# Azure Data Factory vs SSIS

## May the best tool win!

Regis Baccaro

@regbac          http://theblobfarm.wordpress.com

# About.me : Regis Baccaro

Consultant

Developer

Speaker

Author

Data Platform MVP

Farmer

SQL Nexus lead

MCT

# The contestants

# Agenda

| What is Data Factory | Pricing | | | |
|---|---|---|---|---|

| Core Concepts | JSON | Datasets | Pipeline & Activities | Scheduling & Execution |
|---|---|---|---|---|

| Building a pipeline | Azure Portal | Visual Studio | PowerShell | ARM Template | REST API |
|---|---|---|---|---|---|

| SSIS 101 | What is SSIS for? | Benefits | | | |
|---|---|---|---|---|---|

| Comparing SSIS and ADF | Dev | Admin | Deployment | Monitoring | Source & Destinations | Security |
|---|---|---|---|---|---|---|

Data Community

# What is Data Factory

Data Factory is a fully managed cloud-based data integration service that orchestrates and automates the movement and transformation of data

# Biggest mistake / pain points about Data Factory

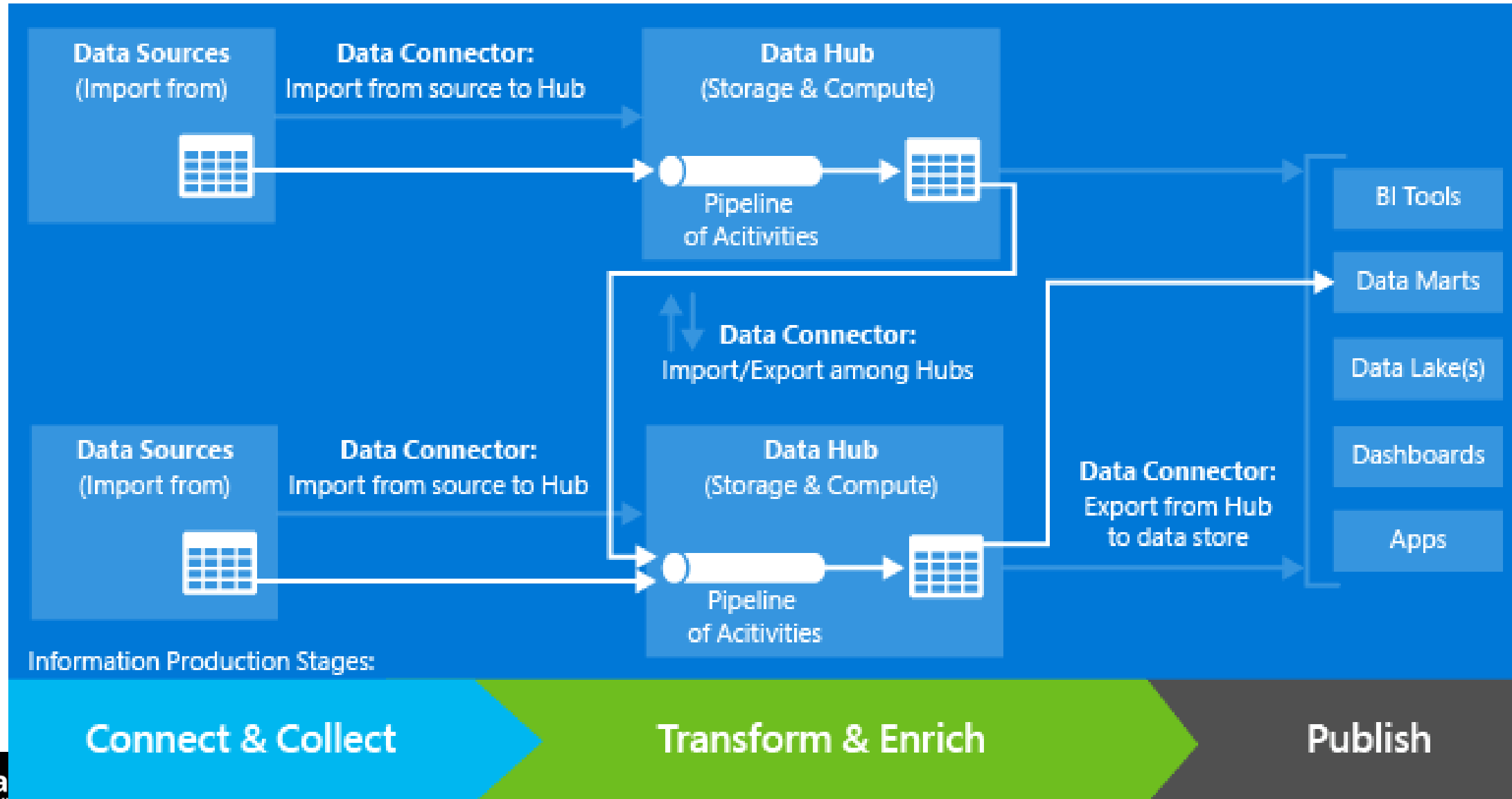Positioning it as SSIS for the Cloud

Adhoc execution

Datatype management

Datasets without time slices

Debugging

Continuous integration

# Azure Data Factory Architecture

No upfront cost

No termination fees

Pay as you g(r)o(w)

Pay for movement and Data usage

A Dataset is a logical description of the data

The mechanism (address, protocol, authentication scheme) to access the data is defined in the Linked Service and referenced in the dataset definition.

Pipeline = logical grouping of activities

Data movement activities

Data transformation activities
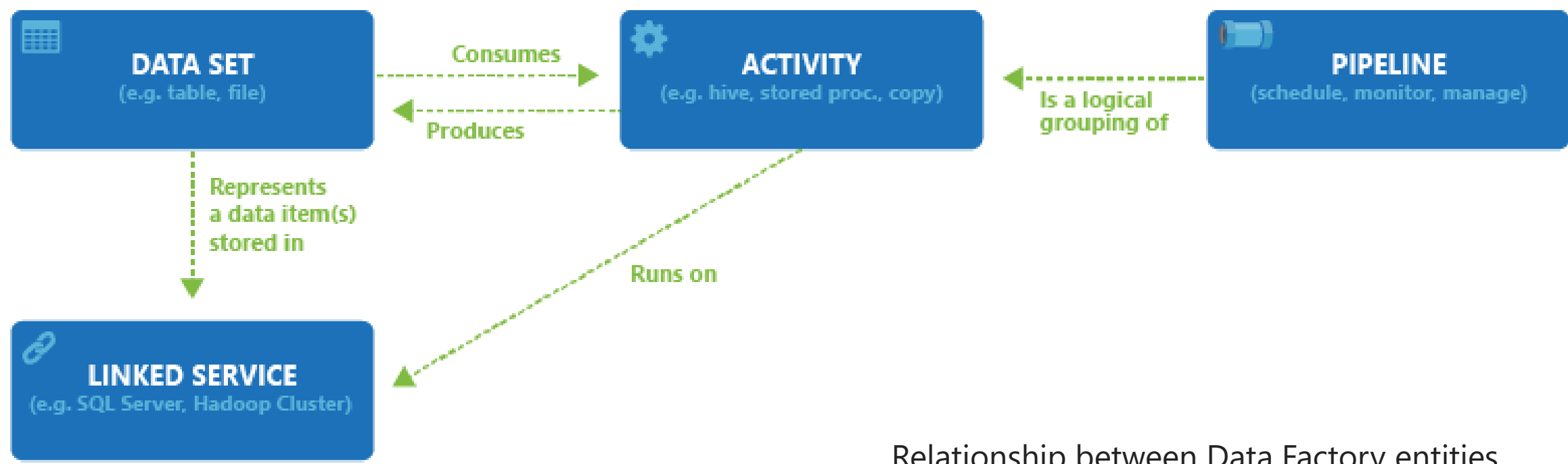


Connect & Collect → Transform & Enrich → Publish

**DATA SET**
(e.g. table, file)

**ACTIVITY**
(e.g. hive, stored proc., copy)

**PIPELINE**
(schedule, monitor, manage)

**LINKED SERVICE**
(e.g. SQL Server, Hadoop Cluster)

Consumes

Produces

Is a logical grouping of

Represents a data item(s) stored in

Runs on

Relationship between Data Factory entities

Core
Concepts

Datasets

Pipeline &
Activities

Scheduling &
Execution

Compute
Linked
Services

2 types of Activities:

Copy & Transform

# Data Movement Activities

## All are Sources / Sink are orange

Azure Blob storage
Azure Data Lake Store
Azure SQL Database
Azure SQL Data Warehouse
Azure Table storage
Azure ~~Document~~ Cosmos DB
Azure Search Index
SQL Server*
Oracle*
MySQL*
DB2*
Teradata*
PostgreSQL*
Sybase*

Cassandra*
MongoDB*
Amazon Redshift
File System*
HDFS*
Amazon S3
FTP
Salesforce
Generic ODBC*
Generic OData
Web Table (table from HTML)
GE Historian*

# Data Transform Activities

| Transformation | |
|---|---|
| **Transformation** | HDInsight [Hadoop] |
| Hive | HDInsight [Hadoop] |
| Pig | HDInsight [Hadoop] |
| MapReduce | HDInsight [Hadoop] |
| Hadoop Streaming | Azure VM |
| Machine Learning activities: Batch Execution and Update | |
| Resource | Azure SQL, Azure SQL Data Warehouse, or SQL Server |
| Stored Procedure | Azure Data Lake Analytics |
| Data Lake Analytics U-SQL | HDInsight [Hadoop] or Azure Batch |
| DotNet | |

# Recurring schedule

# Data slices

Scheduler: Run every hour



| Input Dataset | Activity | Output Dataset |
|---------------|----------|----------------|
| Hourly data **slices** | activity run tumbling windows | Hourly data **slices** |
| 8-9 AM  ready | 8-9 AM | 8-9 AM  ready |
| 9-10 AM  ready | 9-10 AM | 9-10 AM  ready |
| 10-11 AM  ready | 10-11 AM  current | 10-11 AM  current |

**Transformation**

Hive

Pig

MapReduce

Hadoop Streaming

Machine Learning activities: Batch Execution and

Update Resource

Stored Procedure

Data Lake Analytics U-SQL

**Compute Environment**

HDInsight [Hadoop]

HDInsight [Hadoop]

HDInsight [Hadoop]

HDInsight [Hadoop]

Azure VM

Azure SQL, Azure SQL Data Warehouse, or SQL Server

Azure Data Lake Analytics

HDInsight [Hadoop] or Azure Batch

# PORTAL.AZURE.COM

# Linked Service represents

Data Stores = Source or Sink

Compute resource = Data transformation

# Data Stores

Contain credentials and connection information for Sources and Destinations.

An On Premises Data Store **MUST** reference a Data Gateway

# Dataset

Data structure in the data store

**PowerShell**

Get-AzureRmDataFactoryDataset

Get-AzureRmDataFactoryGateway

Get-AzureRmDataFactoryHub

Get-AzureRmDataFactoryLinkedService

Get-AzureRmDataFactoryPipeline

Get-AzureRmDataFactoryRun

Get-AzureRmDataFactorySlice

**Azure Resource Manager Template**

More JSON !!

**REST API**

CURL Tool with REST

Create Web Application in AAD

Assign to Data Factory Contributor Role

Use CURL to communicate with Web Application

```
$cmd = {.\curl.exe -X PUT -H "Authorization: Bearer $accessToken" -H "Content-Type: application/json" --data "@azurestoragelinkedservice.json" https://management.azure..

$results = Invoke-Command -scriptblock $cmd;
```

Data Extract, Transformation and Loading tool

Born in 2005 – enterprise ready but still room for new features

Part of SQL Server license and installation

Rich development tool

Many built-in transformations

Extensible with scripts

| Development tool | Administration tool | Data source & destinations | Data transformations |
| Price | Error handling | Deployment | Monitoring |
| Security | Technology | Requirements | Big Data compatibility |

| SSIS | ADF |
|------|-----|
| SSDT | Azure portal ADF Editor |
| Free | PowerShell |
| | JSON Scripts |
| | Visual Studio |

| | SSIS | ADF |
|---|------|-----|
| Standalone tool | Yes | Yes |
| Powerful GUI | Yes | No |
| Available | Free | Free |
| Prerequisite | SQL Server | Azure Subscription |

| SSIS | ADF |
| --- | --- |
| SSMS | Azure portal |
| PowerShell | PowerShell w/ADF Cmdlets |

Data movement

- Between Cloud data stores: €0,21/hour

- When on-prem is involved: €0,08/hour

Inactive Pipelines:€0,67/month

Re-running activities

- In the Cloud: €1,13 per 1000 re-runs

- On-prem: €2,83 per 1000 re-runs

Inactive Pipelines:€0,67/month

Calculator

https://azure.microsoft.com/en-us/pricing/calculator/?service=data-factory

Basic SSIS

    Free = Express edition – Import & Export wizard

Standard SSIS

    Standard (and BI editions)

Enterprise SSIS (CDC & Advanced adapters)

    Enterprise edition

Still true with SQL Server 2016 SP1

|  | ADF | SSIS |
|---|---|---|
| Licensing |  | Yes |
| Pay for features |  | Yes |
| Pay per usage | Yes |  |

|  | ADF | SSIS |
|---|---|---|
| Azure environment | Yes |  |
| Hardware setup |  | Yes |
| Software setup |  | Yes |
| Administration costs / Data center |  | Yes |

# SSIS

    MSDB

    Project deployment

    Package deployment

    SSIS Catalog

# ADF

    Power shell scripts

    Automization w/ PS Scripts

|  | ADF | SSIS |
|---|---|---|
| Alerts |  | Yes |
| Error Loging | Yes | Yes |
| Error handling |  | Yes |



| Data slices (by update time) SQLAzureTblShipData | | | |
|---|---|---|---|
| ▽ Filter | | | |
| LAST UPDATE TIME | SLICE START TIME | SLICE END TIME | STATUS |
| 10/06/2015, 9:01:14 ... | 10/06/2015, 8:00 P... | 10/06/2015, 9:00 PM... | ✖ Failed |
| 10/06/2015, 8:01:43 ... | 10/06/2015, 7:00 P... | 10/06/2015, 8:00 PM... | ✖ Failed |
| 10/06/2015, 7:01:32 ... | 10/06/2015, 6:00 P... | 10/06/2015, 7:00 PM... | ✖ Failed |
| 10/06/2015, 6:01:18 ... | 10/06/2015, 5:00 P... | 10/06/2015, 6:00 PM... | ✖ Failed |

SSIS Logging

SSIS Catalog reports

Diagram view

Drill through features

Capable GUI

Data slice execution

Data lineage

|  | ADF | SSIS |
|---|---|---|
| Monitoring GUI | Yes | Yes |
| Drillthrough | Yes | Yes |
| Data slice | Yes | |
| Data lineage | Yes | |

SQL Server
Oracle
SAP
Azure
Access
Sybase
PostGresSQL
FoxPro
SharePoint
WebService
…..

▲ Other Sources
⮕ ADO NET Source
CDC Source
Excel Source
Flat File Source
ODBC Source
OLE DB Source
Raw File Source
XML Source

▲ Other Destinations
ADO NET Destination
Data Mining Model Training
DataReader Destination
Dimension Processing
Excel Destination
Flat File Destination
ODBC Destination
OLE DB Destination
Partition Processing
Raw File Destination
Recordset Destination
SQL Server Compact Destination
SQL Server Destination

# All are Sources / Sink are orange

Azure Blob storage
Azure Data Lake Store
Azure SQL Database
Azure SQL Data Warehouse
Azure Table storage
Azure DocumentDB
[Azure Search Index]
SQL Server*
Oracle*
MySQL*
DB2*
Teradata*
PostgreSQL*

Sybase*
Cassandra*
MongoDB*
Amazon Redshift
File System*
HDFS*
Amazon S3
FTP
Salesforce
Generic ODBC*
Generic OData
Web Table (table from HTML)
GE Historian*

|  | ADF | SSIS |
|---|---|---|
| Copy | Yes | Yes |
| C# custom transformations | Yes | Yes |
| Pig and Hive | Yes | Yes* |
| Azure ML Scoring | Yes | With scripting |
| Stored procedure | Yes | Yes |
| Built-in transformations |  | Yes |

| | ADF | SSIS |
|---|---|---|
| Role based | Yes | Yes |

Package Execution Progress

UploadDefects
  Validation ha
  DFT_tblDe
  Validatio
  Start,
  Valid
  Vali
  Sta
  Fin
  Fin

Elapsed time: 00:07:01.656
Elapsed time: 00:07:01.698

Duration

00:07:34

Input Dat

...f-440e-bf0c-3e49c44eb756

Input Datasets 0
tblDefectOnPrem

Output

Package Execution Progress

UploadDefects
  Validation has started
  DFT_tblDefects
  Validation is completed
  Start, 09:52:26
  Validation has star
  Validation is comp
  Start, 09:52:27
  Finished, 09:58:
  Finished, 09:58:

ime: 00:06:08.062
me: 00:06:08.104

Data Factory Name

uration

00:05:03

...85f-440e-bf0c-3e49c44eb756

Input Datasets 0
tblDefectOnPrem

## Package Execution Progress

- UploadDefects
  - Validation has started
  - DFT_tblDefects
  - Validation is completed
  - Start, 10:58:19
  - Validation has started
  - Validation is completed
  - Start, 10:58:19
  - Finished, 11:07:27, E
  - Finished, 11:07:27, E

e: 00:09:08.516

e: 00:09:08.554

uration

0:09:35

85f-440e-bf0c-3e

# Conclusion

Not build for the same purpose

ADF still a V1 product

ADF great for Cloud Data integration for MS and Azure

Get benefits of both

    Hybrid SSIS and ADF

    Cloud based data movement, computing and monitoring

    On premises Data Transformations

# Questions?

Contact me : regis@Baccaro.com

Twitter : @regbac