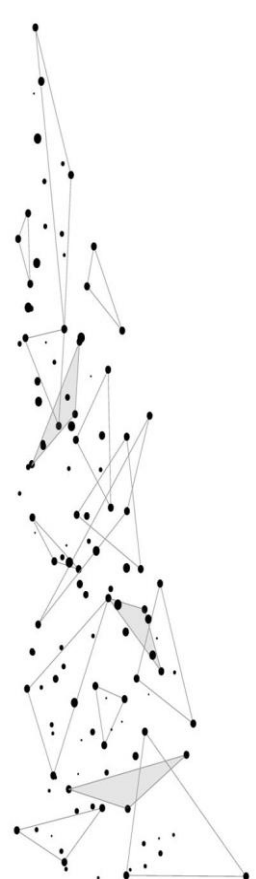
# Working with Azure Data Lake Store

## • Developer

- Azure Portal
- Azure PowerShell
- Azure Cross-platform CLI
- Using Data Lake Tools for Visual Studio
- Azure Storage Explorer

## • Production

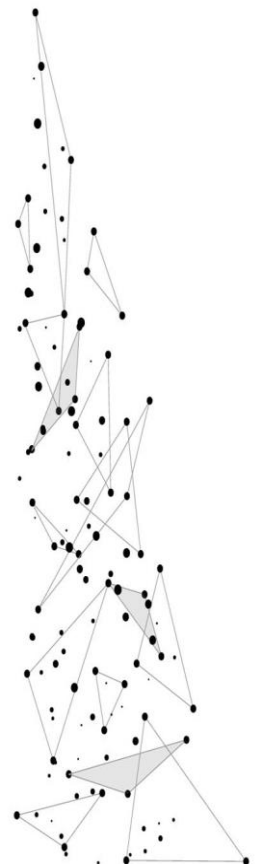
- Azure Portal
- Azure PowerShell
- Azure Cross-platform CLI
- Azure Data Factory
- Apache Sqoop
- WebHDFS-compatible REST interface



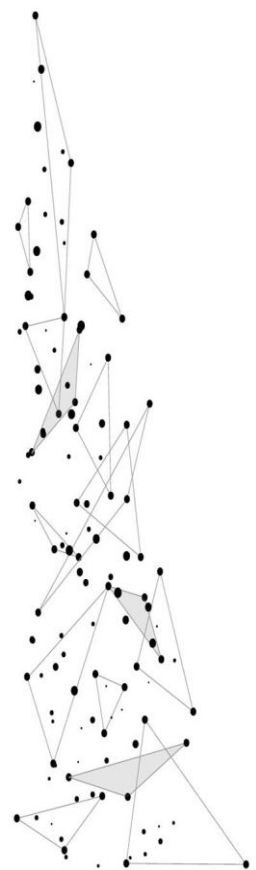
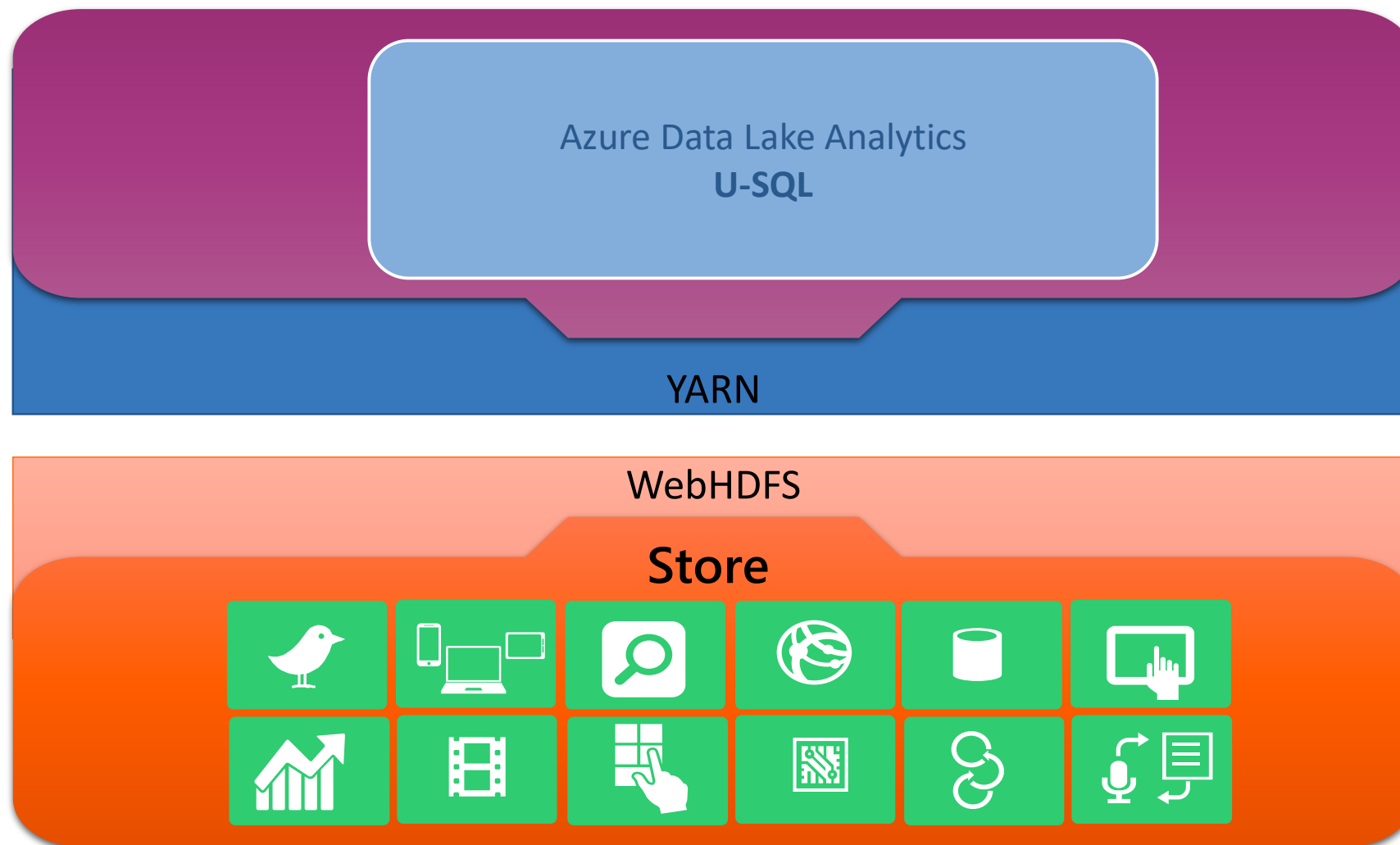
# Azure Data Lake Analytics

A distributed analytics service built on Apache YARN that dynamically scales to your needs

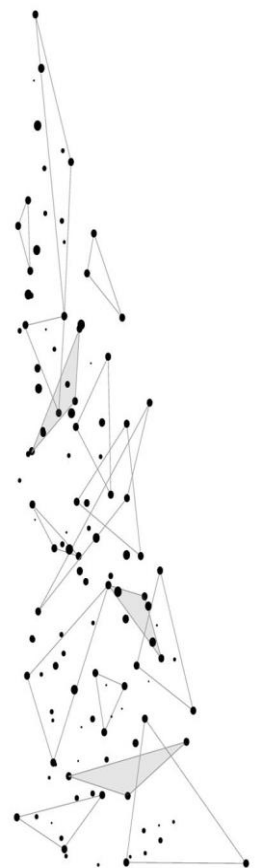
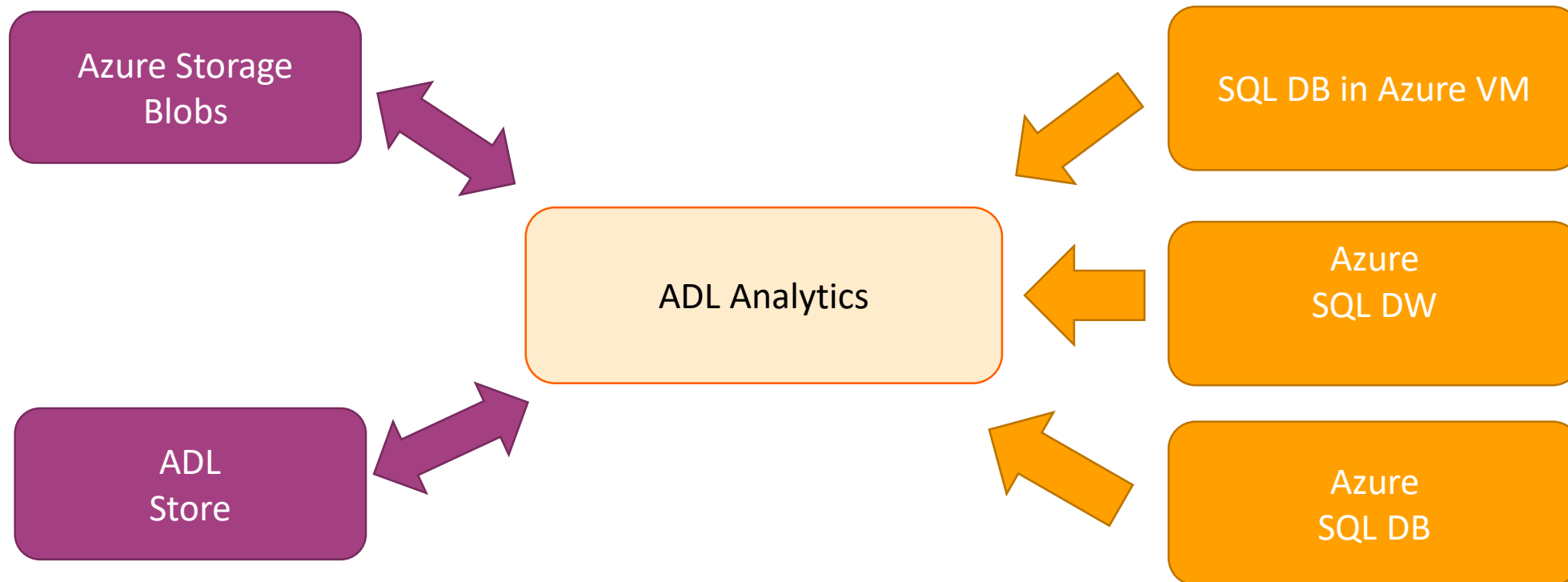
- Pay **PER QUERY** & Scale **PER QUERY**
- **FEDERATED QUERY** across Azure data sources
- Includes **U-SQL**, a language that unifies the benefits of SQL with the expressive power of C#
- No limits to **SCALE**
- Optimized to work with **ADL STORE**



# Azure Data Lake Analytics



# ADLA Sources and Sinks

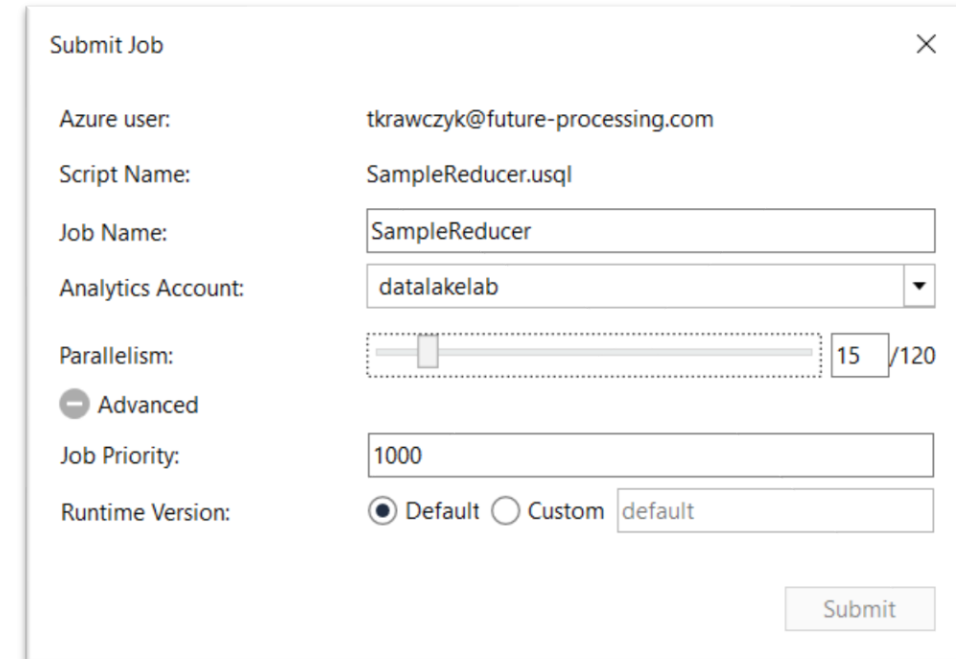


# Data Lake Analytics Runtime

1 ADLAU  $\sim$  A VM with 2 cores  
and 6 GB of memory

- Limited Network Access \*

Parallelism  $N = N$  ADLAUs



Submit Job

Azure user: tkrawczyk@future-processing.com

Script Name: SampleReducer.usql

Job Name: SampleReducer

Analytics Account: datalakelab

Parallelism:  15 / 120

☒ Advanced

Job Priority: 1000

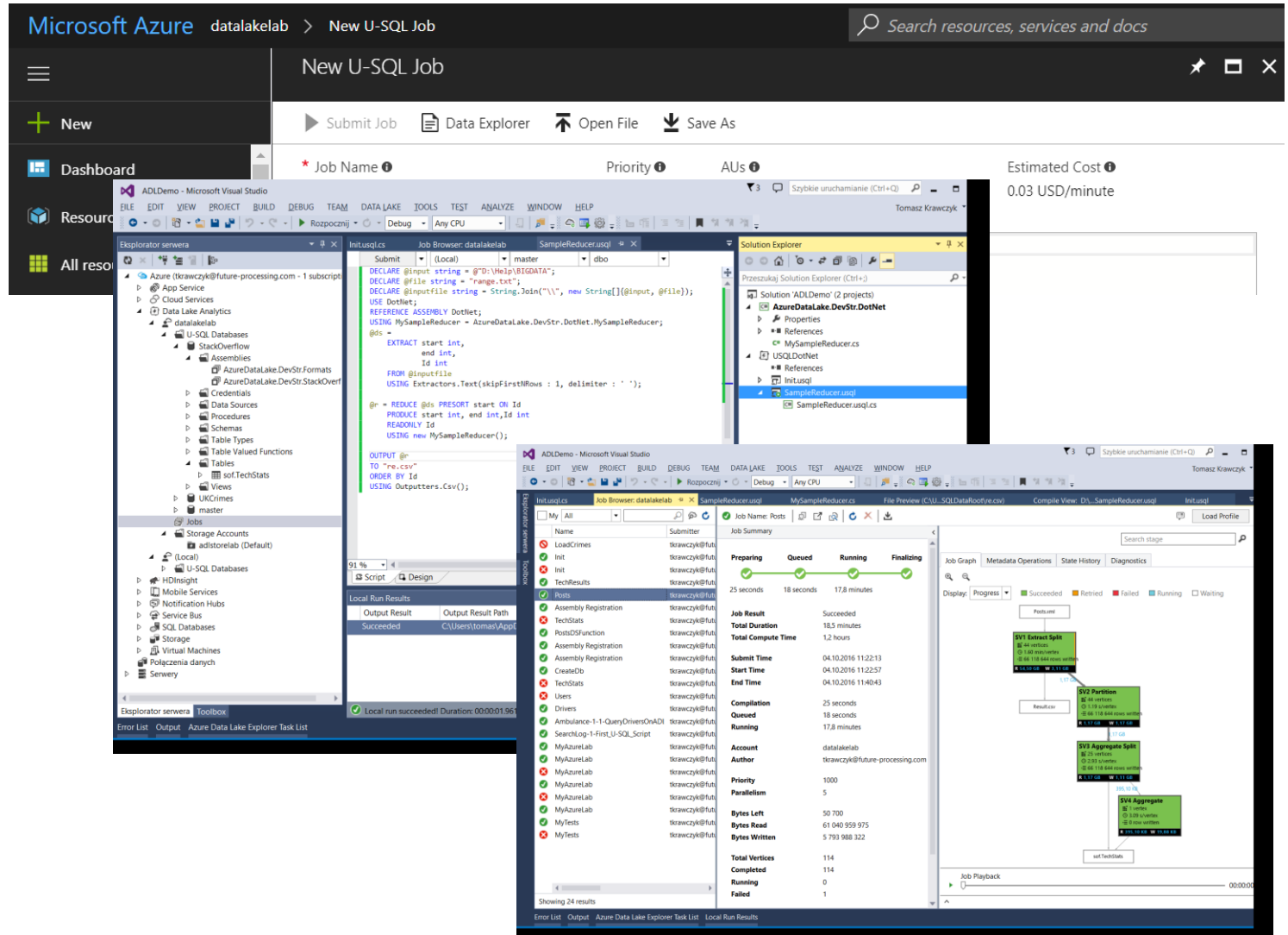
Runtime Version: ☒ Default ☐ Custom default

Submit



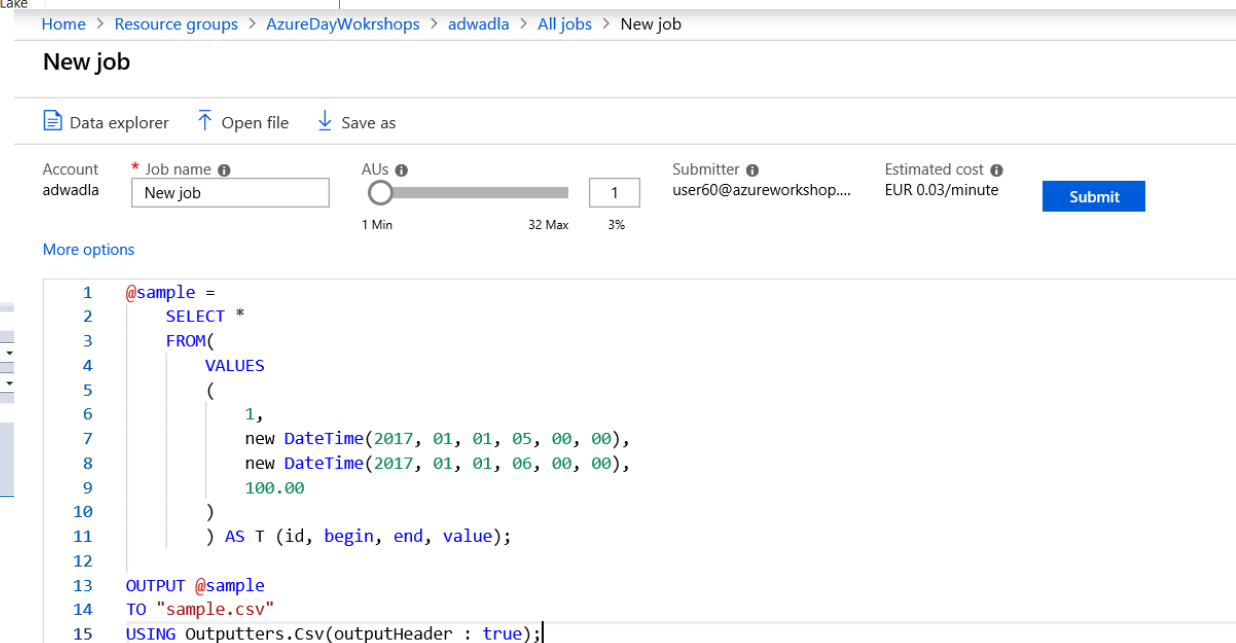
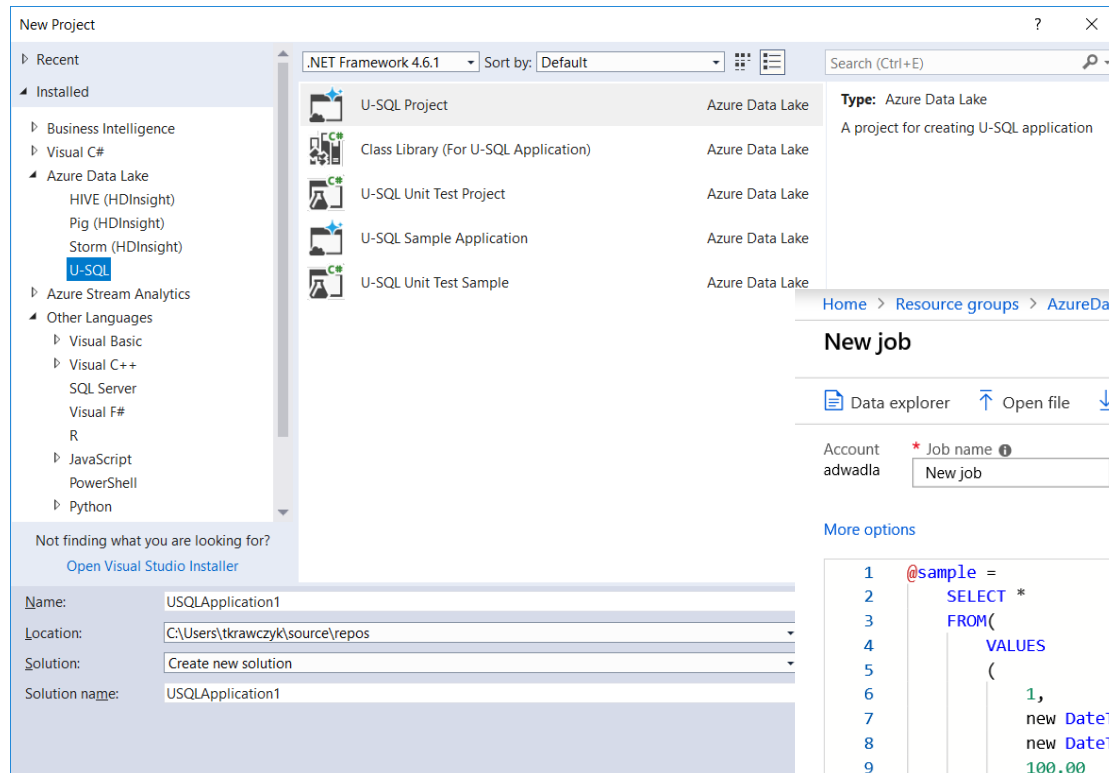
# • Data Lake Analytics Developer Tools

- Azure Portal
- Visual Studio
- ADLA Tools
- Visual Studio Code
- Power Shell



The image displays two overlapping screenshots. The top screenshot shows the Microsoft Azure portal interface for a 'New U-SQL Job'. It includes a search bar, a sidebar with navigation options like 'New', 'Dashboard', and 'Resources', and a main area with a 'Submit Job' button and a 'Data Explorer' link. The bottom screenshot shows the Visual Studio IDE with the 'ADLA Tools' extension installed. The 'Solution Explorer' on the right shows a project named 'ADLDemo' with files like 'AzureDataLake.DevStr.DotNet', 'MySampleReducer.cs', and 'Init.usql'. The 'Job Browser' on the left shows a list of jobs with columns for Name, Submitter, and Status. The 'Job Summary' pane on the right provides details for a specific job, including its duration, compute time, and status.

# • Data Lake Analytics Developer Tools

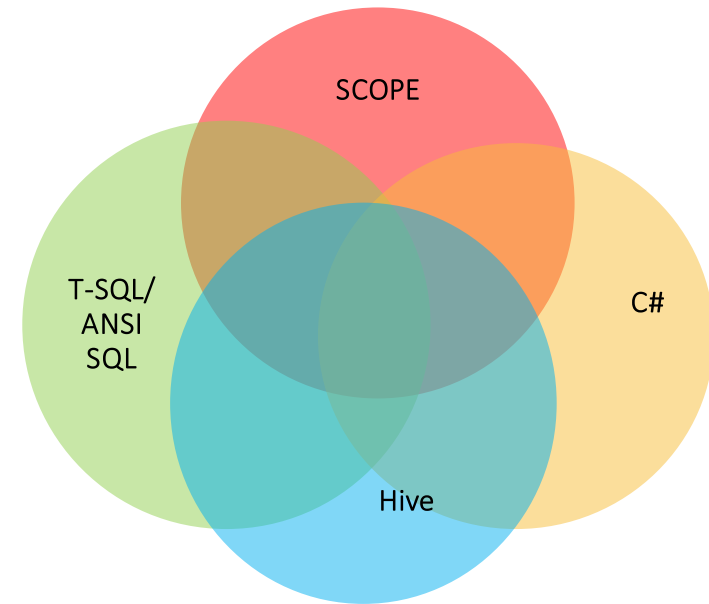


# DEMO

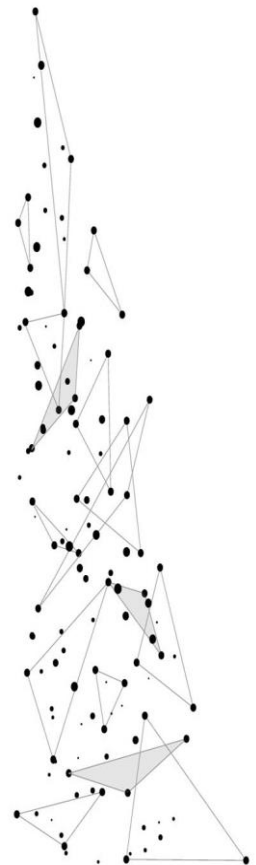
- Create ADLA
- Run Sample Job
- Visual Studio Data Lake Analytics Tools
- Task1, Task2, Task3

# U-SQL - new language for Big Data

- Familiar syntax to millions of SQL & .NET developers
- Unifies declarative nature of SQL with the imperative power of C#
- Unifies structured, semi-structured and unstructured data
- Distributed query support over all data



**SQL DECLARATIVITY + C# EXTENSIBILITY = U-SQL**



# U-SQL Language Overview

```

DECLARE @projectsInput string = @"Projects\{file}.csv";
DECLARE @eventDate DateTime = System.DateTime.Parse("2018/10/15");
DECLARE @numbers int = 2;
REFERENCE ASSEMBLY USQLCSharpDemo;
USING ImageColorsProcessor = USQLCSharpDemo.ImageColorProducer;
@projects =
    EXTRACT project string,
             startDate DateTime,
             endDate DateTime,
             file string
    FROM @projectsInput
    USING Extractors.Csv(skipFirstNRows : 1, quoting : true);
@rs =
    EXTRACT content byte[],
            fileName string
    FROM @imgFiles
    USING new BinaryExtractor();
@assignments =
    SELECT user.ToUpper() AS user,
           new SqlArray<string>(projects.Split(new char[]{'|'},
StringSplitOptions.RemoveEmptyEntries)) AS proj
    FROM @usersprojects;
@agg =
    SELECT project,
           COUNT( * ) AS units
    FROM @details WHERE project.StartsWith("My")
    GROUP BY project;
@myprojects =
    SELECT us.project,
           p.endDate
    FROM @details AS us
    JOIN
        @projects AS p
    ON p.project == us.project
    WHERE user.StartsWith("Me")
    ORDER BY p.endDate DESC
    FETCH 10 ROWS;
OUTPUT @myprojects
TO "myprojects.csv"
USING Outputters.Csv();

```

C# Types

.NET Assemblies

Apply Schema on Read

Extractor

Rowset(s)

External Extractor

SQL Dialect (SELECT FROM)

SQL Aggregation(s)

.NET Methods

SQL Dialect (JOIN, WHERE, ORDER BY ...)

Output

# U-SQL Data Types

## Numeric

- sbyte
- int
- long
- float
- double
- decimal
- short
- byte
- uint
- ulong
- ushort

## Text

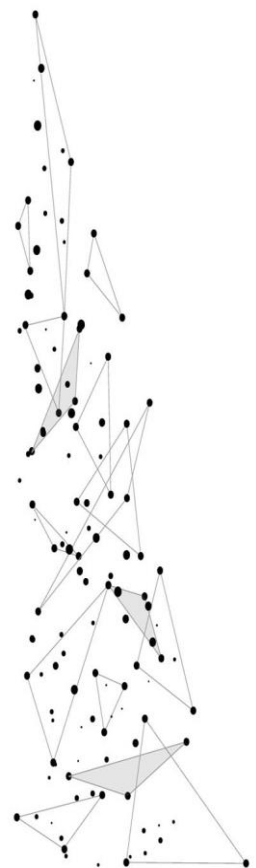
- char
- String (128 kB)
- > Complex
- MAP<k,v>
- ARRAY<v>
- > Miscellaneous
- bool
- Guid
- DateTime
- byte[]
- > RowSets

SQL.ARRAY<T> == IList<T>  
SQL.MAP<T,U> == IDictionary<T,U>

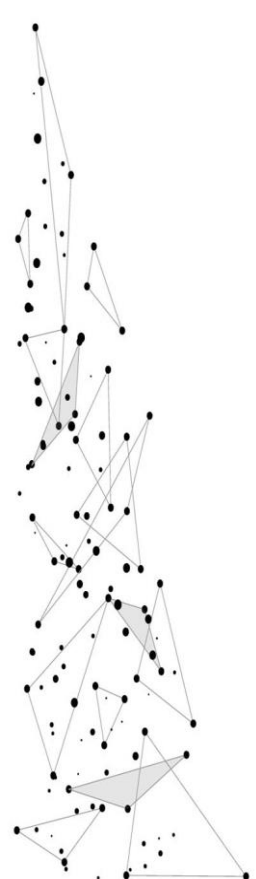
```
@m = SELECT new SqlArray<string>
      (tweet.Split( new char[]{' '}).Where(x =>
        x.StartsWith("@"))) AS mentions
```

```
FROM @t;
```

```
@m = SELECT m.Substring(1) AS m, "mention" AS category
      FROM @m
      CROSS APPLY EXPLODE(mentions)
AS t(m);
```



# U-SQL Declare Variables



```
DECLARE @text1 string = "Big Data";
DECLARE @text2 string = @"Azure as a Big Data Platform";
DECLARE @text3 char = 'a';
DECLARE @text4 string = "BEGIN" + @text1 + "END";
DECLARE @text5 string = string.Format("BEGIN{0}END", @text1);
DECLARE @text6 string = string.Join(" ", new String[]{@text1,
"2017"});

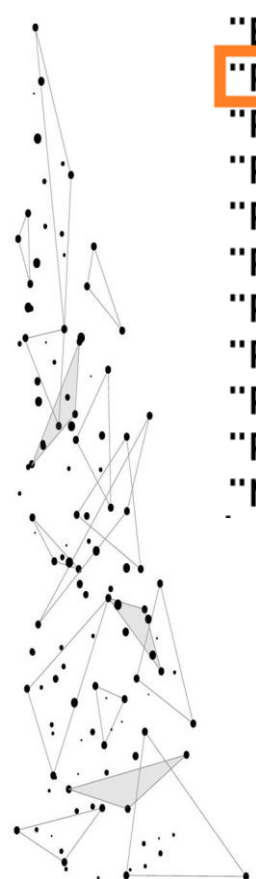
DECLARE @numeric1 sbyte = 0;
DECLARE @numeric2 short = 1;
DECLARE @numeric3 int = 2;
DECLARE @numeric4 long = 3L;
DECLARE @numeric5 float = 4.0f;
DECLARE @numeric6 double = 5.0;

DECLARE @d1 DateTime = System.DateTime.Parse("1979/03/31");
DECLARE @d2 DateTime = DateTime.Now;

DECLARE @misc1 bool = true;
DECLARE @misc2 Guid = System.Guid.Parse("BEF7A4E8-F583-4804-9711-
7E608215EBA6");
DECLARE @misc4 byte [] = new byte[] { 0, 1, 2, 3, 4};
```



# U-SQL EXTRACT



Project	StartDate	EndDate
"BigBang"	2015-01	2016-12
"BioData"	2016-08	2016-10
"Project1"	2016-01	2016-09
"Project2"	2016-11	2016-12
"Project3"	2014-01	2018-12
"Project4"	2016-01	2017-08
"Project5"	2016-01	2017-12
"Project6"	2015-01	2017-02
"Project7"	2015-01	2016-12
"Project8"	2015-08	2016-08
"Project9"	2014-01	2016-12
"NewProject"	2017-12	2018-12

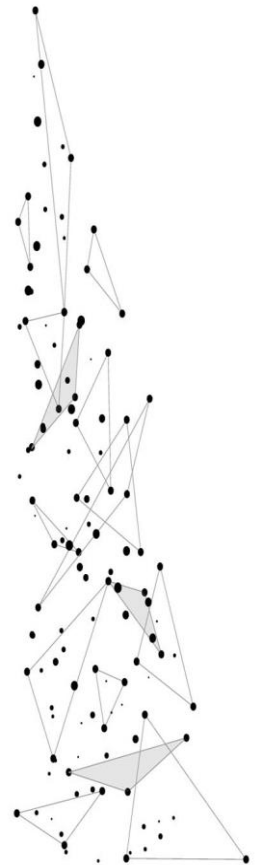
```
DECLARE @projectsInput string =  
@"Projects\projects.csv";  
@projects =  
    EXTRACT project string,  
             startDate DateTime,  
             endDate DateTime  
FROM @projectsInput  
USING Extractors.Csv(  
    skipFirstNRows : 1,  
    quoting : true);
```

Apply Schema on Read

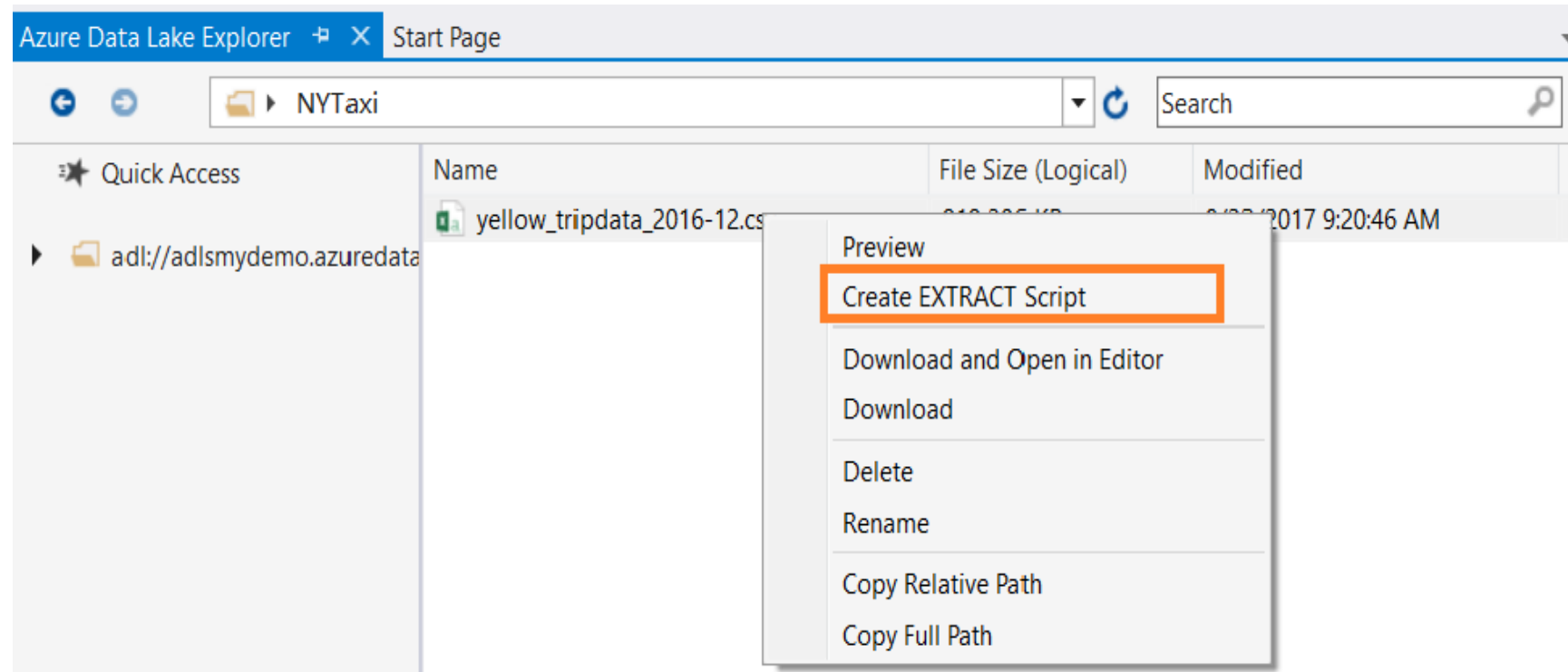
Extractor  
with additional options

# U-SQL EXTRACTORS

- List of **EXTRACTORS** AND **OUTPUTERS**
  - CSV
  - TEXT
  - TSV
  - GZIP (\*)
  - PARQUET (\*)
- **API**
  - IExtractor
  - IOutputter



# U-SQL CREATE EXTRACT SCRIPT



# U-SQL ROWSETS

```
@postCodes =  
    EXTRACT id string,  
            postcode string,  
            latitude string,  
            longitude string  
    FROM @inputPostCodes  
    USING Extractors.Csv(skipFirstNRows:1);
```

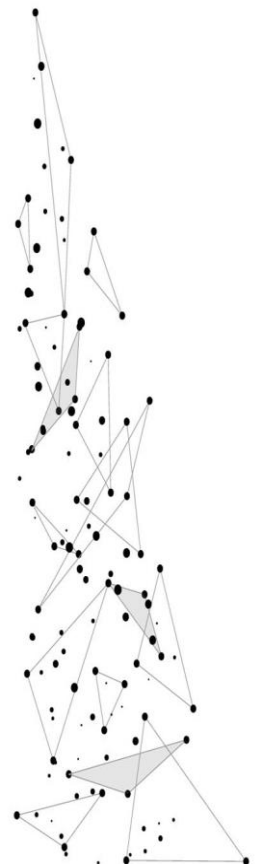
Rowset 1 (@postCodes)

```
@topCities =  
    EXTRACT id int,  
            name string,  
            population string,  
            postcode string  
    FROM @input10topCities  
    USING Extractors.Text(delimiter : ';');
```

Rowset 2 (@topCities)

```
@topCitiesWithGPS =  
    SELECT tc.name,tc.population,  
           pc.latitude,pc.longitude  
    FROM @topCities AS tc  
    JOIN  
        @postCodes AS pc  
    ON pc.postcode == tc.postcode;
```

Rowset 3 (@topCitiesWithGPS)



# U-SQL FILESETS

```
DECLARE @inputCrimes = @"mySamples/UKCrimes/  
{Date:yyyy}-{Date:MM}/{Input}-street.csv";
```

```
@crimes =
```

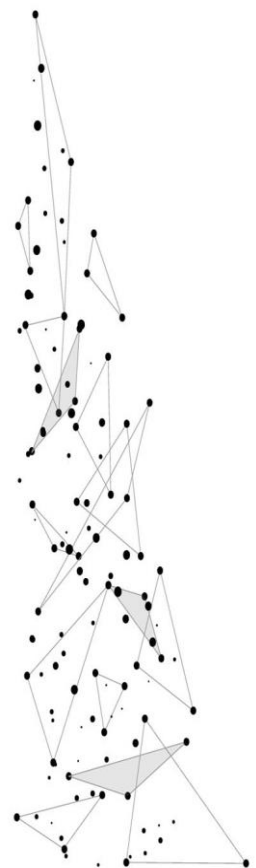
```
    EXTRACT CrimeID string,  
            Month string,  
            ReportedBy string,  
            FallsWithin string,  
            Longitude string,  
            Latitude string,  
            Location string,  
            LSOACode string,  
            LSOAName string,  
            CrimeType string,  
            LastOutcomeCategory string,  
            Context string,  
            Date DateTime,  
            Input string }
```

```
FROM @inputCrimes
```

```
USING Extractors.Csv(silent : false,skipFirstNRows:1);
```

FileSets ( {Date},{Input})

Virtual Columns



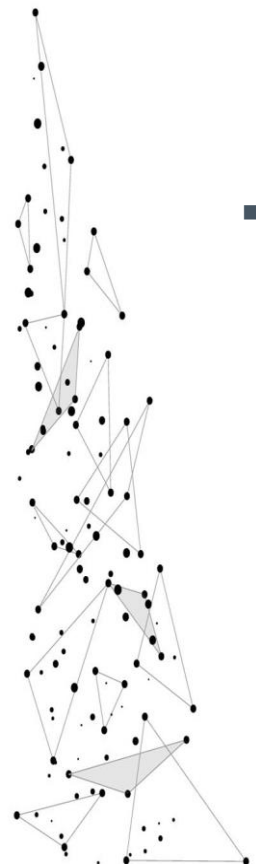
# U-SQL

## FILTERING AND SORTING

- WHERE
  - AND & OR
  - ==, >=, != (C# OPERATOR(s))
  - CONTAINS (C# string)
- ORDER BY
  - ROWSETS
    - requires a FETCH
  - OUTPUTS

## JOINS

- INNER JOIN
- FULL OUTER JOIN
- LEFT OUTER JOIN
- RIGHT OUTER JOIN
- CROSS JOIN
- LEFT SEMIJOIN
- RIGHT SEMIJOIN
- LEFT ANTISEMIJOIN
- RIGHT ANTISEMIJOIN



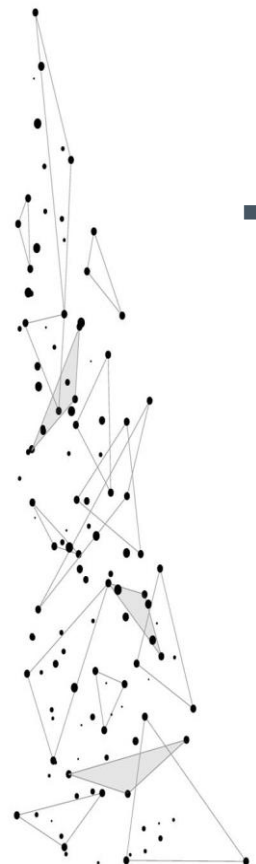
# U-SQL

## FILTERING AND SORTING

- WHERE
  - AND & OR
  - ==, >=, != (C# OPERATOR(s))
  - CONTAINS (C# string)
- ORDER BY
  - ROWSETS
    - requires a FETCH
  - OUTPUTS

## JOINS

- INNER JOIN
- FULL OUTER JOIN
- LEFT OUTER JOIN
- RIGHT OUTER JOIN
- CROSS JOIN
- LEFT SEMIJOIN
- RIGHT SEMIJOIN
- LEFT ANTISEMIJOIN
- RIGHT ANTISEMIJOIN





# U-SQL

## AGGREGATIONS

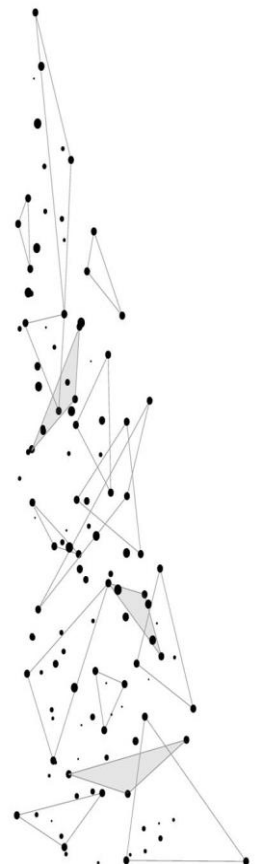
- GROUP BY
- HAVING
- MAX ,MIN ,SUM,COUNT, ARRAY\_AGG

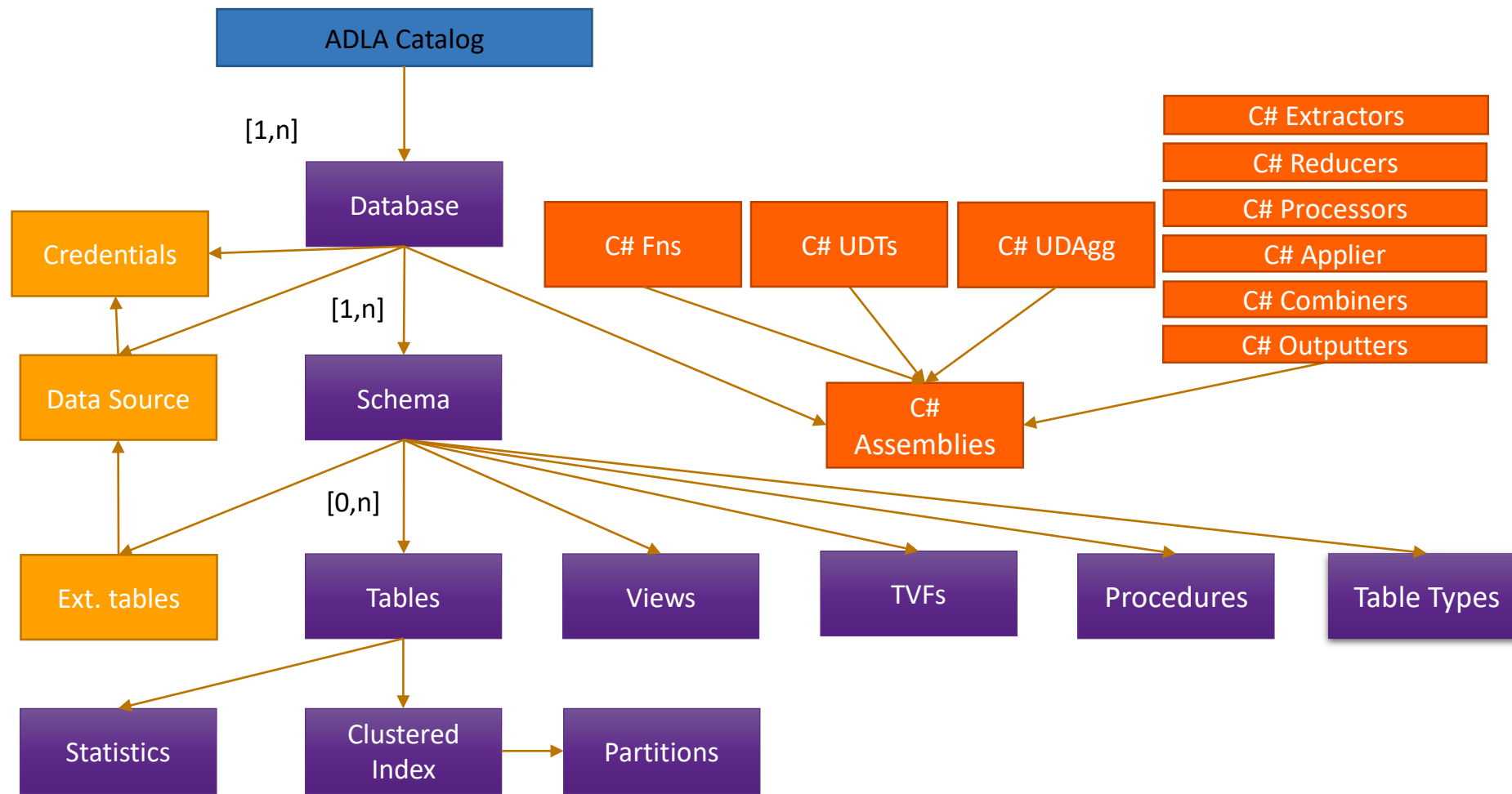
## RANKING FUNCTIONS

- RANK
- DENSE\_RANK
- NTILE
- ROW\_NUMBER

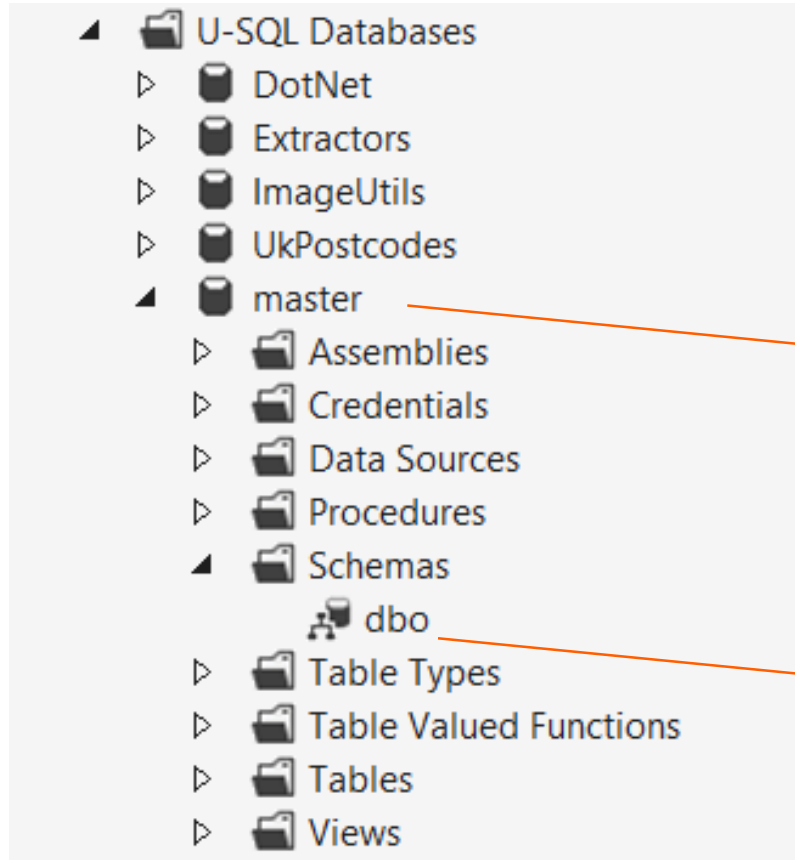
## ANALYTIC WINDOW FUNCTIONS

- CUME\_DIST
- PERCENT\_RANK
- PERCENTILE\_CONT
- PERCENTILE\_DISC
- CUME\_DIST





# U-SQL DATABASES



```
CREATE DATABASE IF NOT EXISTS UKCrimes;  
USE DATABASE UKCrimes;  
CREATE SCHEMA IF NOT EXISTS cr;
```

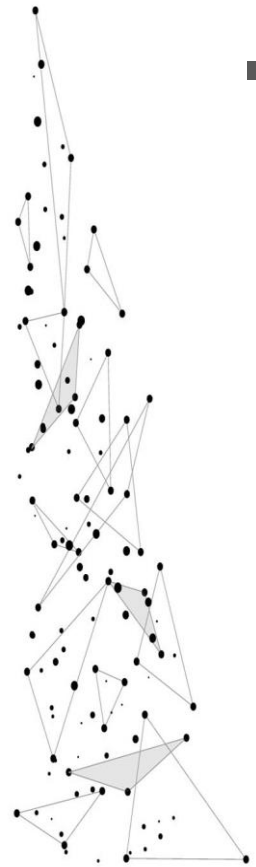
Default data base

Default schema

# U-SQL DATABASES

- **MANAGED TABLES** and EXTERNAL TABLES
- ONLY INSERT
- CONSISTS OF FOUR THINGS:
  - A NAME
  - COLUMNS
  - A CLUSTERED INDEX
  - DISTRIBUTION (PARTITIONING) SCHEME

```
CREATE TABLE T
(
    id int,
    date DateTime,
    INDEX IDX
    CLUSTERED(id)
    PARTITIONED BY (date)
    DISTRIBUTED BY
    HASH(id)
    INTO 4
);
```



# U-SQL VIEWS and FUNCTIONS

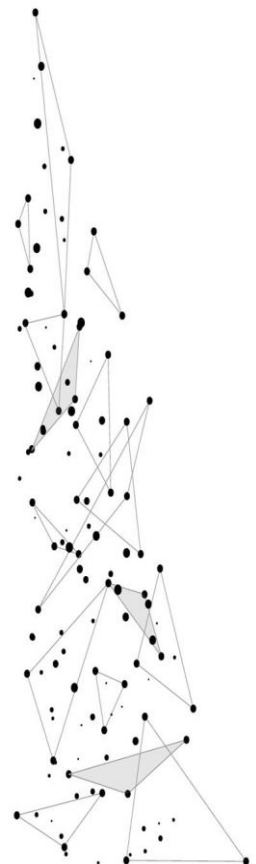
## VIEWS

```
CREATE VIEW IF NOT EXISTS vCrimes
    AS
    EXTRACT CrimeID string,
            Month string,
            Date DateTime,
            Input string
    FROM @"\\UKCrimesCities\\{Date:yyyy}-
        {Date:MM}\\{Input}-street.csv"
    USING Extractors.Csv(silent : false,
        skipFirstNRows : 1);
```

## FUNCTIONS

```
CREATE FUNCTION tvf_Crimes(@input string)
    RETURNS @result TABLE(CrimeID string,
        Month string)
    AS
    BEGIN
        @crimes =
            EXTRACT CrimeID string,
                    Month string
            FROM @input
            USING Extractors.Csv(silent : false,
                skipFirstNRows:1);

        @result = SELECT CrimeID,
                        Month
                        Input FROM @crimes;
    END;
```



# U-SQL STORED PROCEDURES

```
DROP PROCEDURE IF EXISTS DemoDb.dm.SampleSP;  
CREATE PROCEDURE DemoDb.dm.SampleSP(@startDate DateTime,  
@endDate DateTime, @outputName string)  
BEGIN
```

```
    @sample =  
        SELECT *  
        FROM(  
            VALUES  
            (  
                1,  
                new DateTime(2017, 01, 01, 05, 00, 00),  
                new DateTime(2017, 01, 01, 06, 00, 00),  
                100.00  
            )  
        )  
    )  
AS T (id, begin, end, value);
```

```
    @rs =  
        SELECT id,  
            begin,  
            end,  
            value  
        FROM @sample  
        WHERE begin >= @startDate AND end <= @endDate;
```

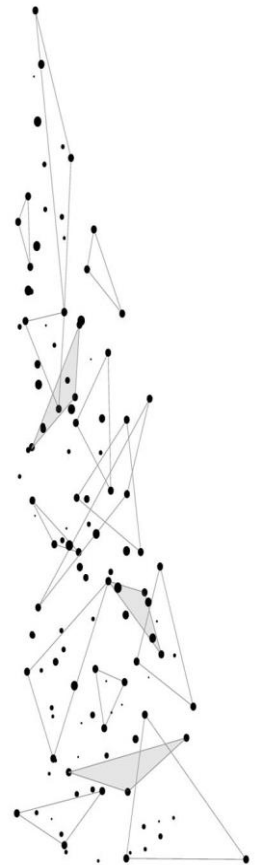
```
    OUTPUT @rs  
    TO @outputName  
    USING Outputters.Csv();  
END;
```



```
DECLARE @startDate DateTime =  
    new DateTime(2017, 01, 01, 05, 00, 00);  
DECLARE @endDate DateTime =  
    new DateTime(2017, 01, 01, 07, 00, 00);  
DECLARE @outputName string = "sp_result.csv";  
DemoDb.dm.SampleSP  
(  
    @startDate,  
    @endDate,  
    @outputName  
);
```

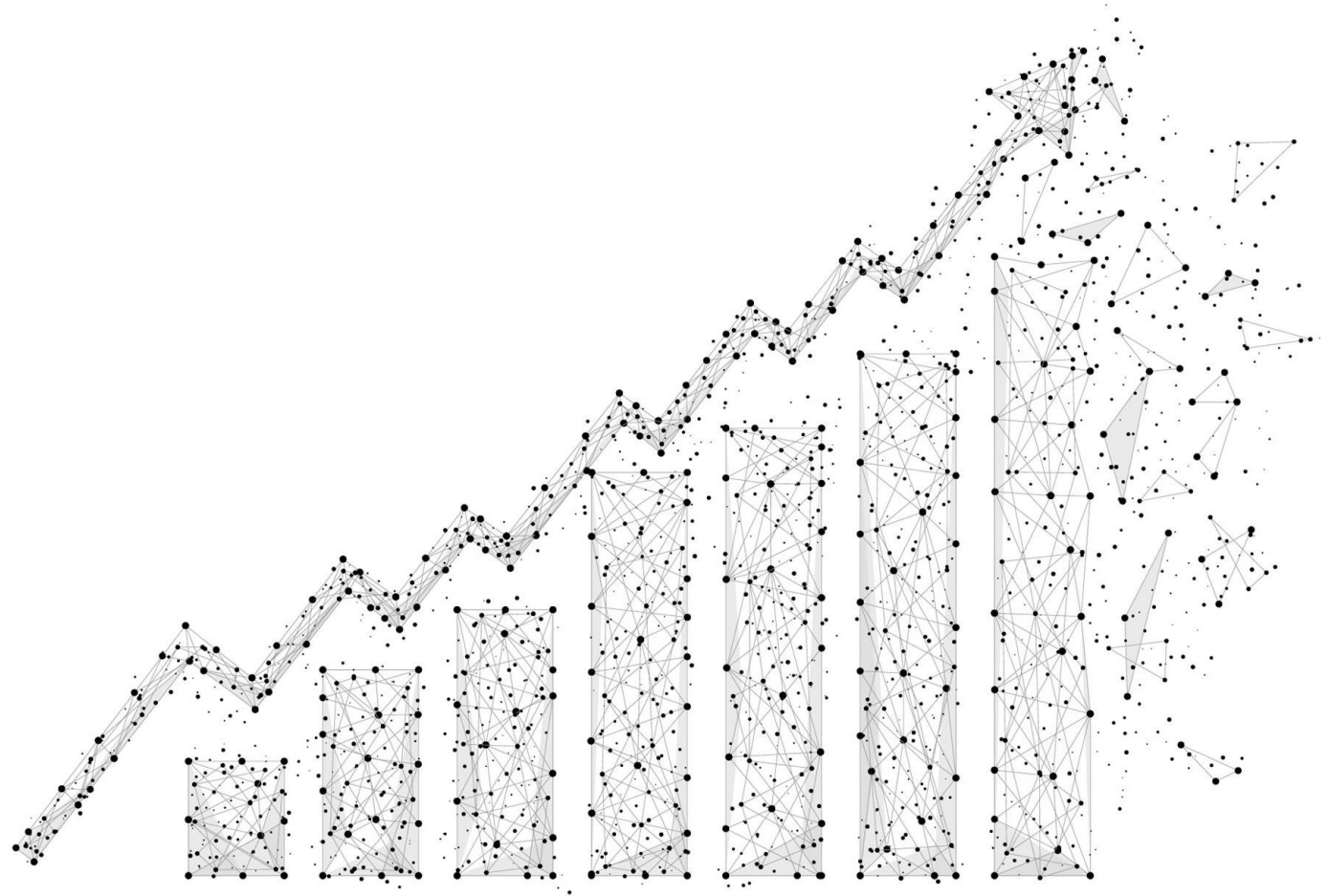
# U-SQL TOP 5'S SUPRICES FOR SQL DEVS

- U-SQL is case sensitive
- AS is not as
  - C# keywords and SQL keywords overlap
- = != ==
  - C# expression language
- null IS NOT NULL
  - C# nulls are two-valued (HasValue)
- No UPDATE, DELETE, nor MERGE



# DEMO

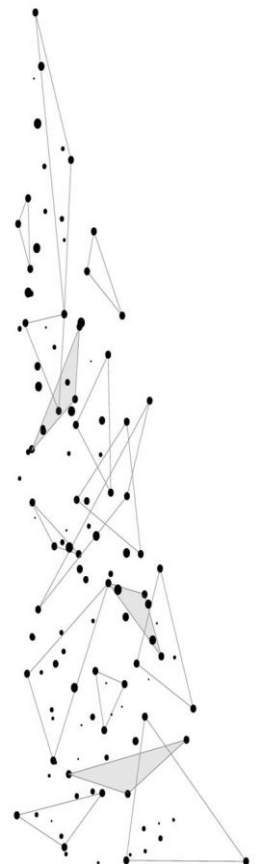
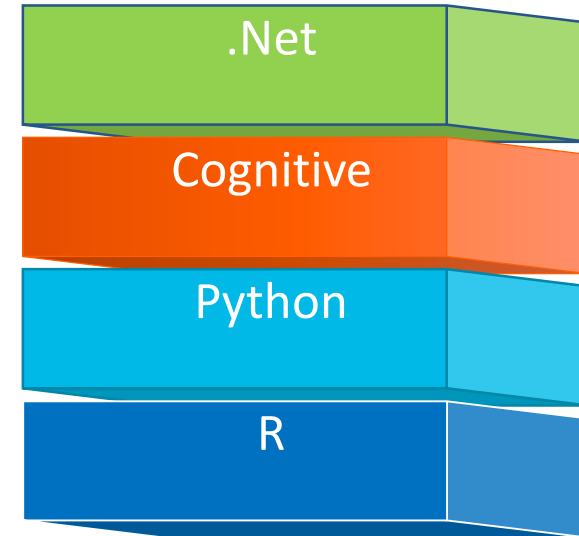
- Samples\Array
- Samples\Map
- Demo000IISLogs
- Demo001UKCrimes
- Task4





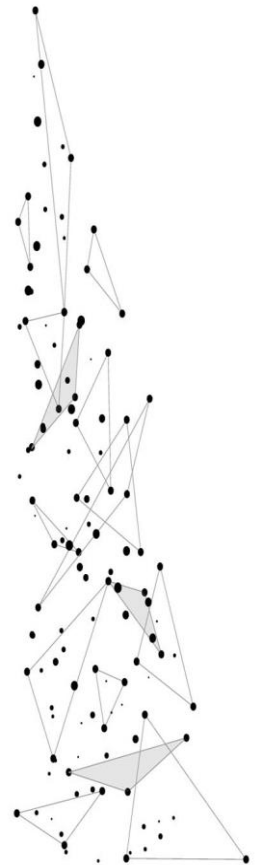
# U-SQL EXTENSIONS

- .NET Extensions
- Built-in Cognitive capabilities
- Extensions for Massively Parallel processing
  - R Language
  - Python



# U-SQL .NET EXTENTIONS

- **C# Functions/Methods**
  - C# UDTs
  - C# UDAggs
- **Extractors USQL: EXTRACT**
  - Converts files into rowset
- **Reducers USQL: REDUCE**
  - Take n rows and produce m rows (normally  $m < n$ )
- **Processors USQL: PROCESS**
  - Take one row and produce one row
  - Pass-through versus transforming
- **Appliers USQL: CROSS APPLY**
  - Take one row and produce 0 to n rows
  - Used with OUTER/CROSS APPLY
- **Combiners USQL: COMBINE**
  - Combines rowsets (like a user-defined join)
- **Outputters USQL: OUTPUT**
  - Converts rowset into file



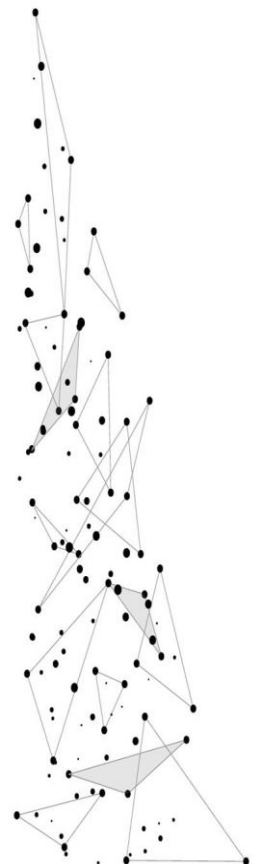
# U-SQL .Net Framework and System Assemblies

U-SQL uses the .NET Framework version 4.5  
(code in a 64-bit (x64) format)

- Preloaded System Assemblies
  - mscorlib.dll
  - System.dll
  - System.Core.dll
  - System.Data.dll
  - **Microsoft.Analytics.Interfaces.dll**
  - **Microsoft.Analytics.Types.dll**

Example:

```
REFERENCE SYSTEM ASSEMBLY [System.XML];
```



# U-SQL Creating Extensions

## ■ Methods

- Class + static method

```
public class IpConverter
{
    public static string ToIp4Format(string ip)
    {
        return IPAddress.Parse(ip).MapToIPv4().ToString();
    }
}
```

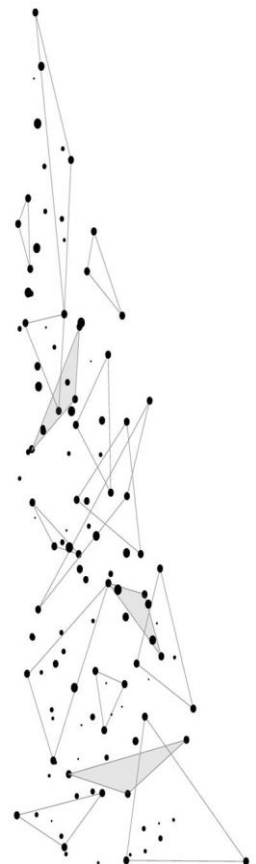
## ■ Extractors, Reducers, Processors, Appliers, Combiners, Outputters

- New project type: Class Library (For U-SQL Applications)
- Abstract classes
  - IExtractor IReducer, IProcessor, IApplier, ICombiners, IOutputters

```
public class NameReverseProcessor : IProcessor
{
    public override IRow Process(IRow input, IUpdatableRow output)
    {
        var s = input.Get<string>("name");
        output.Set<string>("reversed", Reverse(s));
        return output.AsReadOnly();
    }
    private static string Reverse(string s)
    {
        char[] charArray = s.ToCharArray();
        Array.Reverse(charArray);
        return new string(charArray);
    }
}
```

# U-SQL Registering .Net Assemblies

```
DECLARE @AssemblyPath string = @"Assemblies/";  
DECLARE @AssemblyExt string =  
@AssemblyPath+"ADLAEExt.dll";  
USE 4Developers;  
DROP ASSEMBLY IF EXISTS ADLAEExt;  
CREATE ASSEMBLY ADLAEExt FROM @AssemblyExt;
```



### Assembly Registration

ADLA Account: datalakelab  
Database: StackOverflow

Load assembly from path: D:\Repos\My\AzureDataLake\scr\ADLSamples ...

Assembly Name: AzureDataLake.DevStr.Formats

☒ Replace assembly if it already exists  
Replace assembly (and its included files if any) if it is already registered.  
Note that replacing an existing assembly may remove code that other assemblies depend on.

☐ Managed Dependencies

Name	Path
------	------

Add

☐ Additional Files

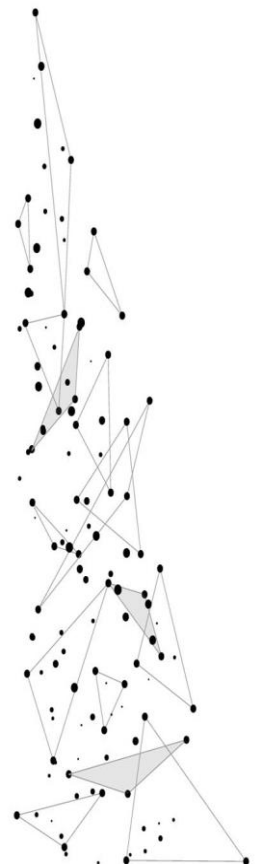
Parallelism:  1 / 120

☐ Advanced

Input Status: Ready

Submit

# U-SQL Using .Net Assemblies



```
USE DATABASE [4Developers];  
REFERENCE ASSEMBLY [ADLAEExt];  
  
USING IpConverter = ADLAEExt.Utills.IpConverter;  
  
@ds =  
    SELECT IpConverter.ToIp4Format(c_ip) AS Ip,  
           date.Date AS Date  
    FROM @iisLogs;
```

Use data base (optional)

Add reference to assembly

Create alias

Use method from assembly

# U-SQL .Net Extractor

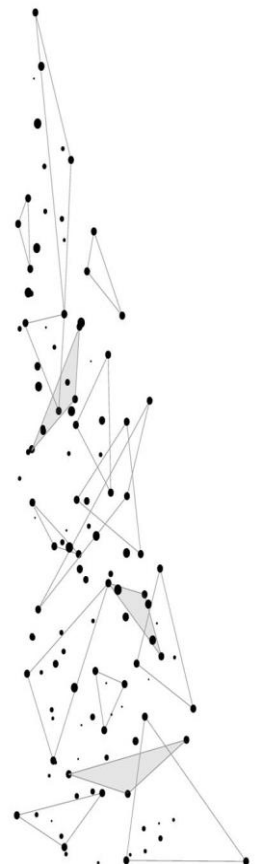
```
[SqlUserDefinedExtractor(AtomicFileProcessing = true)]
public class BinaryContentExtractor : IExtractor
{
    public override IEnumerable<IRow> Extract(IUnstructuredReader
input, IUpdatableRow output)
    {
        using (var ms = new MemoryStream())
        {
            input.BaseStream.CopyTo(ms);
            var content = ms.ToArray();
            output.Set(0, content);
            yield return output.AsReadOnly();
        }
    }
}
```

WHOLE FILE

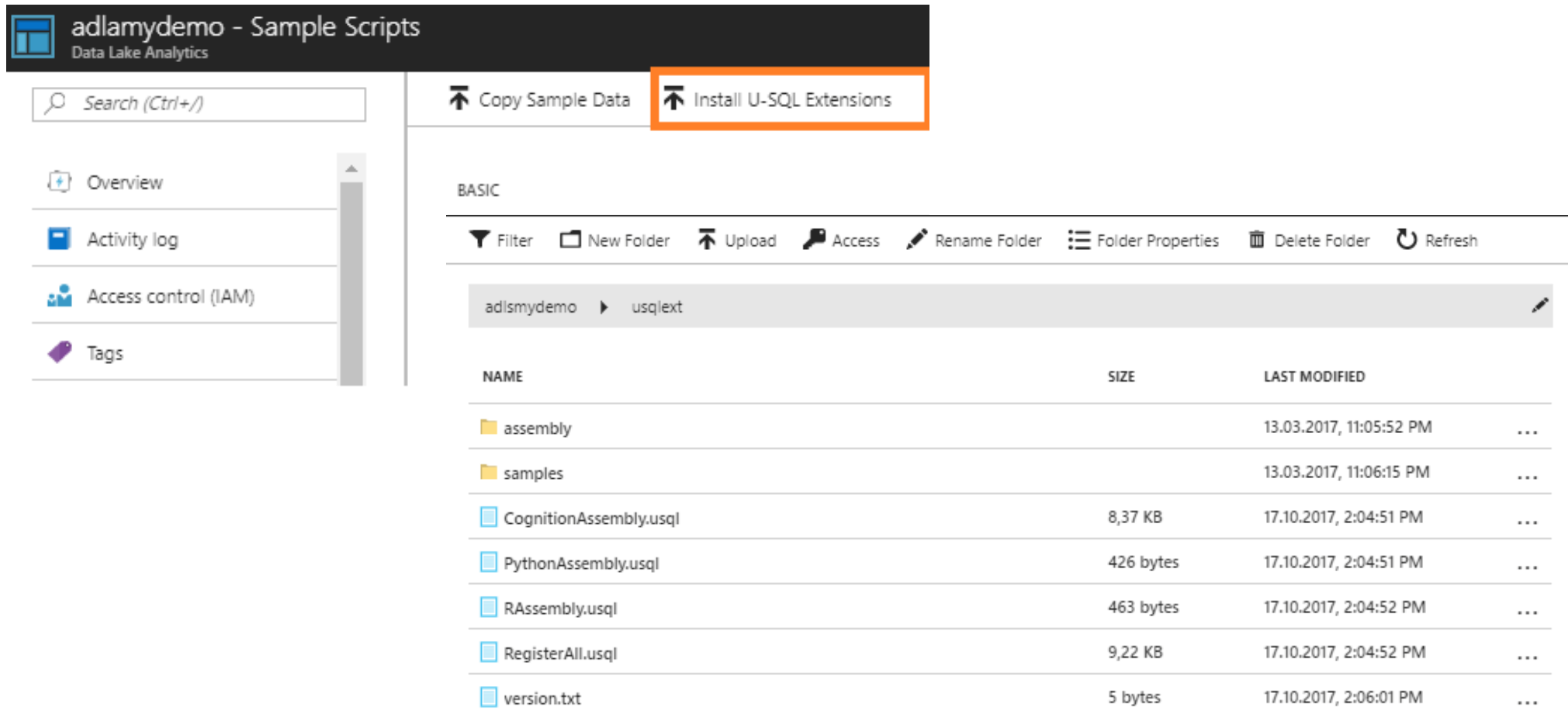
RETURN CONTENT IN FIRST  
COLUMN



```
@rs =
    EXTRACT content byte[],
            fileName string
    FROM @imgFiles
    USING new BinaryExtractor();
```



# U-SQL Python, R Language, Cognitive



adlmydemo - Sample Scripts  
Data Lake Analytics

Search (Ctrl+/)

Overview

Activity log

Access control (IAM)

Tags

Copy Sample Data

**Install U-SQL Extensions**

BASIC

Filter New Folder Upload Access Rename Folder Folder Properties Delete Folder Refresh

adlmydemo ▶ usqltext

NAME	SIZE	LAST MODIFIED
assembly		13.03.2017, 11:05:52 PM ...
samples		13.03.2017, 11:06:15 PM ...
CognitionAssembly.usql	8,37 KB	17.10.2017, 2:04:51 PM ...
PythonAssembly.usql	426 bytes	17.10.2017, 2:04:51 PM ...
RAssembly.usql	463 bytes	17.10.2017, 2:04:52 PM ...
RegisterAll.usql	9,22 KB	17.10.2017, 2:04:52 PM ...
version.txt	5 bytes	17.10.2017, 2:06:01 PM ...



# U-SQL Python, R Language, Cognitive

- Register Assemblies
  - CognitionAssembly.usql
  - PythonAssembly.usql
  - RAssembly.usql
  - RegisterAll.usql

```
CREATE DATABASE IF NOT EXISTS master;  
USE DATABASE master;
```

```
DROP ASSEMBLY IF EXISTS [ExtPython];  
CREATE ASSEMBLY IF NOT EXISTS [ExtPython]  
FROM @"/usqlext/assembly/python/ExtPy.dll"  
WITH ADDITIONAL_FILES =  
(  
    @"/usqlext/assembly/python/ExtPy.pdb",  
    @"/usqlext/assembly/python/UsqlPythonInvokePackage.zip",  
    @"/usqlext/assembly/python/UsqlPythonDeployPackage.zip",  
    @"/usqlext/assembly/python/version.python"  
);
```

```
CREATE DATABASE IF NOT EXISTS master;  
USE DATABASE master;
```

```
DROP ASSEMBLY IF EXISTS ExtR;  
CREATE ASSEMBLY IF NOT EXISTS ExtR  
FROM @"/usqlext/assembly/R/ExtR.dll"  
WITH ADDITIONAL_FILES = (  
    @"/usqlext/assembly/R/ExtR.pdb",  
    @"/usqlext/assembly/R/DynamicInterop.dll",  
    @"/usqlext/assembly/R/RDotNet.dll",
```

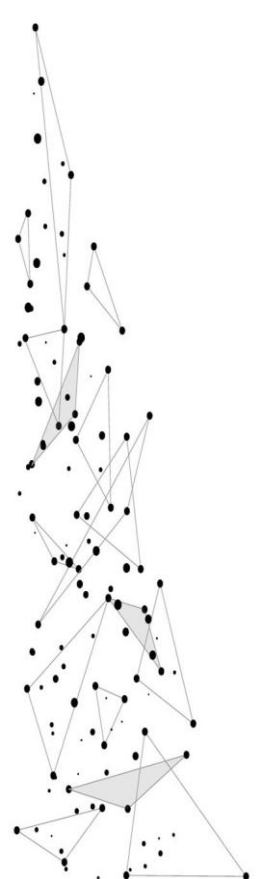
# U-SQL Cognitive Extensions

- Built-in Cognitive capabilities
  - Imaging: Detect faces
  - Imaging: Detect emotion
  - Imaging: Detect objects (tagging)
  - Imaging: OCR
  - Text: Key Phrase Extraction
  - Text: Sentiment Analysis



# U-SQL R and Python Extensions

## ■ BASED ON REDUCER(s)



```
@PyOutput =  
  REDUCE @Extended  
  ON Par  
  PRODUCE Par int,  
           SepalLength double,  
           SepalWidth double,  
           PetalLength double,  
           PetalWidth double,  
           Species string,  
           SepalRatio double,  
           PetalRatio double  
  
  USING new Extension.Python.Reducer(pyScript : @myPyScript);
```

Python Reducer

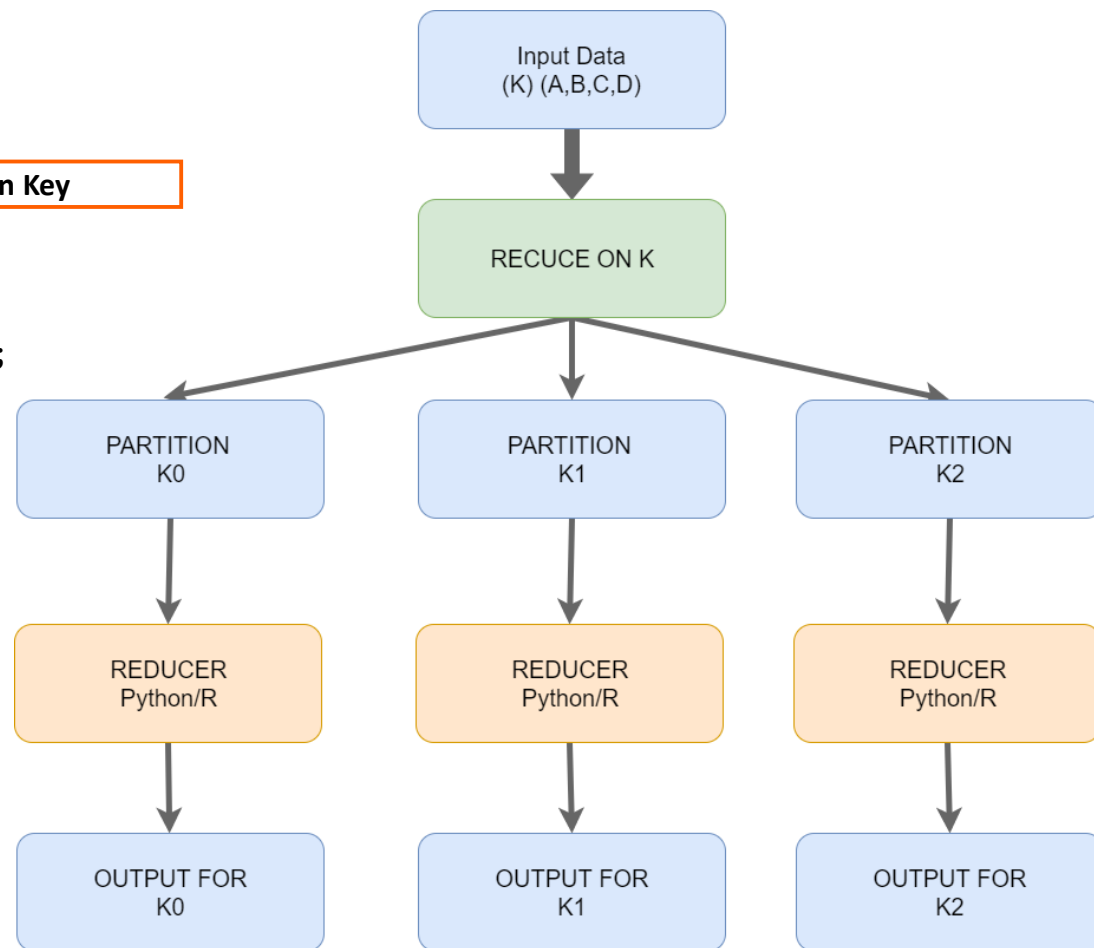
```
@ROutput =  
  REDUCE @PyOutput  
  ON Par  
  PRODUCE Par,  
           RowId int,  
           ROutput string  
  
  READONLY Par  
  USING new Extension.R.Reducer(command : @myRScript,  
                                 rReturnType:"charactermatrix",  
                                 stringsAsFactors:true);
```

Python Reducer

# U-SQL REDUCER

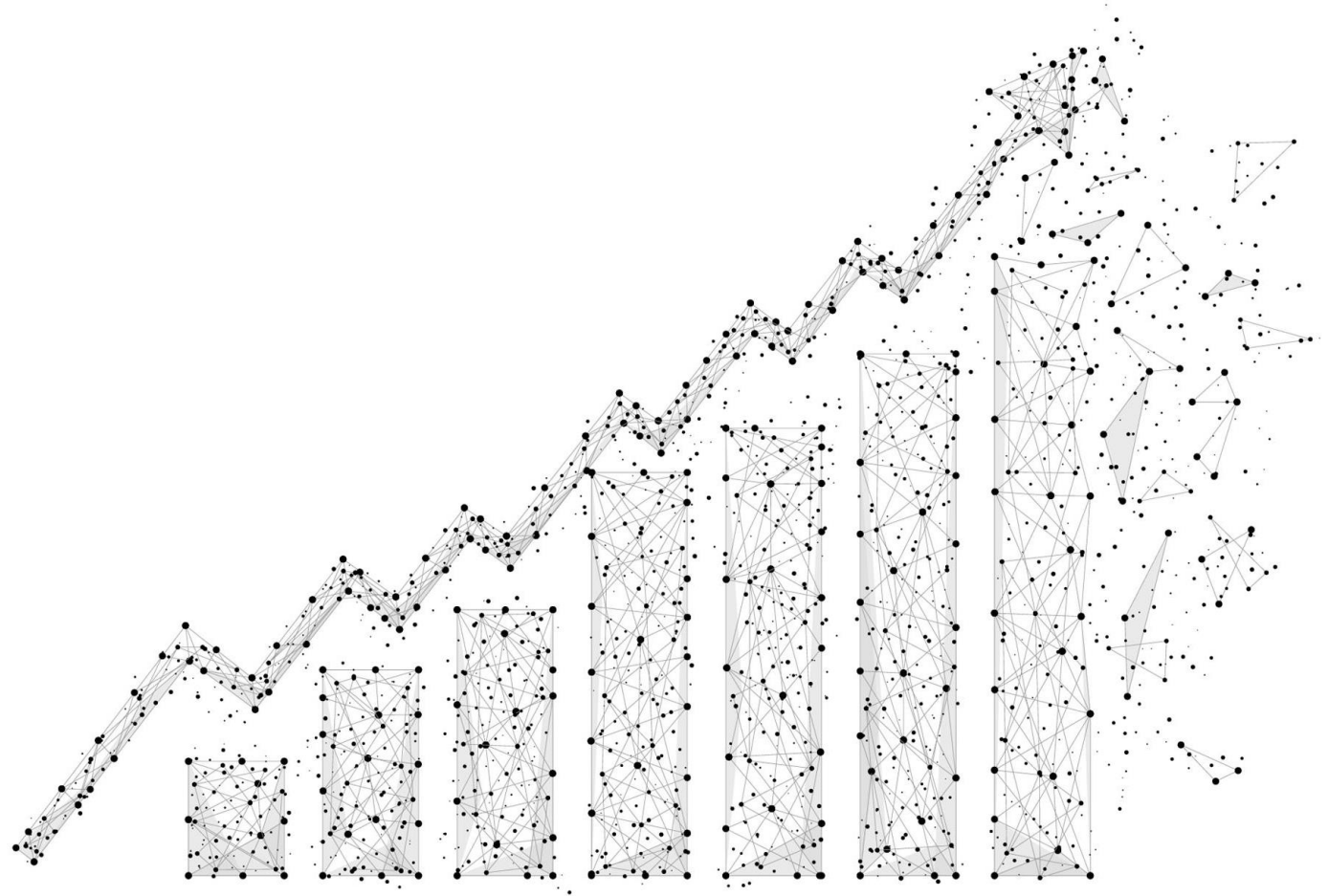
```
@ROutput =  
  REDUCE @PyOutput  
  ON Par  
  PRODUCE Par,  
    RowId int,  
    ROutput string  
  READONLY Par  
  USING new Extension.R.Reducer(command : @myRScript);
```

Partition Key



# DEMO

- Samples\Array
- Samples\Map
- Demo003Stackoverflow
- Demo004Cognitive
- Demo005PythonDeploy
- Demo006GetADLUInfo
- Task5
- Task6
- Task7

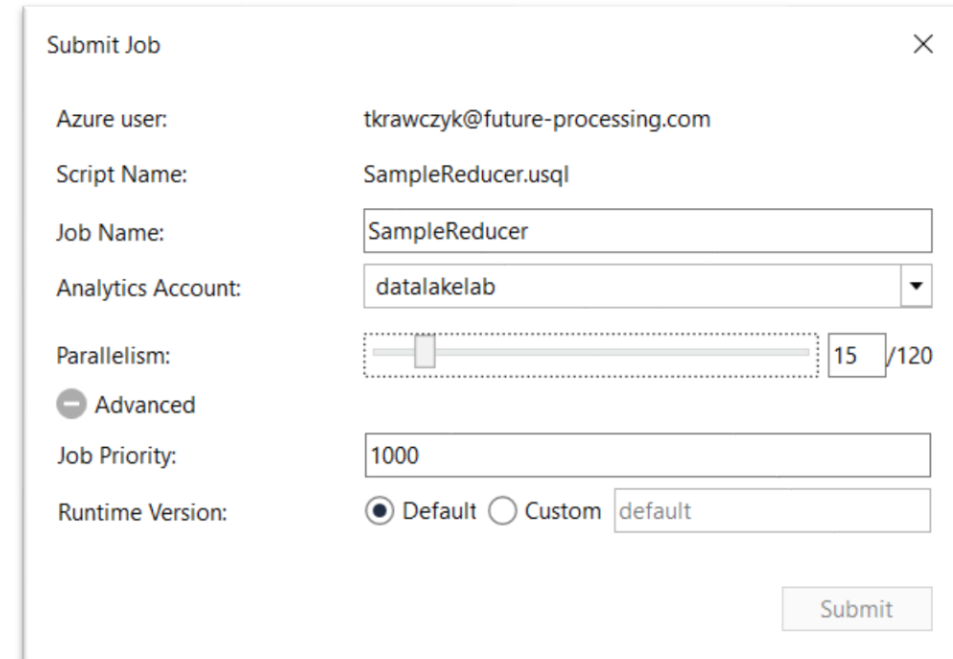


# Data Lake Analytics Runtime

1 ADLAU  $\sim$  A VM with 2 cores  
and 6 GB of memory

- Limited Network Access \*

Parallelism  $N = N$  ADLAUs



Submit Job

Azure user: tkrawczyk@future-processing.com

Script Name: SampleReducer.usql

Job Name: SampleReducer

Analytics Account: datalakelab

Parallelism:  15 / 120

☒ Advanced

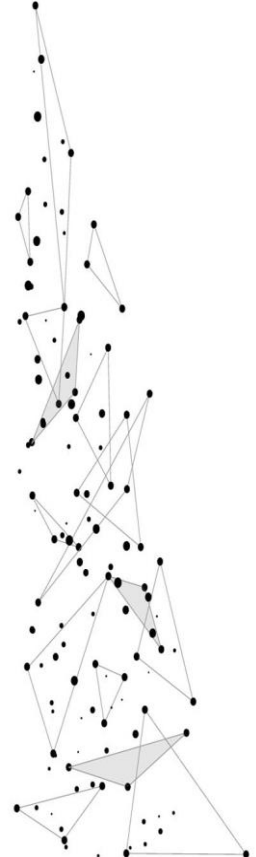
Job Priority: 1000

Runtime Version: ☒ Default ☐ Custom default

Submit

# Data Lake Analytics Runtime

## ADLUA (VERTEX) = Virtual Machine



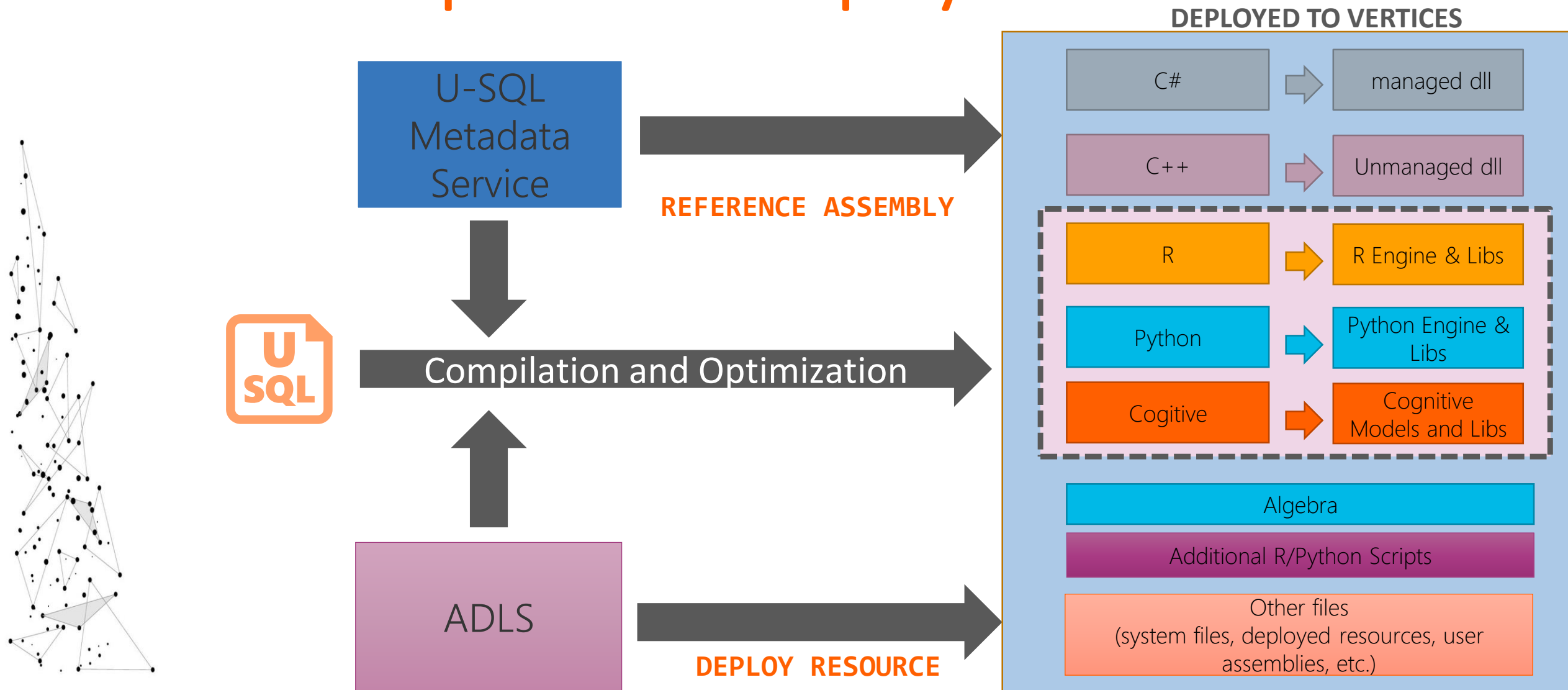
```
public static string GetVMInfo()
{
    var sb = new StringBuilder();
    var myManagementClass = new
        ManagementClass("Win32_PerfRawData_Counters_HyperVDynamicMemoryIntegrationService");
    var myManagementCollection =
        myManagementClass.GetInstances();
    var myProperties =
        myManagementClass.Properties;
    var myPropertyResults =
        new Dictionary<string, object>();

    foreach (var obj in myManagementCollection)
    {
        foreach (var myProperty in myProperties)
        {
            myPropertyResults.Add(myProperty.Name,
                obj.Properties[myProperty.Name].Value);
        }
    }

    foreach (var myPropertyResult in myPropertyResults)
    {
        var item = $"{myPropertyResult.Key}:{myPropertyResult.Value}";
        sb.AppendLine(item);
    }
    return sb.ToString();
}
```

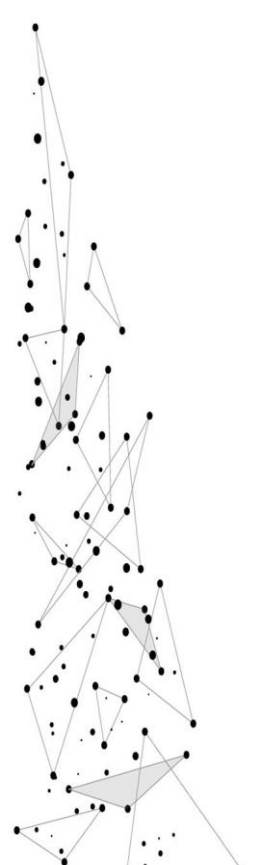
```
Processor: Intel(R) Xeon(R) CPU E5-2673 v3 @ 2.40GHz;
Microsoft Windows Server 2012 R2 Datacenter
OSArchitecture: 64-bit
SystemDevice: \Device\HarddiskVolume1
SystemDirectory: C:\Windows\system32
SystemDrive: C:
TotalSwapSpaceSize:
TotalVirtualMemorySize: 8314420
TotalVisibleMemorySize: 2096692
```

# ADLA Compilation & Deployment

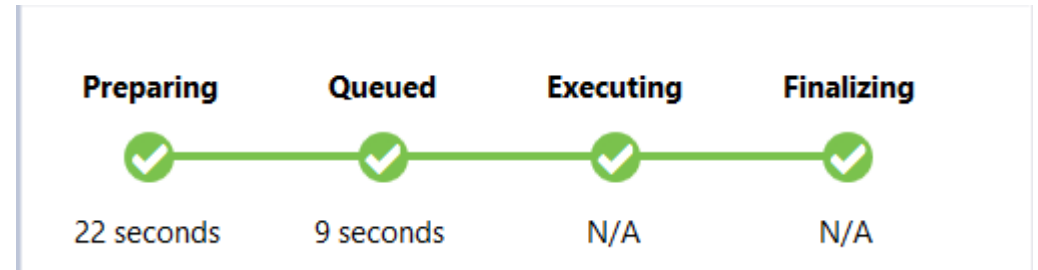




# ADLA Job Lifecycle States

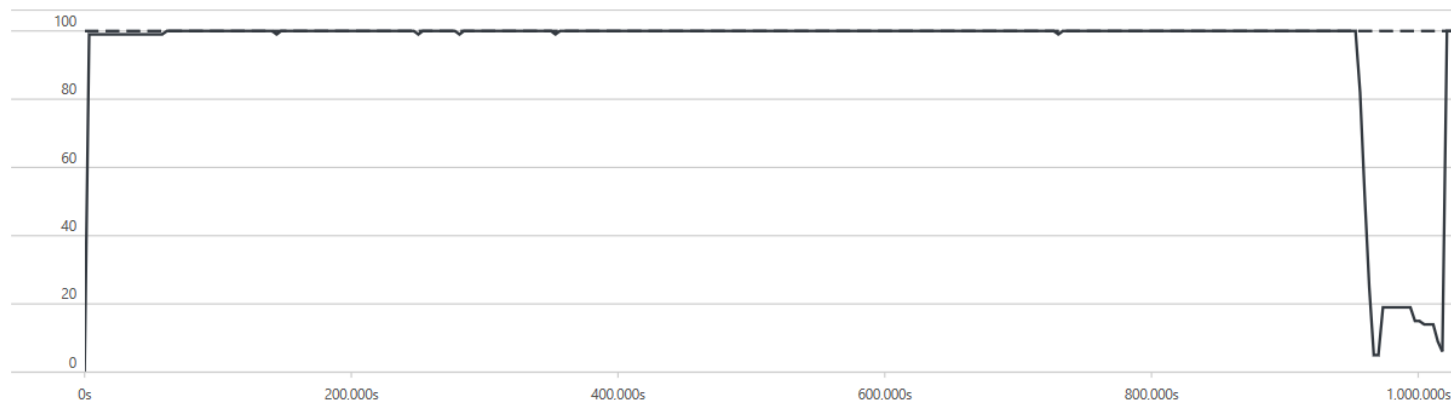
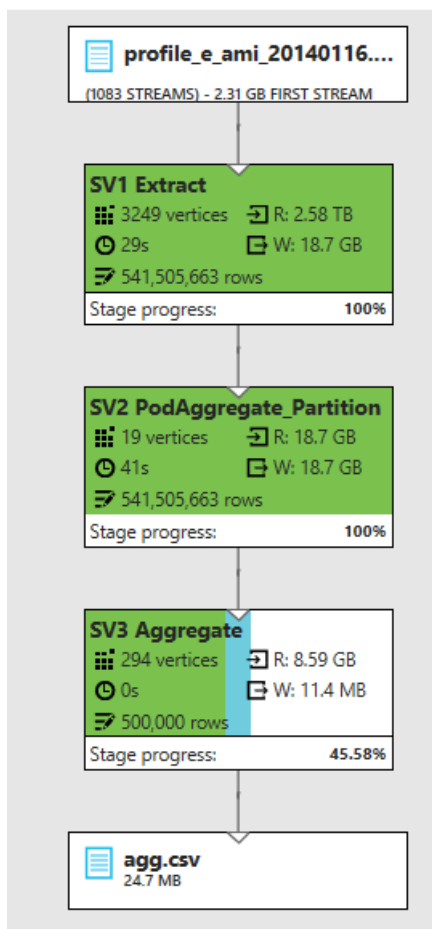
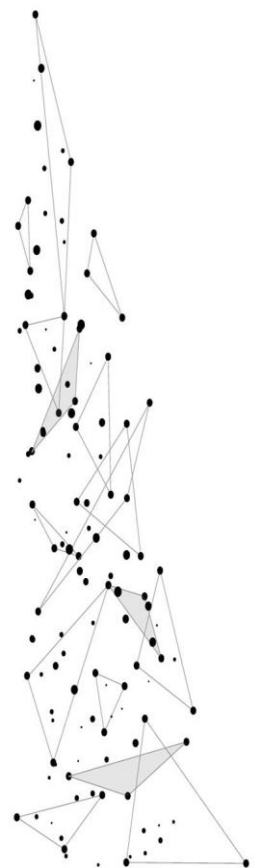


```
@output =  
    SELECT  
        MAX(Duration) AS DurationMax,  
        MIN(Duration) AS DurationMin,  
        AVG(Duration) AS DurationAvg,  
        SUM(Duration) AS DurationSum,  
        VAR(Duration) AS  
        DurationVariance,  
        STDEV(Duration) AS DurationStDev,  
    FROM @searchlog  
    GROUP BY Region  
    HAVING DurationMin > 1;
```



Result = Succeeded | Failed | Cancelled

# Azure Data Lake Analytics Optimization



Actual	
Used	Allocated
AUs allocated	100
AU-hours	28.57
Run time	17min 8s
Estimated cost	USD 42.85
Efficiency	N/A

Balanced	
Used	Allocated
AUs allocated	1105
AU-hours	45.24
Run time	2min 27s
Estimated cost	USD 67.86
Efficiency	60%
<a href="#">Select</a>	

Fast	
Used	Allocated
AUs allocated	1381
AU-hours	50.4
Run time	2min 11s
Estimated cost	USD 75.61
Efficiency	54%
<a href="#">Select</a>	

Custom	
Used	Allocated
AUs allocated	1
AU-hours	27.14
Run time	1d 3hr
Estimated cost	USD 40.71
Efficiency	100%
<a href="#">Select</a>	

# Azure Data Lake vs Spark

Why



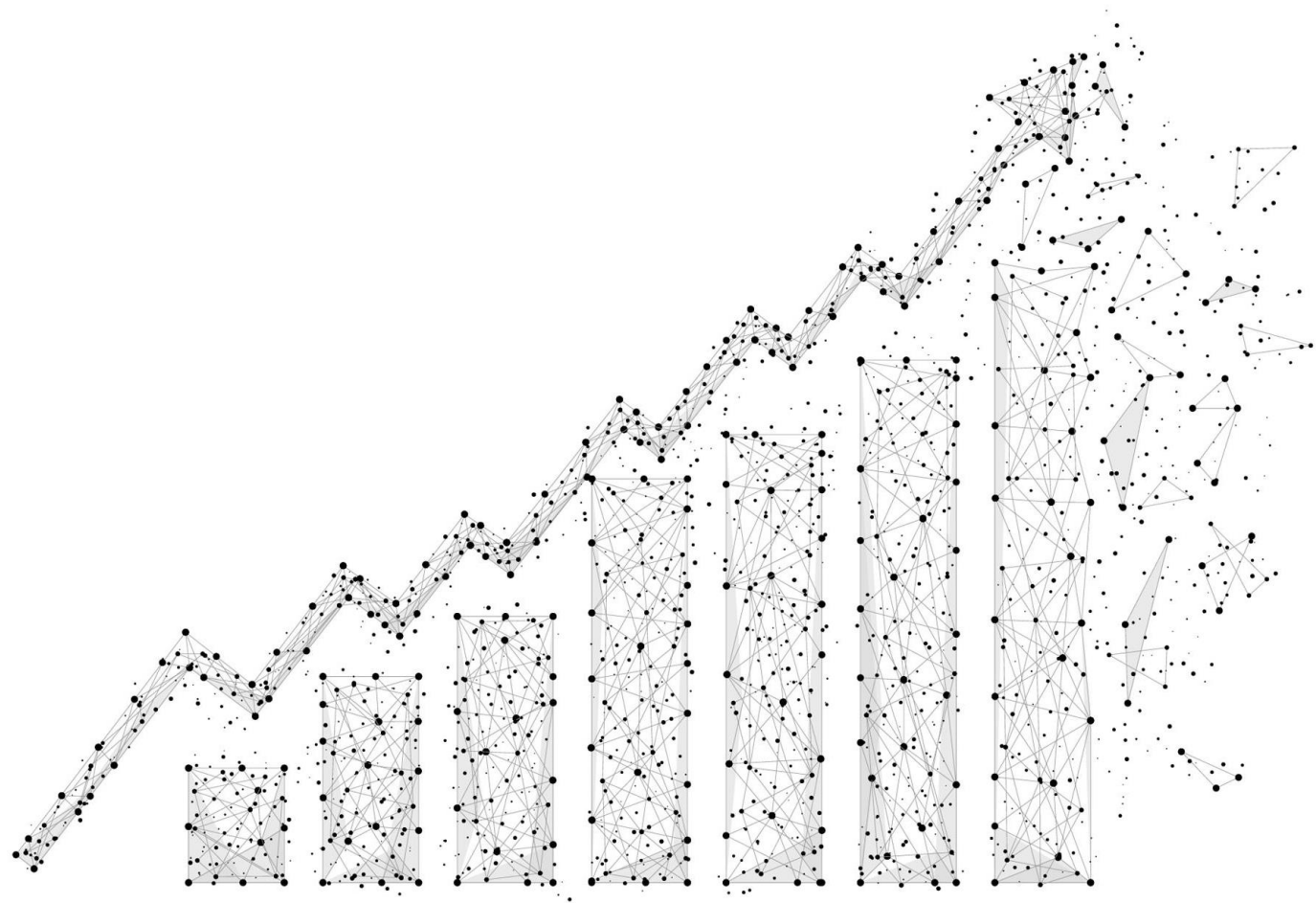
Python is a king of data science

Data sources and sinks

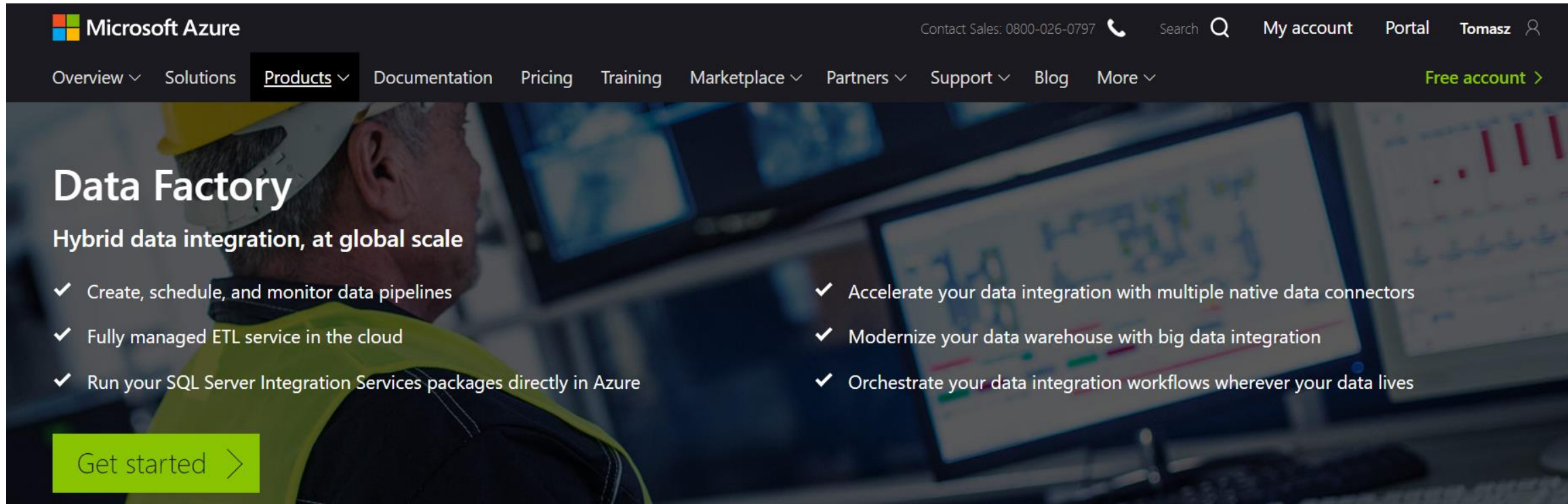
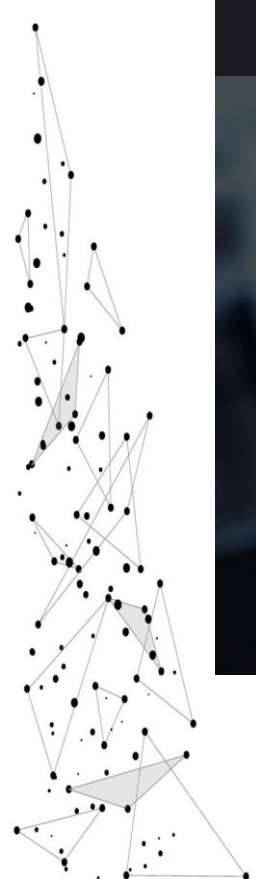


# DEMO

- Demo003Stackoverflow ADLU



# Azure Data Factory (V2)



The screenshot shows the Microsoft Azure Data Factory landing page. The header includes the Microsoft Azure logo, contact information (0800-026-0797), a search bar, and links for 'My account', 'Portal', and 'Tomasz'. The main navigation bar lists 'Overview', 'Solutions', 'Products' (which is highlighted), 'Documentation', 'Pricing', 'Training', 'Marketplace', 'Partners', 'Support', 'Blog', and 'More'. A 'Free account' link is also present. The main content area features a large image of a person in a hard hat looking at a computer screen. The text 'Data Factory' is prominently displayed, followed by the tagline 'Hybrid data integration, at global scale'. Below this, there are six bullet points with checkmarks, arranged in two columns. The first column lists: 'Create, schedule, and monitor data pipelines', 'Fully managed ETL service in the cloud', and 'Run your SQL Server Integration Services packages directly in Azure'. The second column lists: 'Accelerate your data integration with multiple native data connectors', 'Modernize your data warehouse with big data integration', and 'Orchestrate your data integration workflows wherever your data lives'. A green 'Get started' button with a right arrow is located at the bottom left of the main content area.

Microsoft Azure

Contact Sales: 0800-026-0797 Search My account Portal Tomasz

Overview Solutions **Products** Documentation Pricing Training Marketplace Partners Support Blog More

Free account >

## Data Factory

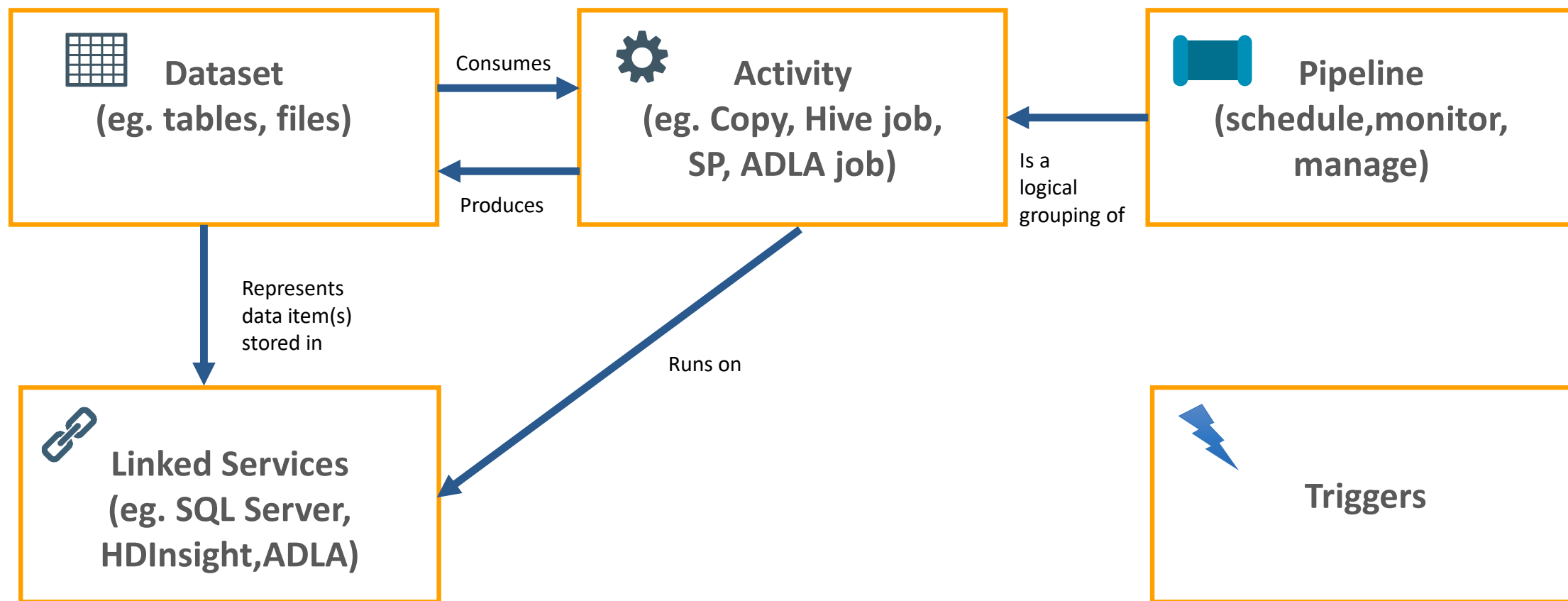
Hybrid data integration, at global scale

- ✓ Create, schedule, and monitor data pipelines
- ✓ Fully managed ETL service in the cloud
- ✓ Run your SQL Server Integration Services packages directly in Azure
- ✓ Accelerate your data integration with multiple native data connectors
- ✓ Modernize your data warehouse with big data integration
- ✓ Orchestrate your data integration workflows wherever your data lives

Get started >

<https://azure.microsoft.com/en-us/services/data-factory/>

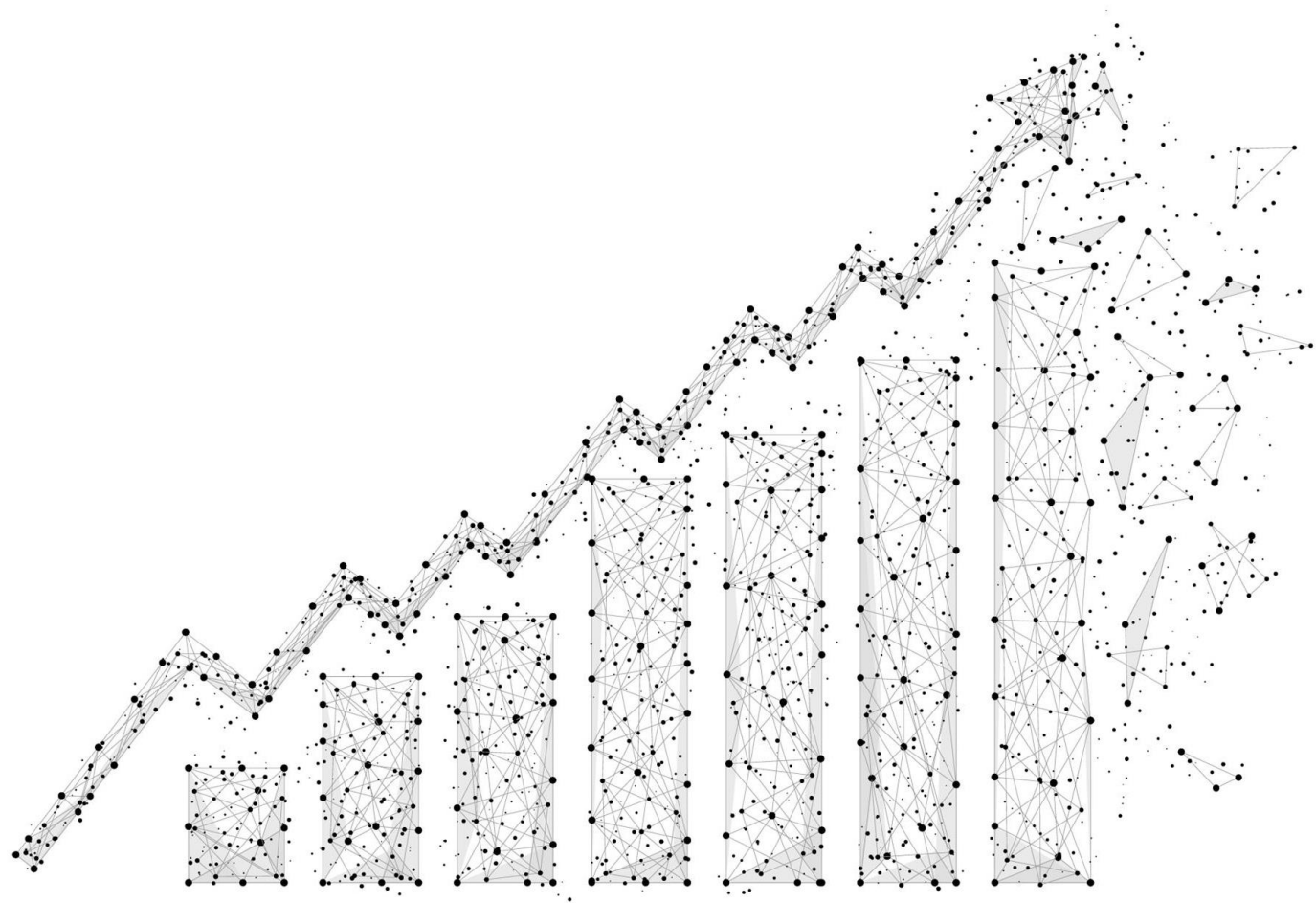
# Azure Data Factory Concept





# DEMO

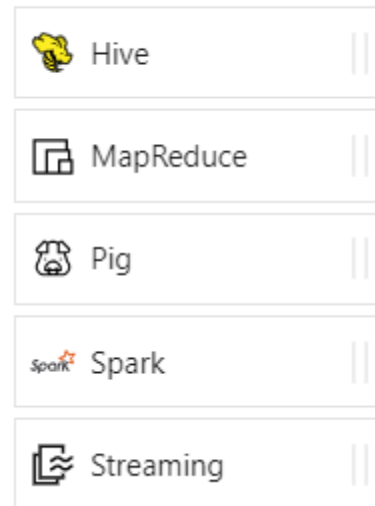
- Demo ADF
- Task8



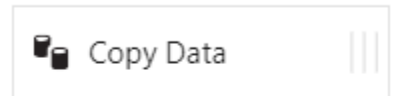
# Azure Data Factory - Data activities

- Data movement activities : **Copy Activity**
- Data transformation activities :
  - HDInsight (Hive, Pig, MapReduce ,Hadoop Streaming ,Spark)
  - **Azure Databricks Notebooks, Jars, Python**
  - Machine Learning activities: Batch Execution and Update Resource (\*)
  - **Stored Procedure**
  - **Data Lake Analytics U-SQL**
  - DotNet (**Custom**)

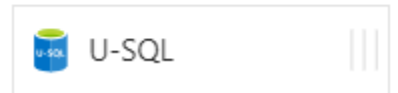
## HDInsight



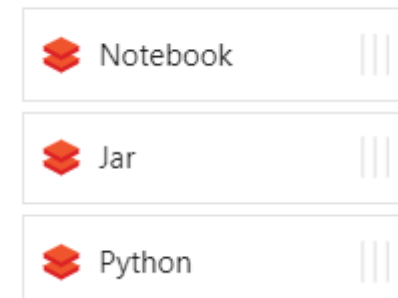
## Move & Transform



## Data Lake Analytics



## Databricks



## Machine Learning





# Resources

## My examples (and demos)

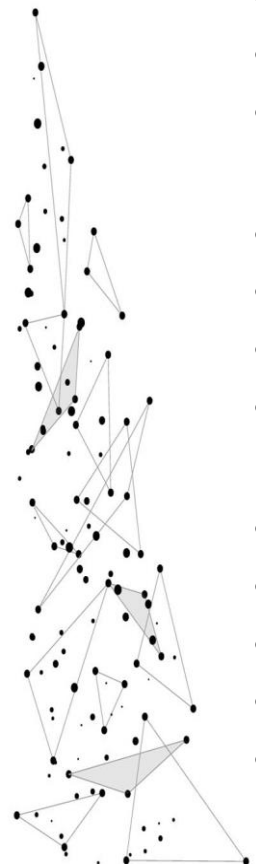
- <https://github.com/FP-DataSolutions/AzureDataLake>
- [https://github.com/cloud4yourdata/CommunityEvents/tree/master/DataCommunity201812\\_AzureDataLake](https://github.com/cloud4yourdata/CommunityEvents/tree/master/DataCommunity201812_AzureDataLake)
- <https://github.com/cloud4yourdata/CommunityEvents/tree/master/4DevKatowice2018/USQL>
- <https://github.com/cloud4yourdata/usql/>

## Blogs and pages

- <http://aka.ms/AzureDataLake>
- <http://usql.io>
- <http://blogs.msdn.microsoft.com/mrys/>
- <http://blogs.msdn.microsoft.com/azuredatalake/>

## Documentation, articles ...

- [http://aka.ms/usql\\_reference](http://aka.ms/usql_reference)
- <https://azure.microsoft.com/enus/documentation/services/data-lake-analytics/>
- <https://msdn.microsoft.com/en-us/magazine/mt614251>
- <https://channel9.msdn.com/Search?term=U-SQL#ch9Search>
- [https://www.youtube.com/results?search\\_query=U-SQL](https://www.youtube.com/results?search_query=U-SQL)



# THANK YOU!

[tkrawczyk@future-processing.com](mailto:tkrawczyk@future-processing.com)