



Azure Data Lake & Factory

Technical Deep Dive –Day 3

Sashank Pappu
Azure Data & AI Specialist

Shiva Priya
Azure Data Engineer



Sashank Pappu

Azure Data & AI Specialist.

Certified from MIT on 'Advanced Analytics and Big Data Challenges'. Technology Enthusiast and Data Specialist. Certified Designer in Computer Vision

Certifications

- Microsoft Certified Power BI Developer.
- Microsoft Certified Azure Machine Learning Developer from Microsoft Virtual Academy
- Certified on Practical Data Analytics with Microsoft Cortana Intelligence Suite.

Role

Currently working as Independent Consultant on Azure Data & AI and Enterprise Trainer on Azure Data Analytics , Cortana Intelligence Suite and Cognitive Sciences



Experience

9+ years of total experience as Data Specialist and Data Engineer.

Training

Expert **facilitator and corporate behavioral trainer** with 5 years of experience in designing programs, content development , Data Analytics and Machine Learning.

- You can see me on – [Sashankpappu.com](https://sashankpappu.com)

What we can learn together ??

- Azure Data Ecosystem
- Storage (Blob Vs Data Lake)
- Data Lake Storage
- Data Lake Analytics
- Data Factory

URL : <https://github.com/aagasp/azuredataecosystem>

Questionnaire



To be taken after each section .



Request you to ask the questions in Chat window .



Will respond to them after a logical break point in each topic.



This would make the workshop interesting and approachable .



Request not to use Phones or get distracted for next 3 hours .

Expected Outcome!



Azure Data
Ecosystem – When to
use what !



Architectural
Standards & Best
Practices in Azure .



Hands on Project
ready confidence.

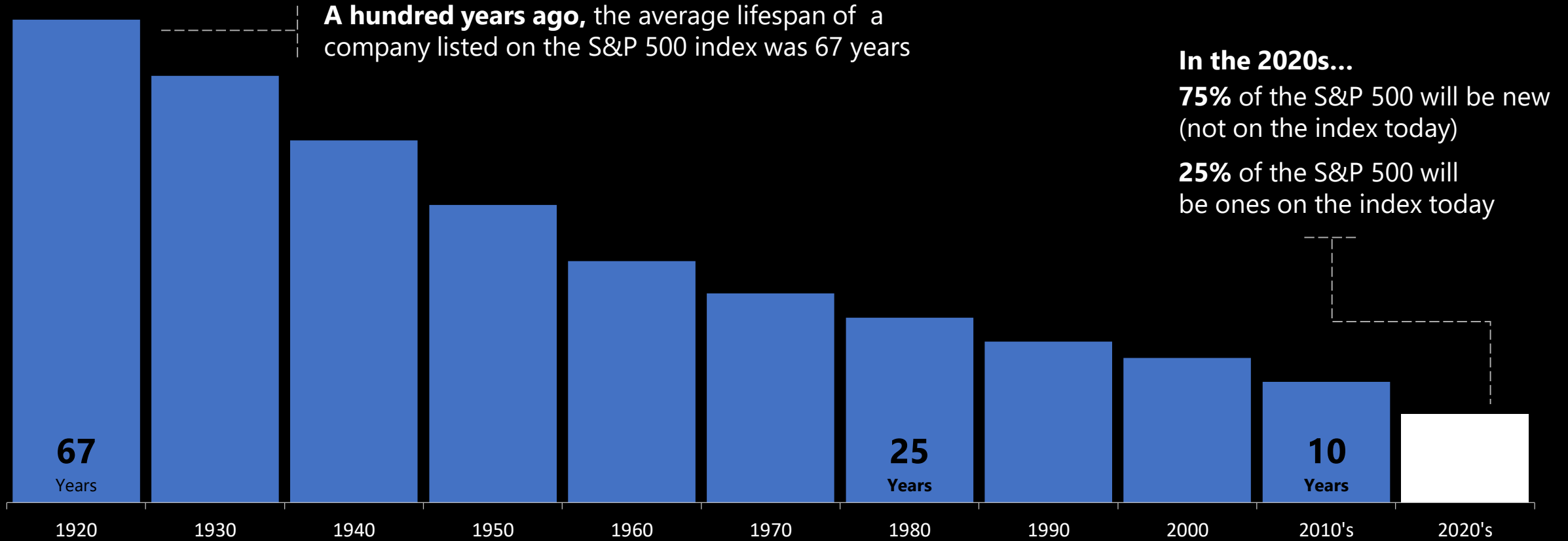


Tackling Real world
challenges in Azure
Data Platform



Real world use cases
and architectures for
huge clients.

Time to adapt is shrinking



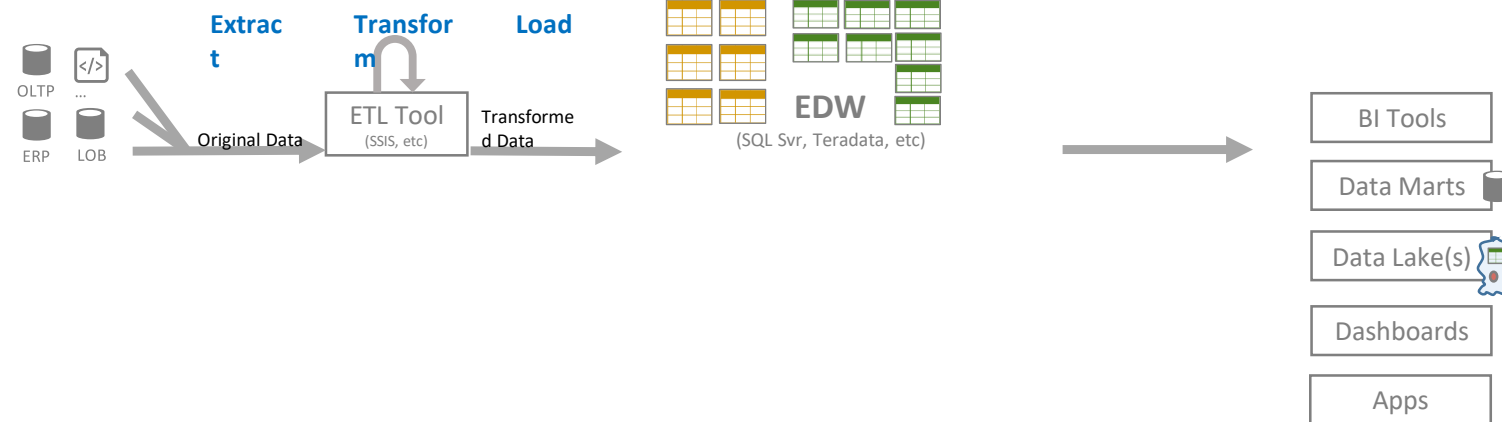
Source: BBC

Experts predict that in the next decade, **only 25%** of the companies currently listed on the S&P 500 will remain there, and the other **75% will be replaced** by new companies.

Approach of Data is
Changing !!!

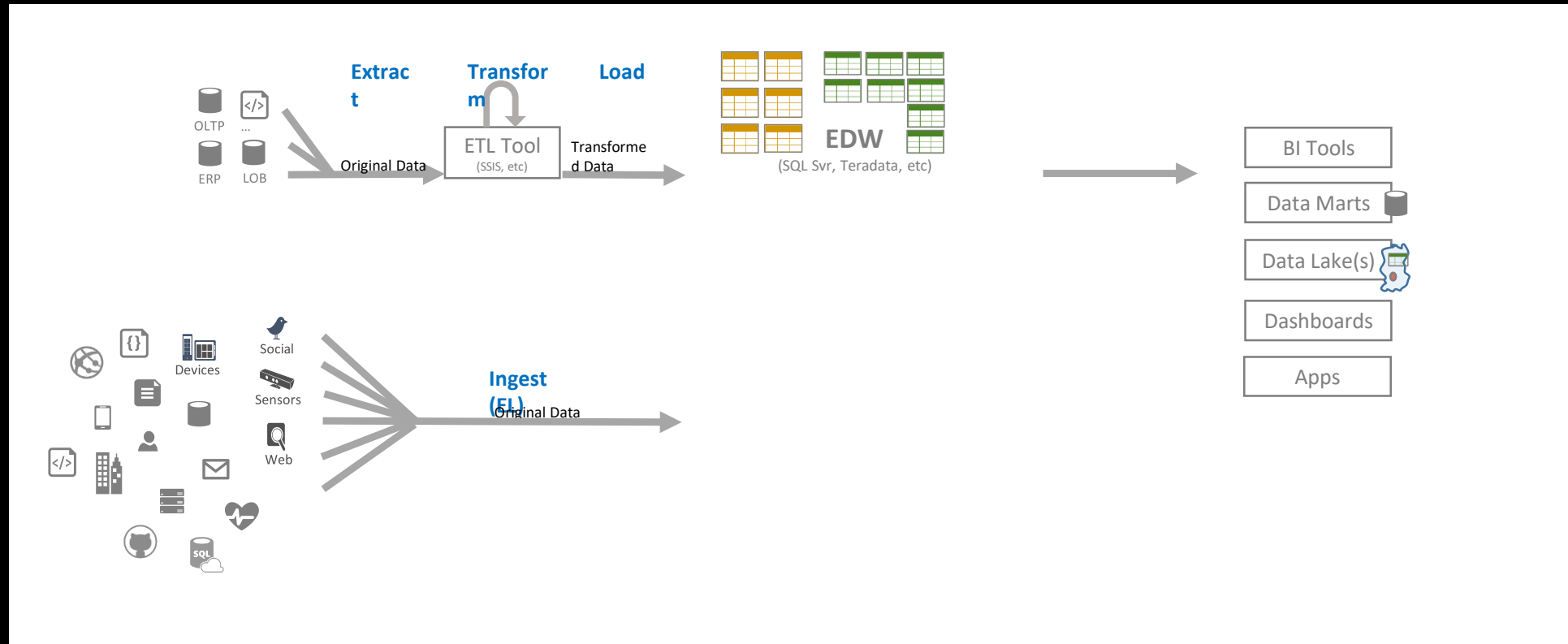
Evolving Approaches to Analytics

EDWs to Data Lakes



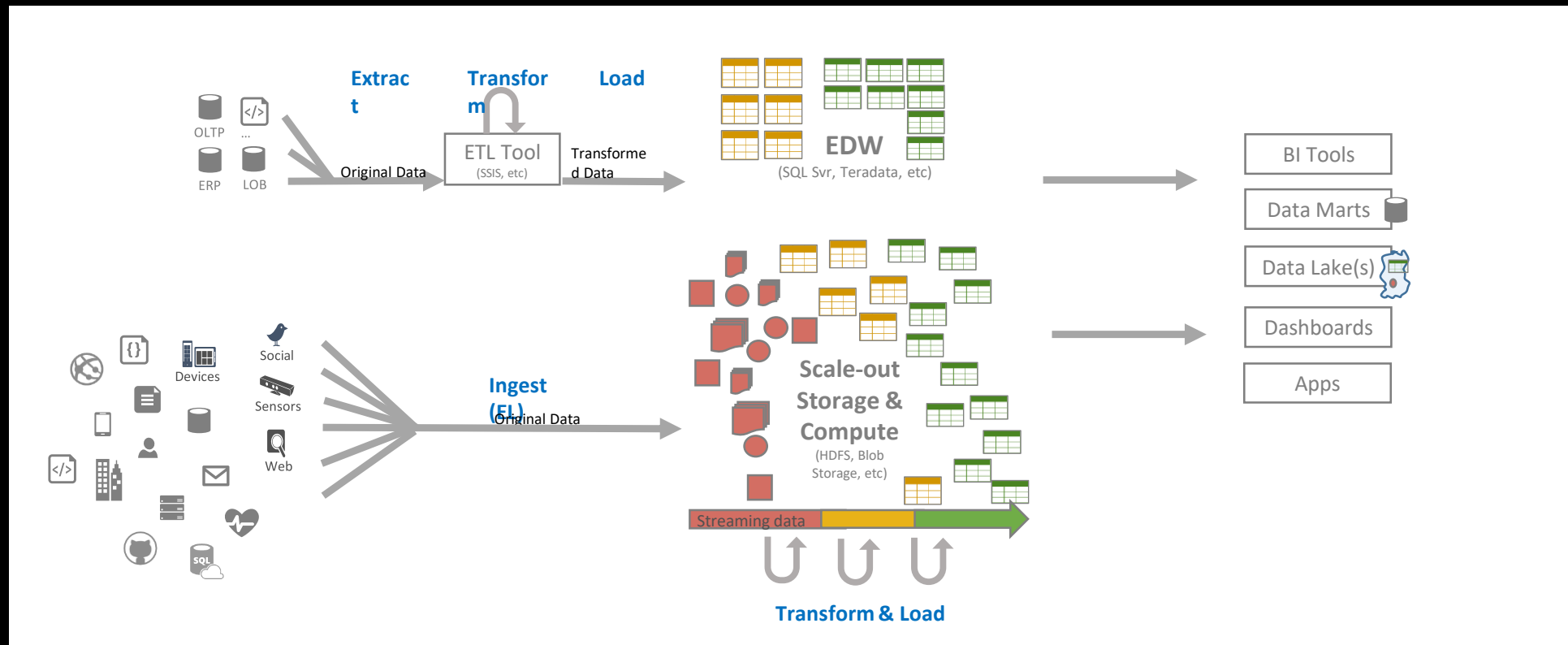
Evolving Approaches to Analytics

EDWs to Data Lakes



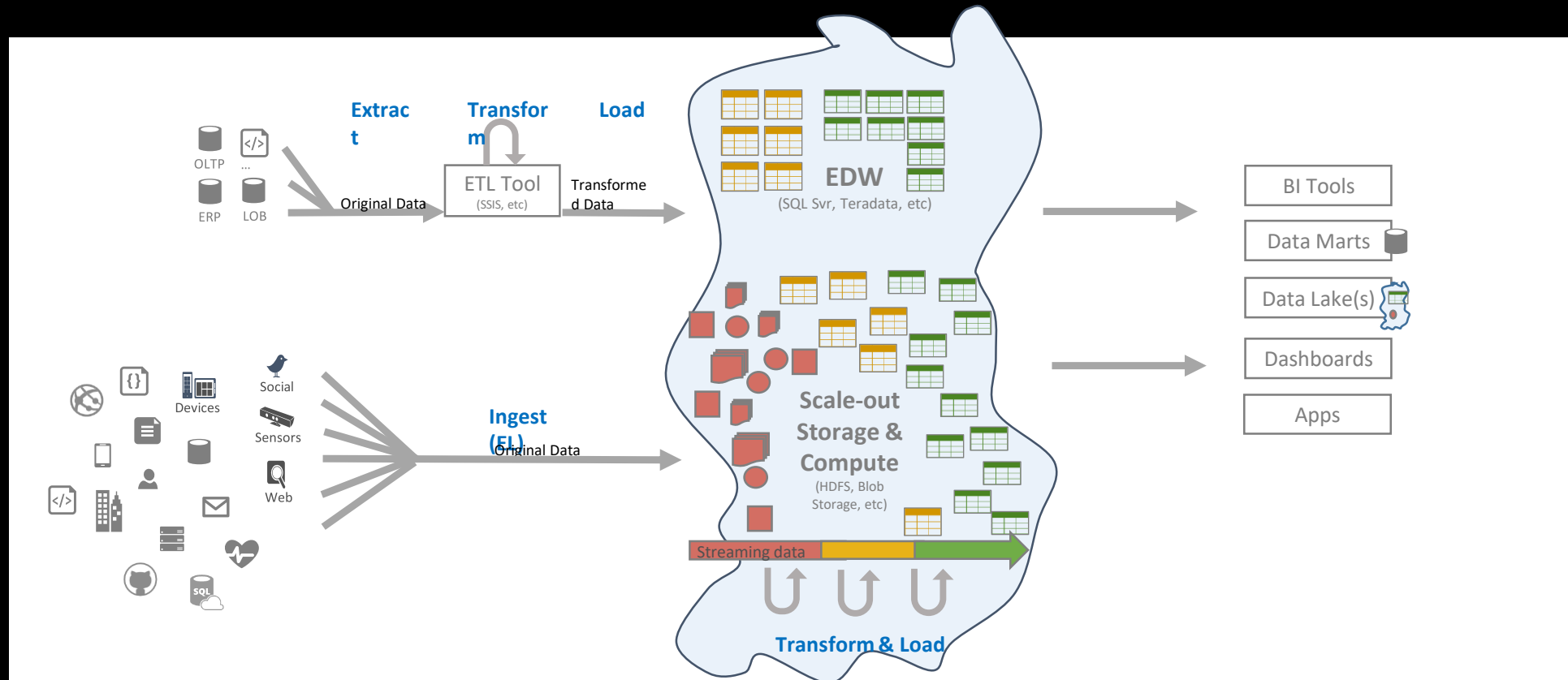
Evolving Approaches to Analytics

EDWs to Data Lakes



Evolving Approaches to Analytics

EDWs to Data Lakes



Data Lake - Customer Challenges

Data Silos



Data spans sources
Inefficiency in
colocation

Analytics



Open interfaces to
data Variety of
analytics tools

Perf & Scale



Storage bottlenecks
IoT sources – Small
writes
Price-performance
Data grows
independently

Security



Compliance challenges
Effectively control access
Corporate policies

The Intelligent Lake

Petabytes of Data + Intelligence



Data Lake Store

No-limits data Lake

- Petabyte-sized files
- Trillions of objects
- Scalable I/O for parallel analytics
- Enterprise-grade Security



Data Lake Analytics

Analytics job service

- Scale instantly
- Scale per job
- Infused with Artificial Intelligence

Azure Data Lake Storage Gen 1 Vs BLOB

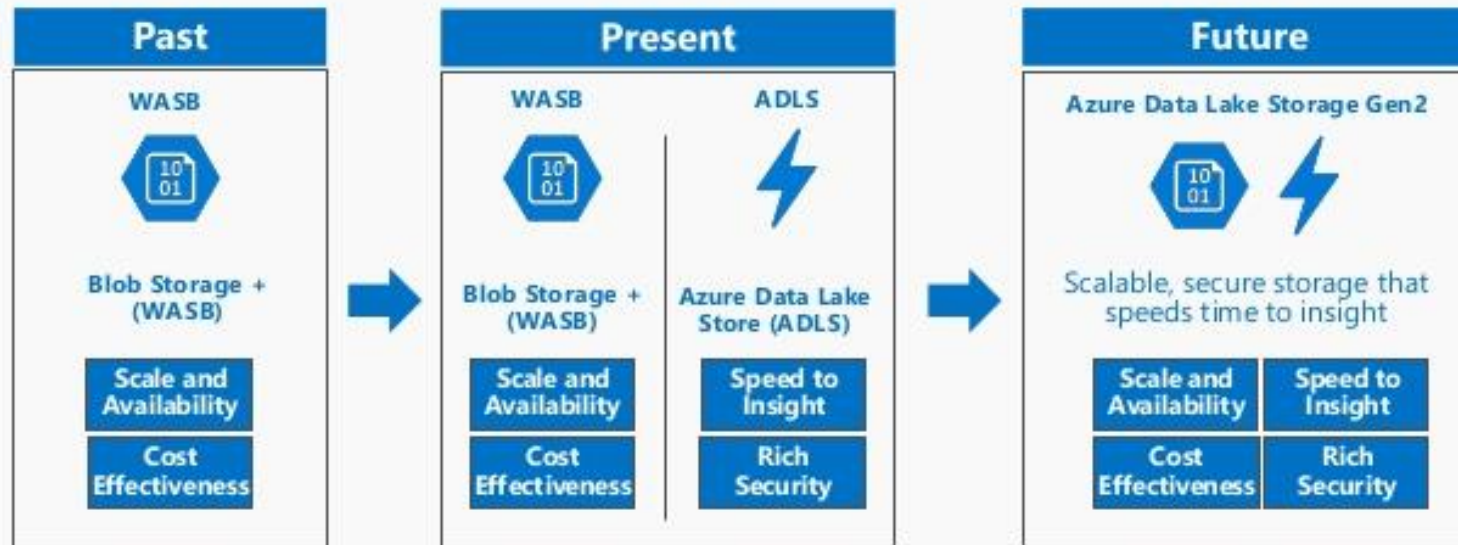
Azure Blob Storage	Azure Data Lake
General purpose storage for storing applications and user data	Specialized storage for storing Big Data which is optimized for Analytics
Object based flat system	Hierarchical system with folders
Security based on shared secret keys and signed URLs	Security based on Azure AD authentication
LRS, ZRS, GRS and RA-GRS	LRS, ZRS, GRS and RA-GRS
SDKs available for .NET, Java, Python, Node.js, C++, Ruby	.NET, Java, Node.js

ADLS Gen1 + BLOB = ADLS Gen 2

Use	Data Lake Store	Blob Storage	Data Lake Gen 2
Hot/Cold Storage Tiers	No	Yes	Yes
Redundant Storage	No	Yes	Yes
AD Security	Yes	No	Yes
HDFS Compatible	Yes	No	Yes

ADLS Gen1 + BLOB = ADLS Gen2

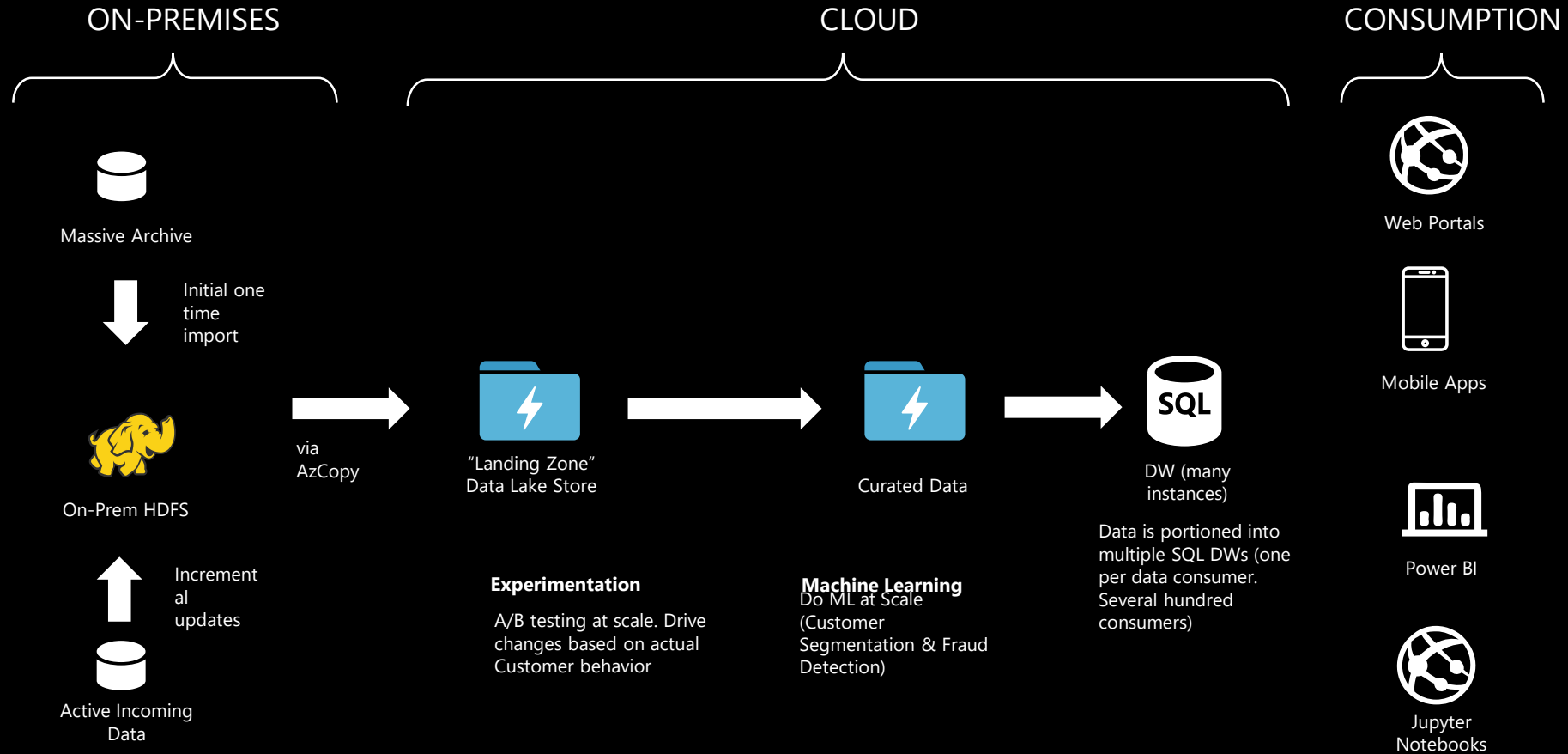
Evolving Data Lake Strategy



Azure Data Lake Storage Gen2: Single Data Lake Store that combines the performance and innovation of ADLS with the scale and rich feature set of Blob Storage

Azure Data Lake Use Cases

Retail Scenario Implementation



TYPES OF DATA



Structured Data



Un-structured Data



Images



Emails

DATA LAKE ZONE



Acquire/ Ingest



Apply Metadata,
Protect sensitive



Data Quality and
Validation



Data Quality and
Validation

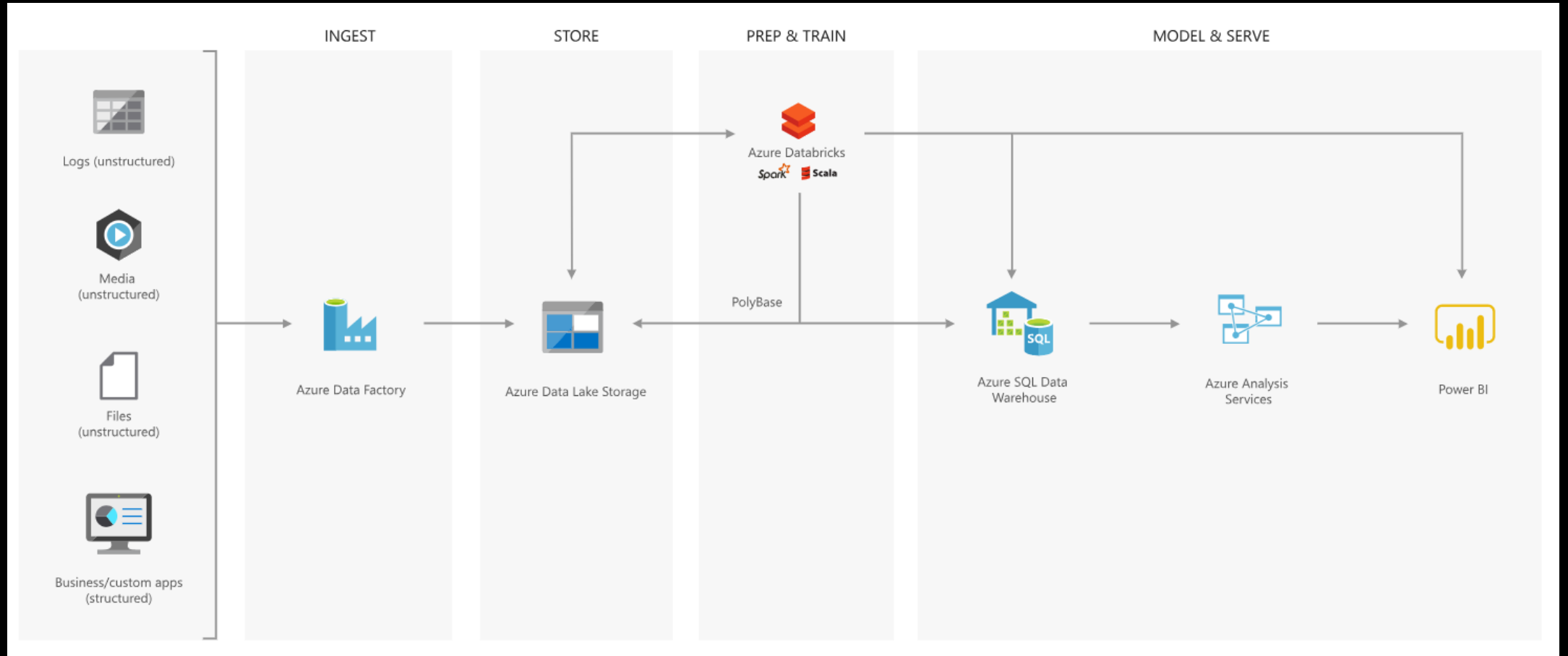
CONSUMER OUTPUT



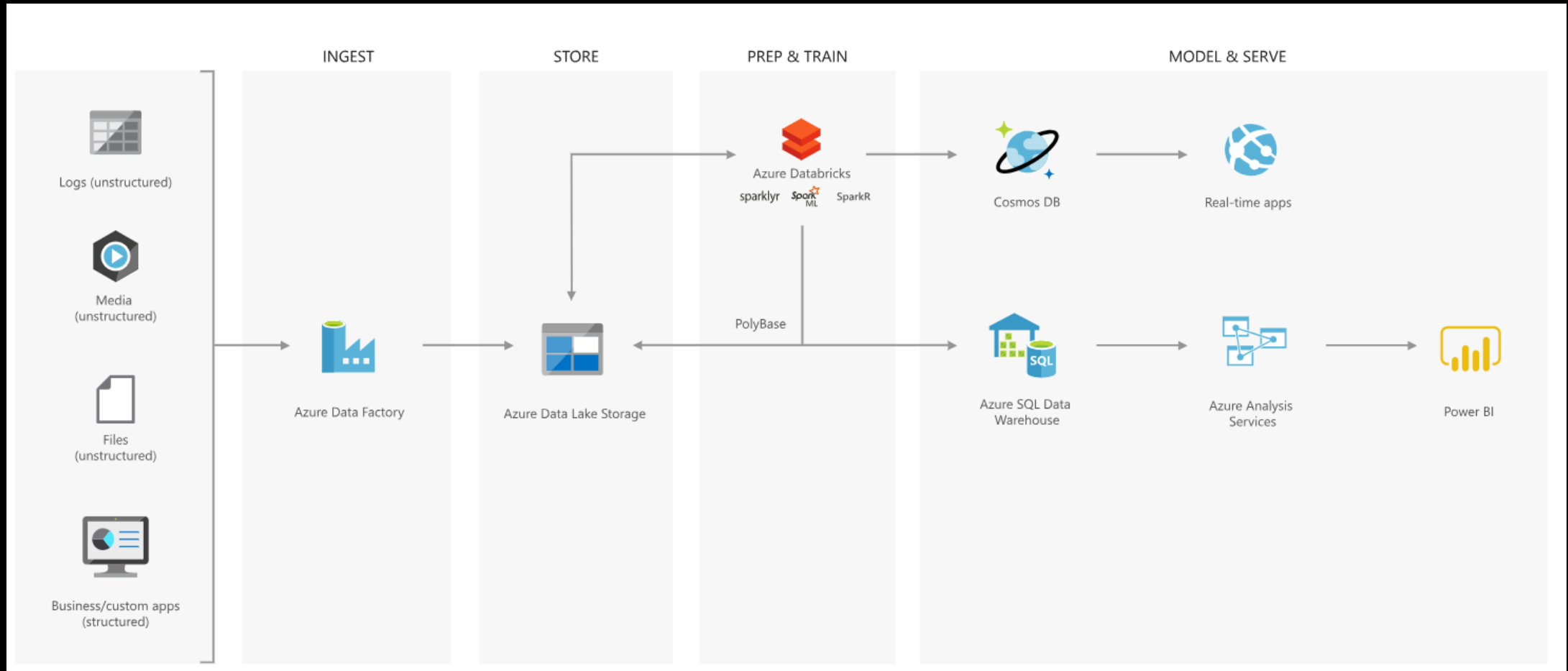
Data as and when
required in formats
suitable to consumers

Data Visualization
Data Prep Tools

Best Practices – Modern Data Warehouse



Advanced Analytics Approach

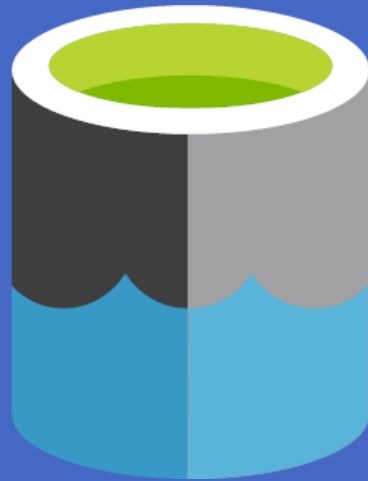


Azure Data Lake Store



Introducing Azure Data Lake Store

A hyper-scale repository for big data analytics workloads



Store **ANY DATA** in its native format

HADOOP FILE SYSTEM (HDFS) for the cloud

ENTERPRISE GRADE

No limits to **SCALE**

Optimized for analytic workload
PERFORMANCE

The Data Lake Approach

Ingest all data
regardless of
requirements

Store all data
in native format
without schema
definition

Do analysis
Using analytic
engines like Hadoop
and ADLA



Any data

Unstructured
Semi-structured
Structured

Devices



Video



Clickstream



Web



Social



Sensors



Relational



LOB
applications

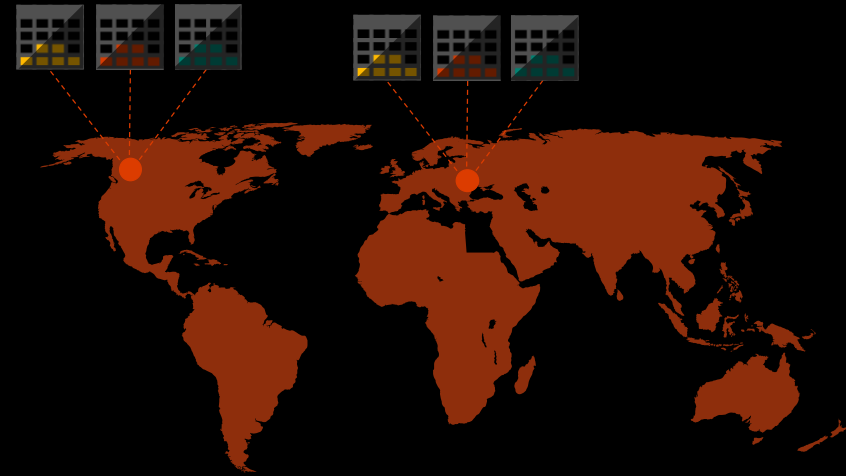


Durable and highly available

Automatically replicates your data

Three copies within a single region

Highly available

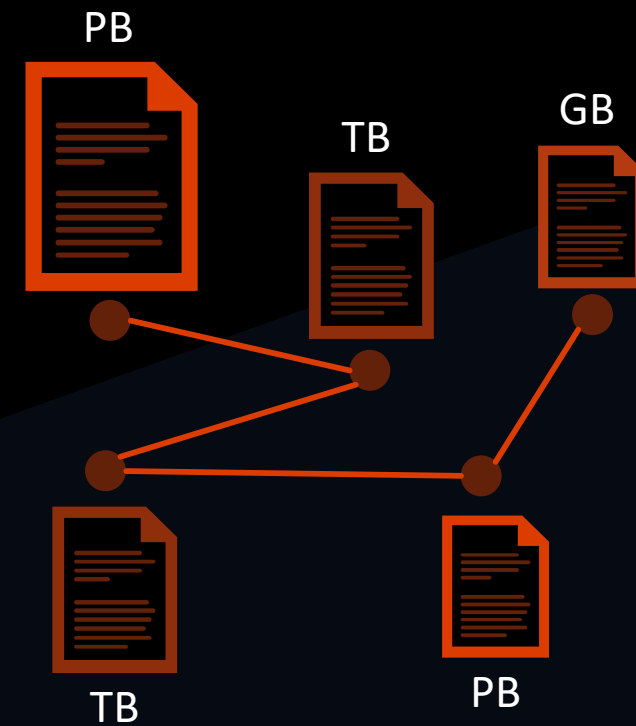


Unlimited storage

Unlimited account sizes

Individual file sizes from gigabytes to petabytes

No limits to scale



Hands on – Data Lake Gen2

Creation of Data Lake Gen2 Account

Security Considerations

Data Protection



Protect the data at rest



Protect the data in transit



Support browser cross-domain access



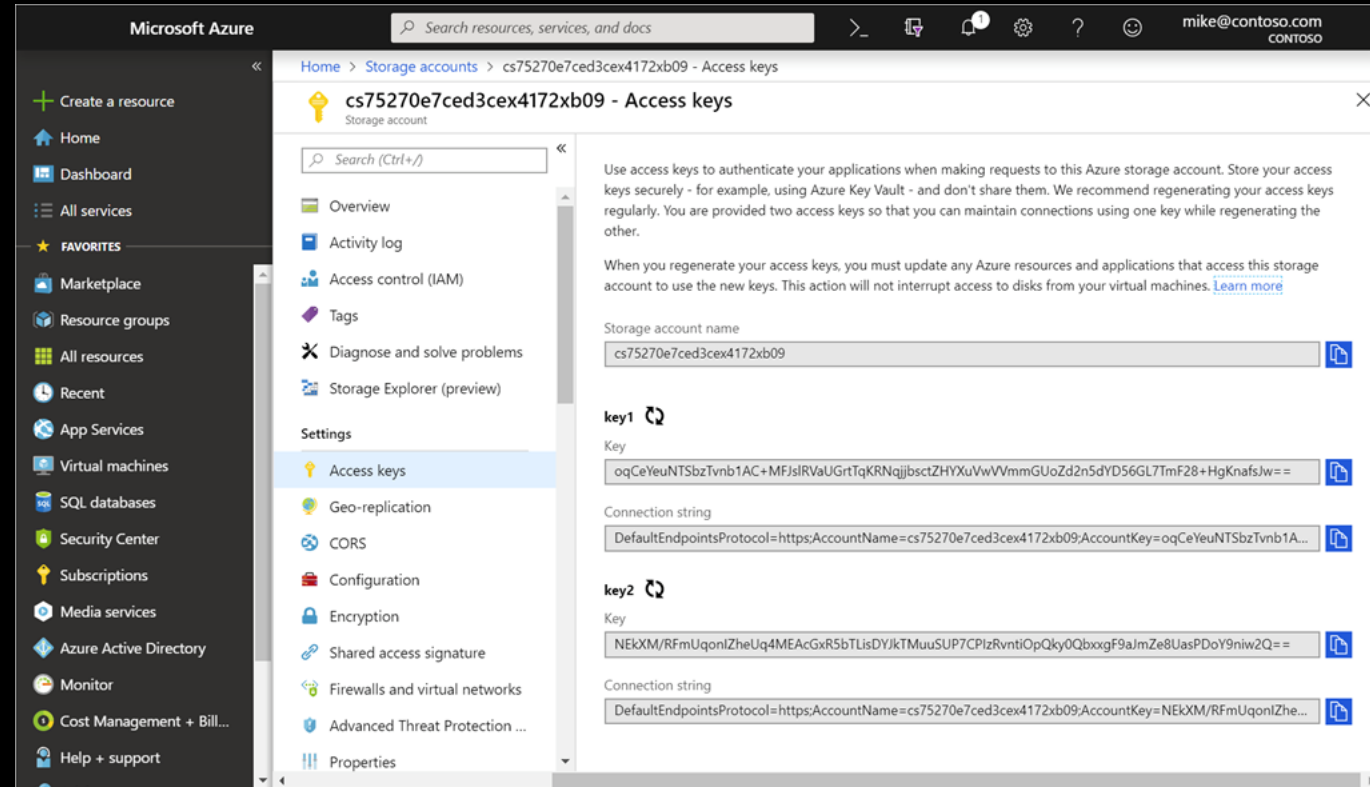
Control who can access data



Audit storage access

Understand storage account keys

- Data is generated or consumed by custom applications
- The applications are written in various languages.
- Create authorized apps in Active Directory to control access to the data in blobs

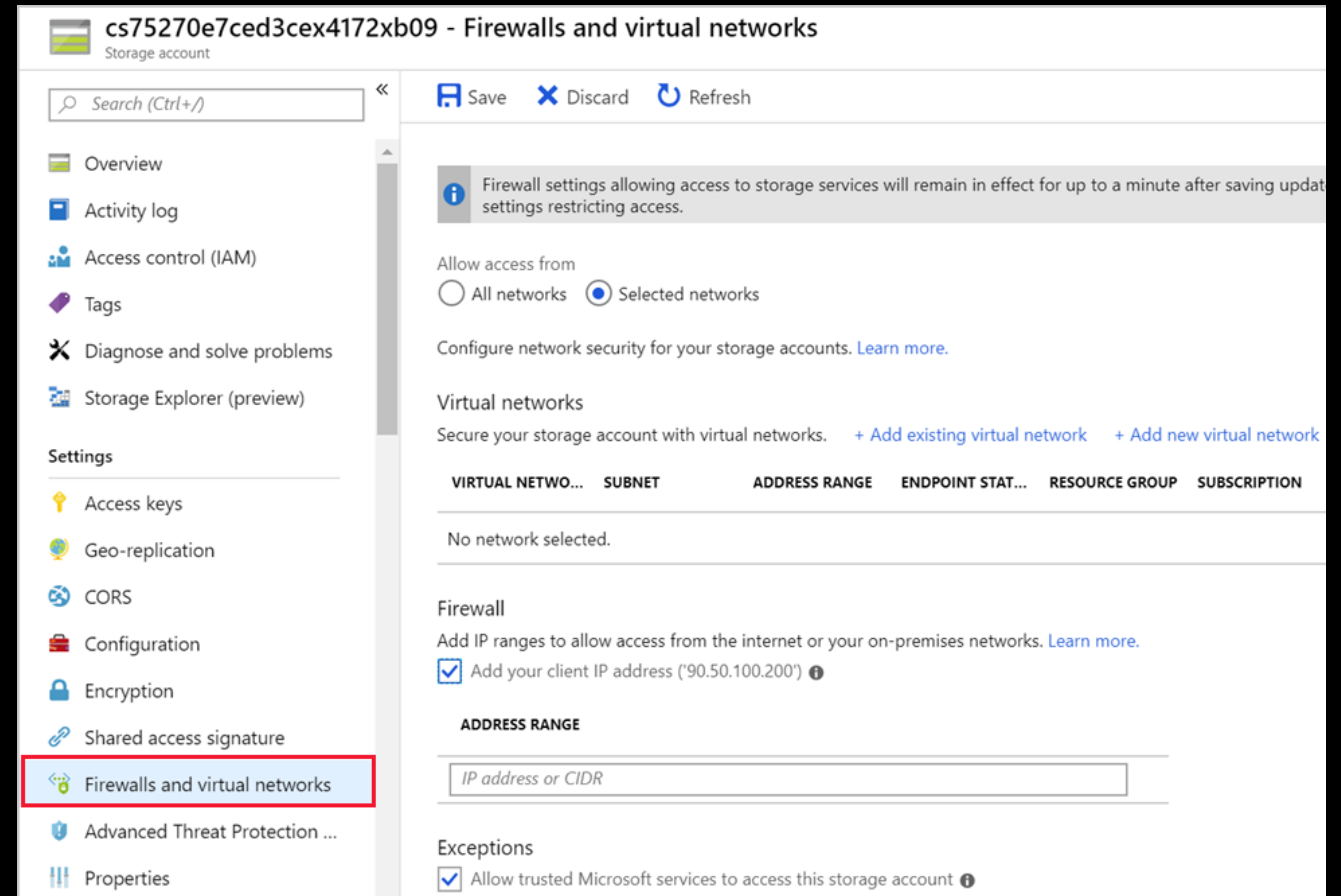


Understand shared access signatures

- You shouldn't share storage account keys with external third-party applications
- For untrusted clients, use a *shared access signature* (SAS)
- **Types of shared access signatures**
 - *service-level* shared access signature
 - *account-level* shared access signature

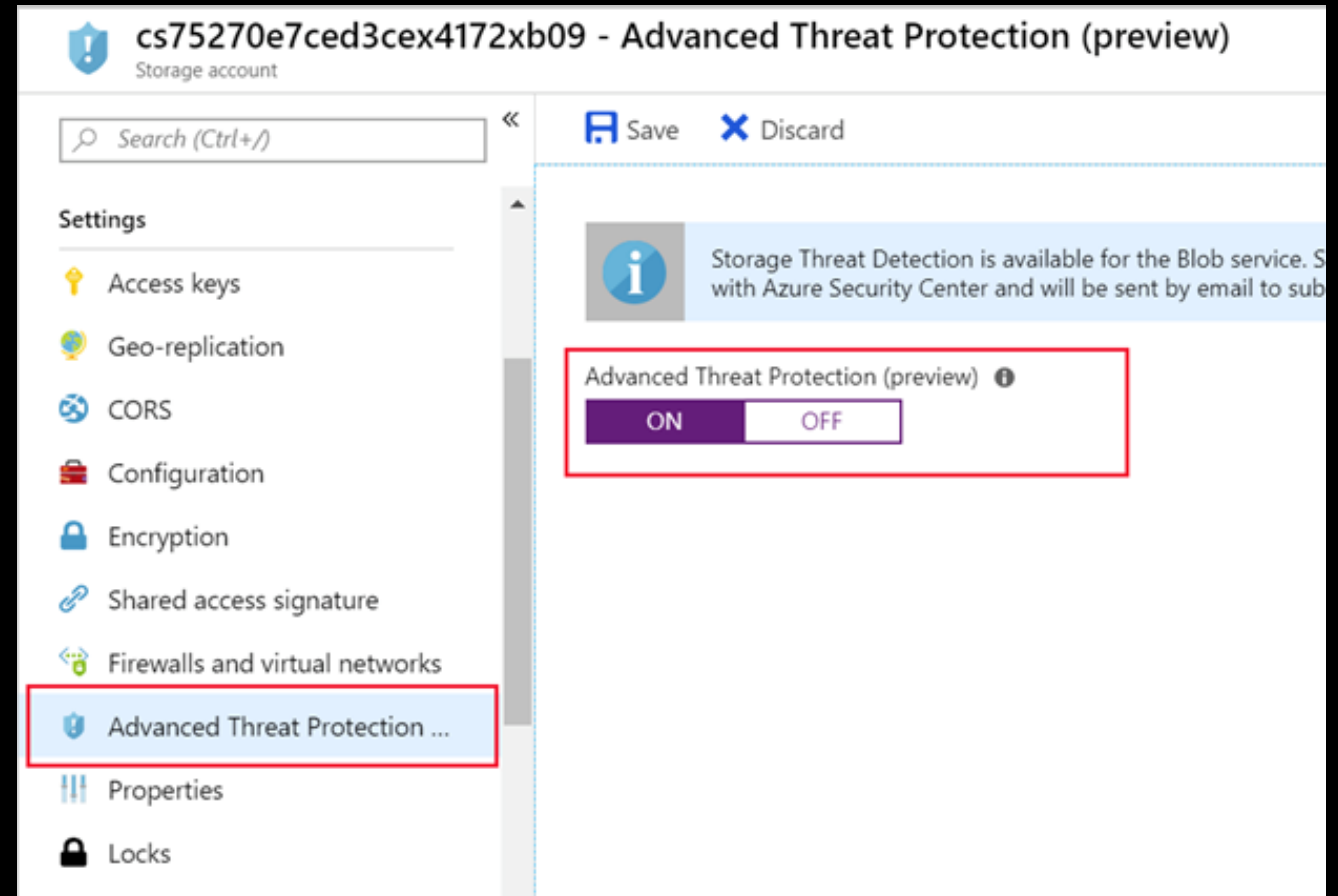
Control network access to your storage account

- Storage accounts accept connections from clients on any network
- To limit access to selected networks, you must first change the default action.
- You can restrict access to specific IP addresses, ranges, or virtual networks



Understand Advanced Threat Protection for Azure Storage

- Often any protection only shows you that an intrusion *has already occurred*.
- What you really want is a way to be notified when suspicious activity is happening
- Advanced Threat Protection, detects anomalies in account activity
- It then notifies you of potentially harmful attempts to access your account.



Azure Data Lake Analytics

Quick Recap

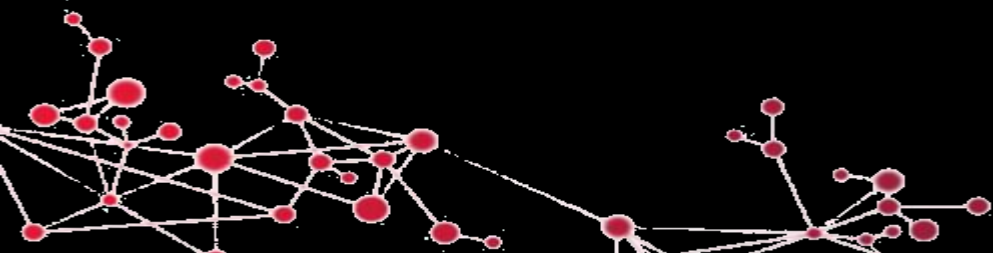
ADL VS SQL DB

Azure Data Lake

ADL Analytics + ADL Store

SQL DB

SQL Query + DB Storage



U-SQL

A 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

History

History

Bing needed to...

- Understand user behavior

And do it...

- At massive scale
- With agility and speed
- At low cost

So they built ...

- Cosmos

Cosmos

- Batch Jobs
- Interactive
- Machine Learning
- Streaming

Thousands of Developers

Pricing

REGION:

CURRENCY:

East US 2

Indian Rupee (₹)

Pricing Details

Azure Data Lake Analytics query service is currently in preview and the prices below reflects a 50% preview discount.

USAGE	PRICE
Analytics Unit	₹1.02 / Minute
Completed Job	₹1.50 / Job

Finally, each Azure Data Lake Analytics account has configurable quotas limiting the number of AUs that can be assigned to jobs and the number of concurrent jobs.

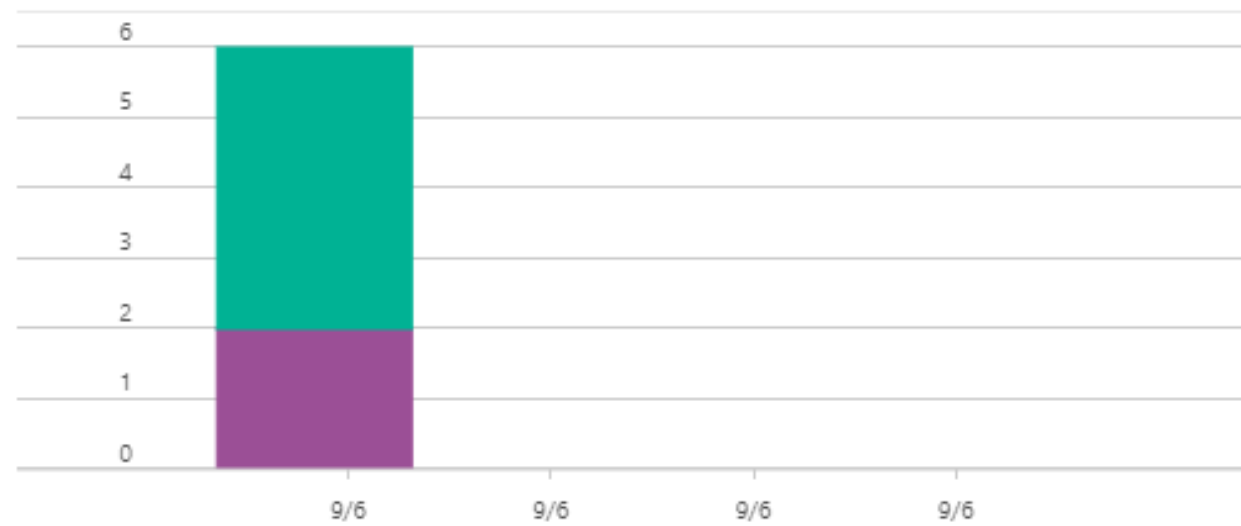
UNIT	LIMIT
# of AU per job	Up to 20 *
# of concurrent jobs per account	Up to 3 jobs *

* To increase level please [contact us](#)

Support & SLA

Azure Data Lake Analytics (ADLA) Demo

Jobs



TOTAL

6

SUCCEEDED

4

FAILED

2

CANCELLED

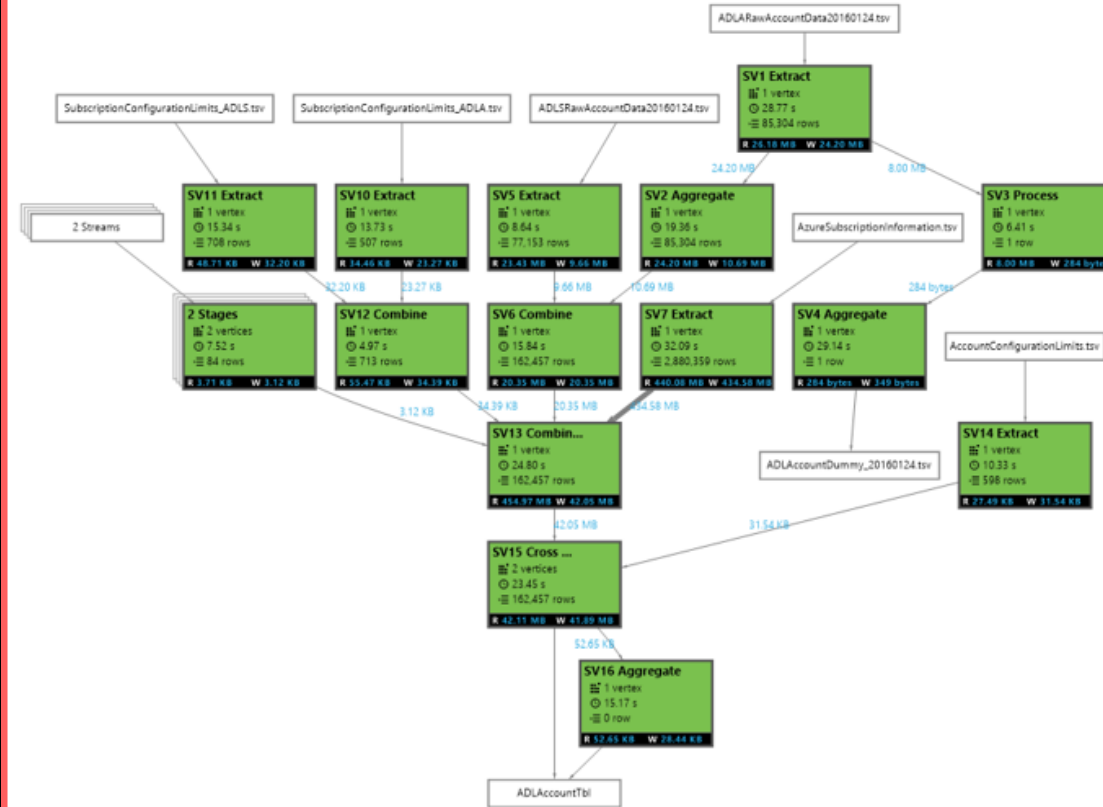
0

ACTIVE

0

Diagnostics

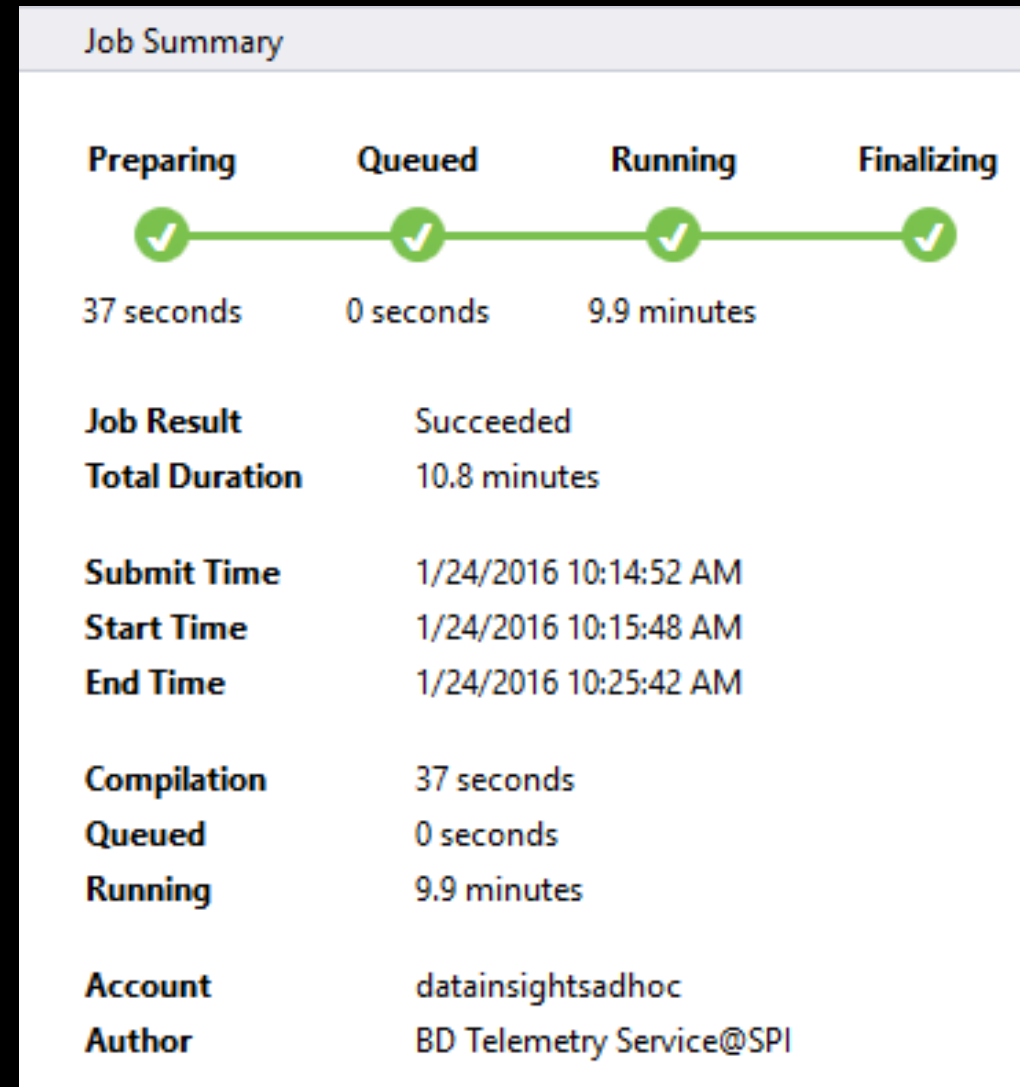
■ Failed ■ Running ■ Waiting Expand ALL



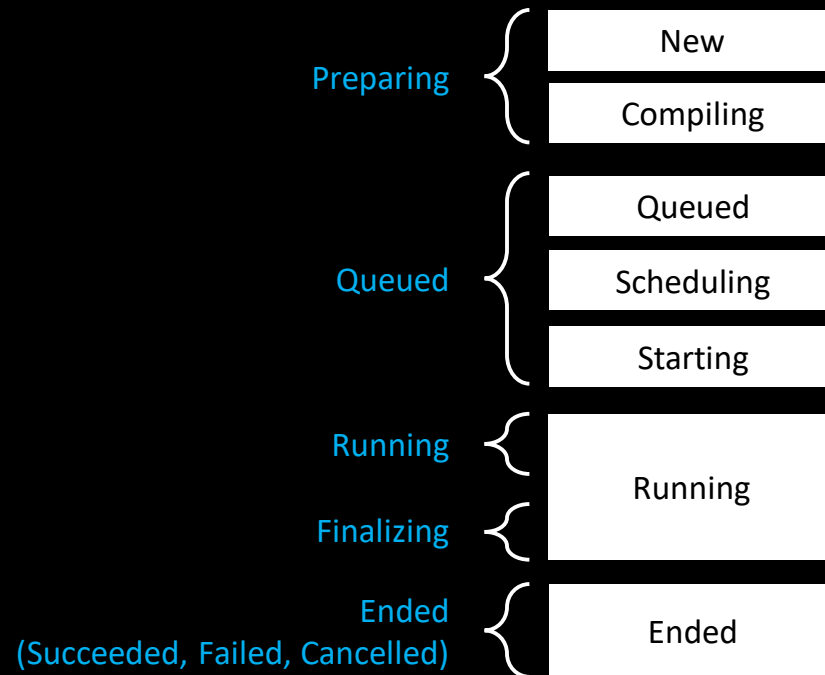
Job Scheduling

States, Queue, Priority

Job Status



UX



The script is being compiled by the Compiler Service

All jobs enter the queue.

Are there enough ADLAUs to start the job?

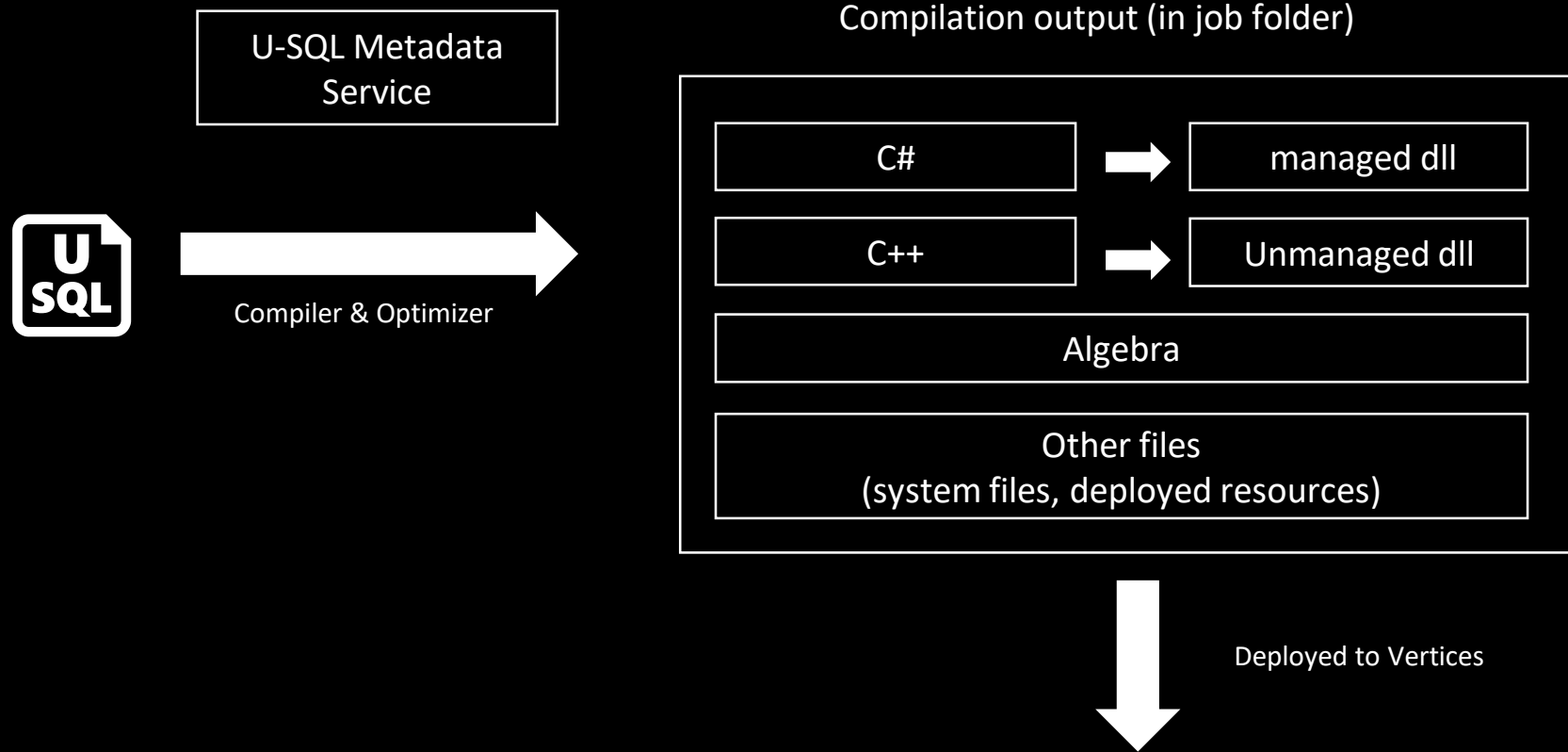
If yes, then allocate those ADLAUs for the job

The U-SQL runtime is now executing the code on 1 or more ADLAUs or finalizing the outputs

The job has concluded.

U-SQL Job Compilation

U-SQL Compilation Process



The Job Folder

Inside the Default ADL Store:

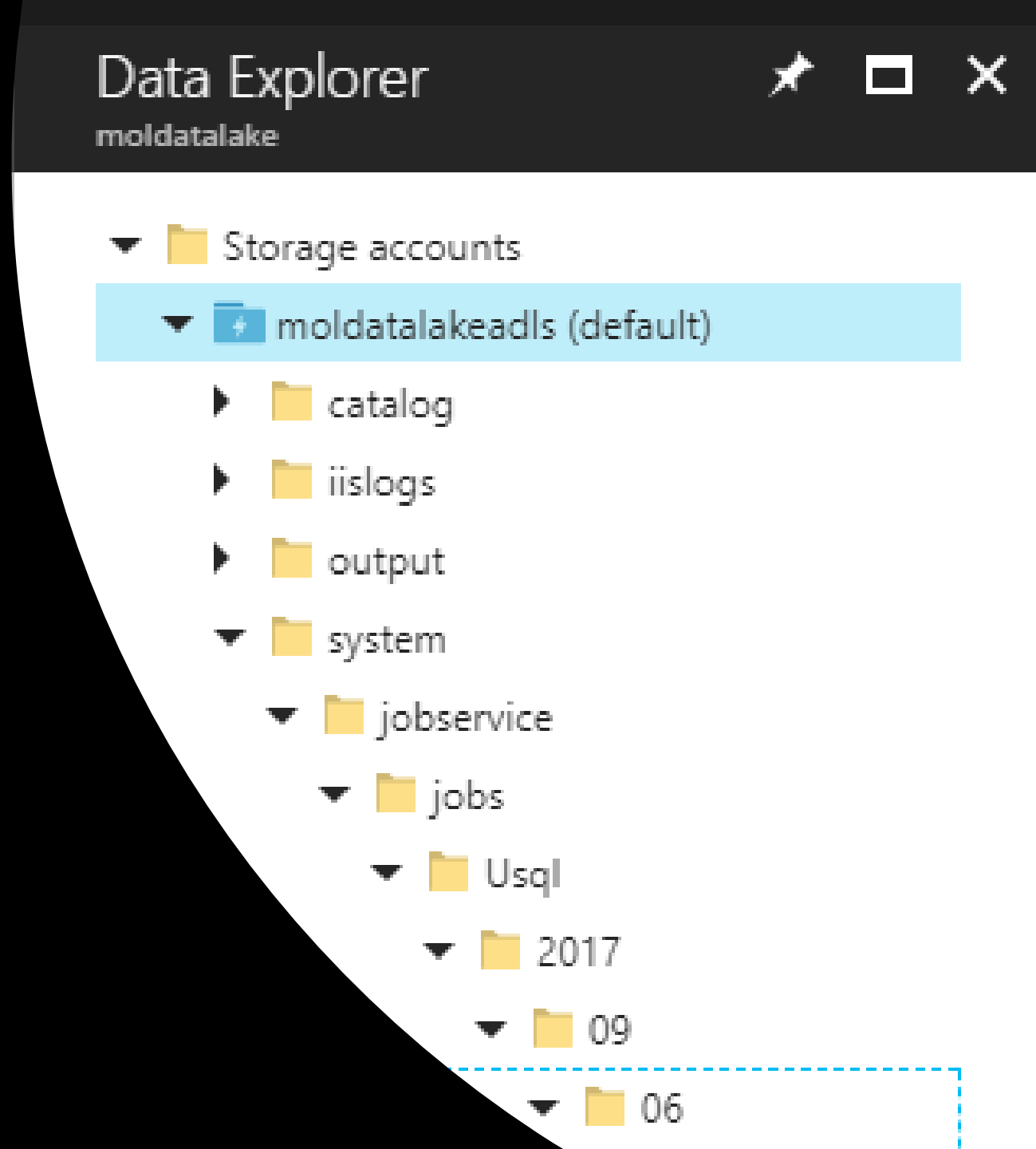
/system/jobservice/jobs/Usql/YYYY/MM/DD/hh/mm/JOBID

/system/jobservice/jobs/Usql/2016/01/20/00/00/17972fc2-4737-48f7-81fb-49af9a784f64



Job folder Structure

- /system/jobservice/jobs/Usql/YYYY/MM/DD/hh/mm/JOBID



ADLAUs

Azure
Data
Lake
Analytics
Unit

Parallelism $N = N$ ADLAUs

1 ADLAU \sim

A VM with 2 cores and 6
GB of memory

Vertex Execution



A VERY BIG FILE

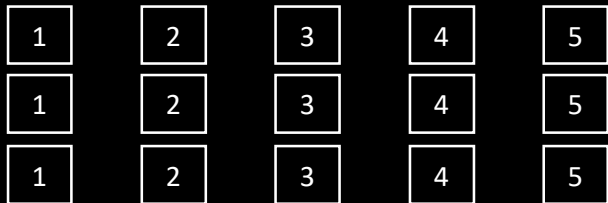
Store Basics

Files are split apart into **Extents**.

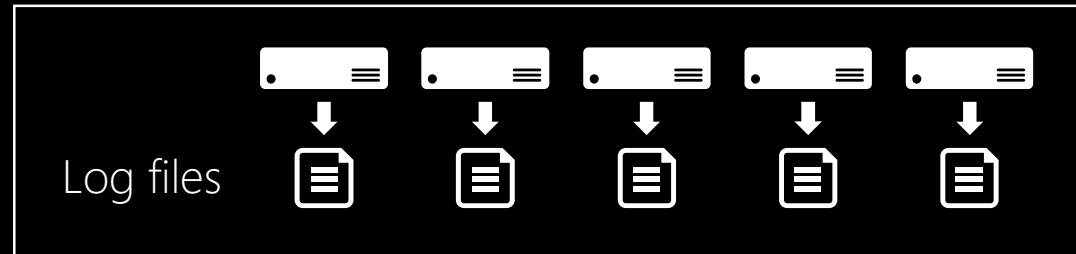
Extents can be up to 250MB in size.

For availability and reliability, extents are replicated (3 copies).

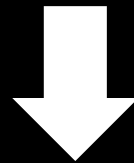
Enables parallelized read



Parallel writing



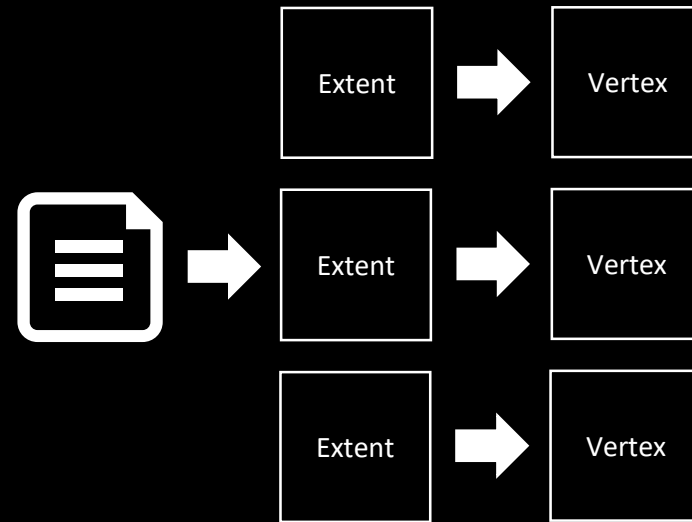
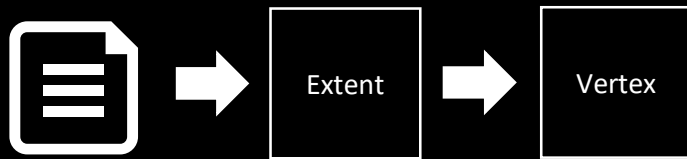
Front-end
machines for a
web service



Simultaneous
uploads



As file size increases, more opportunities for parallelism



ADLA File Management

Hands on

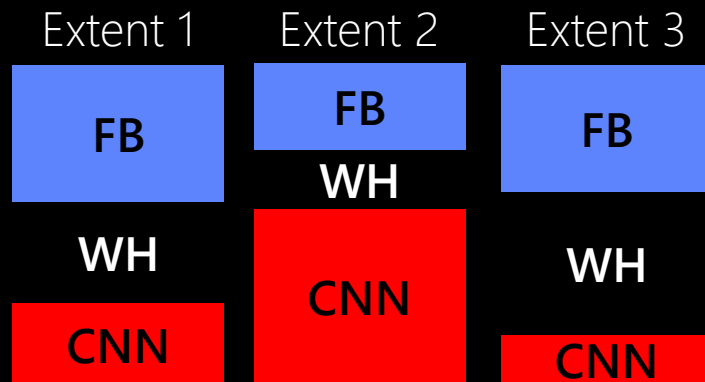
The importance of partitioning input data

Search engine clicks data set

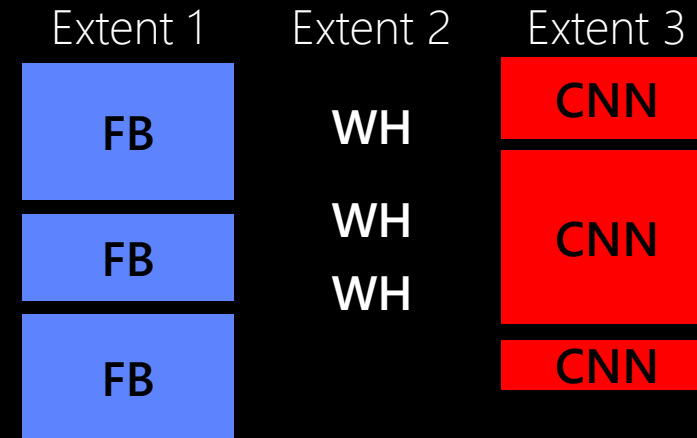
A log of how many clicks a certain domain got within a session

SessionID	Domain	Clicks
3	cnn.com	9
1	whitehouse.gov	14
2	facebook.com	8
3	reddit.com	78
2	microsoft.com	1
1	facebook.com	5
3	microsoft.com	11

Data Partitioning Compared



File:
Keys (Domain) are scattered across
the extents

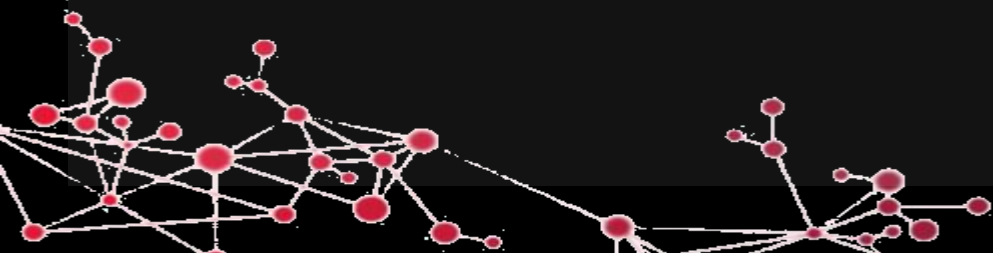


U-SQL Table partitioned on Domain
The keys are now "close together"
also the index tells U-SQL exactly
which extents contain the key

Click Data :

```
• CREATE TABLE SampleDBTutorials.dbo.ClickData  
(  
    SessionId    int,  
    Domain       string,  
    Clicks       int,  
    INDEX idx1 //Name of index  
    CLUSTERED (Domain ASC) //Column to cluster by  
);
```

```
INSERT INTO SampleDBTutorials.dbo.ClickData  
SELECT *  
FROM @clickdata;
```



Find all the rows for cnn.com

// Using a File

```
@ClickData =  
    SELECT  
        Session int,  
        Domain string,  
        Clicks int  
FROM “/clickdata.tsv”  
USING Extractors.Tsv();
```

```
@rows = SELECT *  
    FROM @ClickData  
    WHERE Domain == “cnn.com”;
```

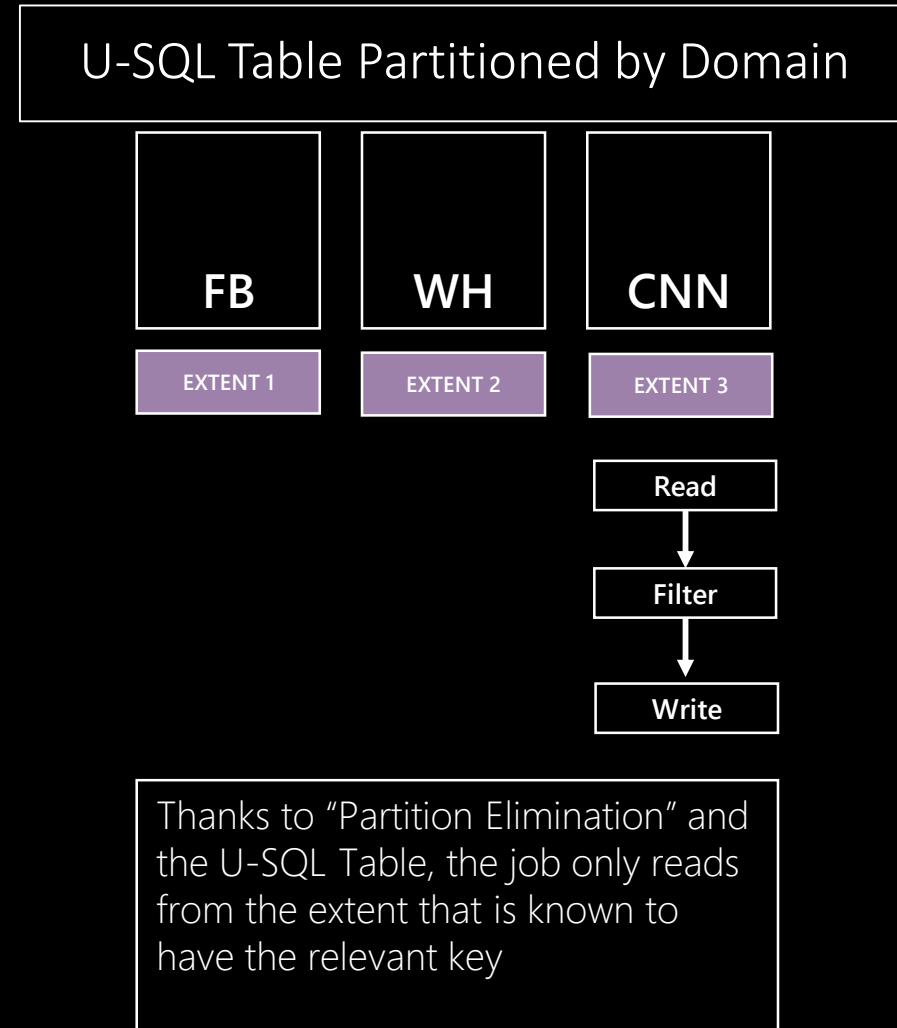
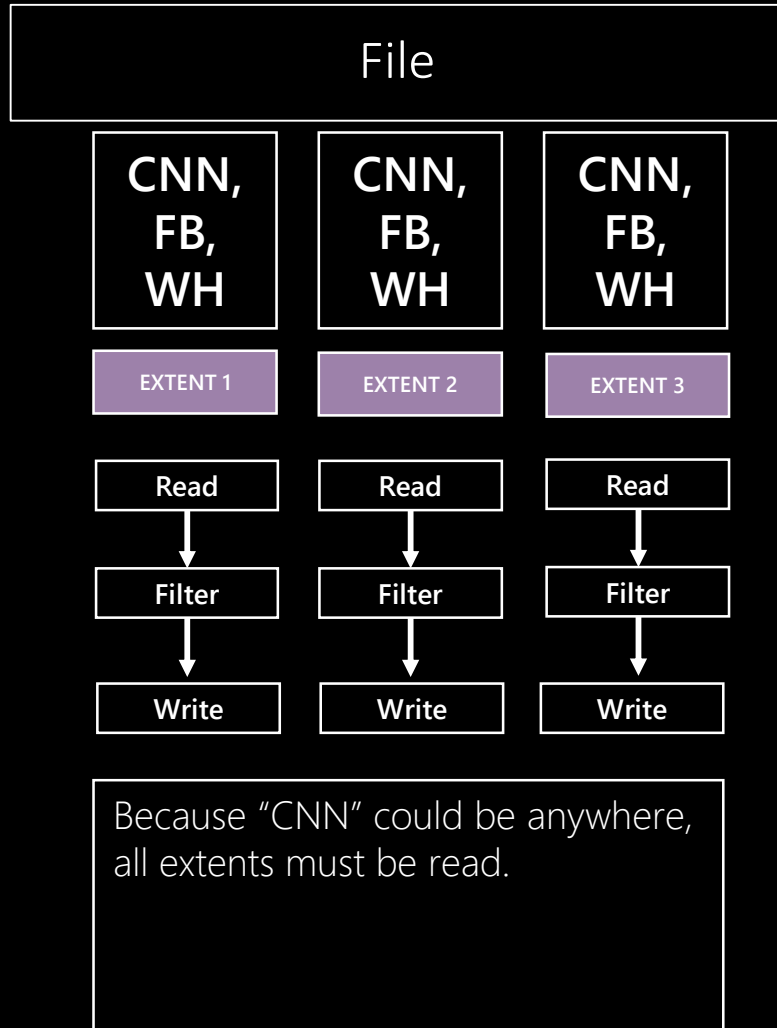
```
OUTPUT @rows  
    TO “/output.tsv”  
    USING Outputters.tsv();
```

// Using a U-SQL Table partitioned by Domain

```
@ClickData =  
    SELECT *  
    FROM MyDB.dbo.ClickData;
```

```
@rows = SELECT *  
    FROM @ClickData  
    WHERE Domain == “cnn.com”;
```

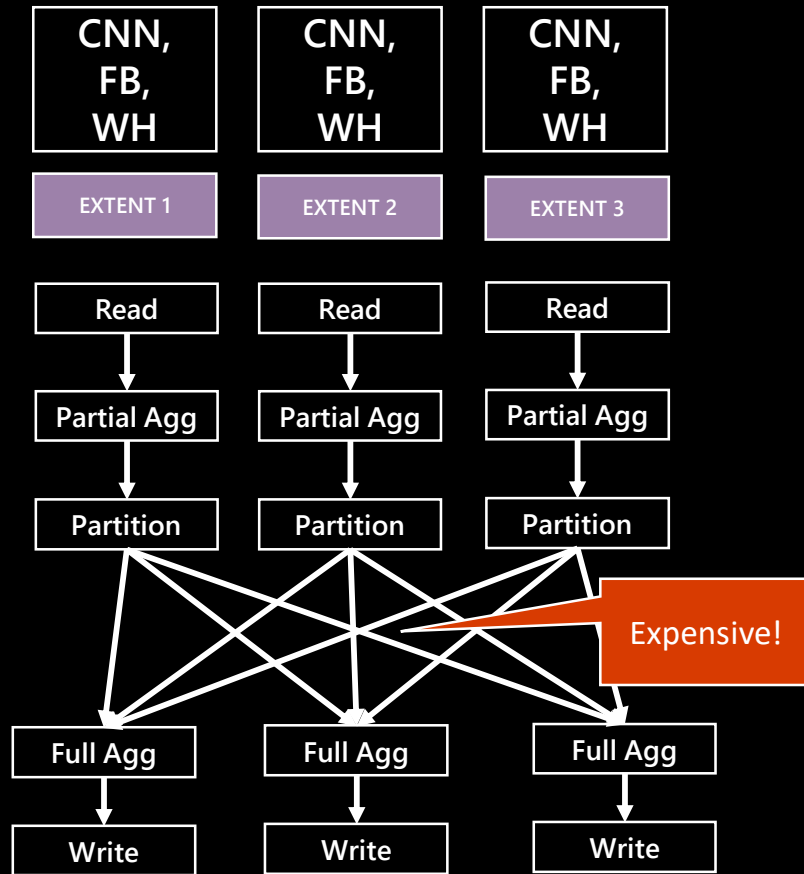
```
OUTPUT @rows  
    TO “/output.tsv”  
    USING Outputters.tsv();
```



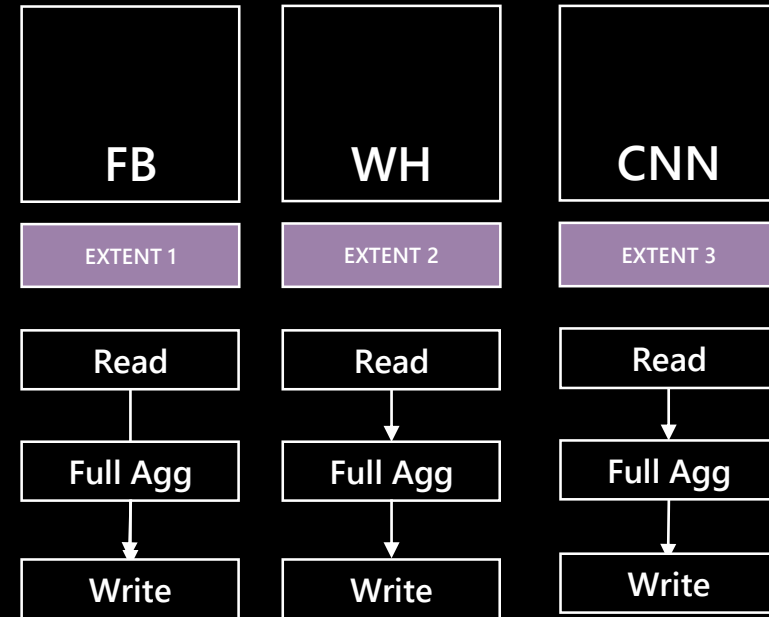
How many clicks per domain?

```
@rows = SELECT
            Domain,
            SUM(Clicks) AS TotalClicks
FROM @ClickData
GROUP BY Domain;
```

File



U-SQL Table Partitioned by Domain



Hands on Lab -ADLA

Creating Data Catalog Objects

How about a Break !!!

Data Factory



INTRODUCTION

AZURE DATA FACTORY



What is **Azure Data Factory**

Introduction to Azure Data Factory Service, a data integration service in the cloud

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



You can create **data integration solutions** using the **Data Factory** service that can ingest data from various data stores, **transform/process** the data, and publish the result data to the data stores.

Azure Data Factory

Enables enterprises

To ingest data from multiple on-premises and cloud sources easily, and gets your data where it needs to go. Prepare and partition your data as you ingest it, or apply pre-processing steps.

Allows you to create data pipelines


That move and transform data, and then run the pipelines on a specified schedule (hourly, daily, weekly, etc.).

Provides rich visualizations

To display the lineage and dependencies between your data pipelines, and monitor all your data pipelines from a single unified view.

Produce trusted data





Data enrichment
Data Integration
Data Migration
Data Transformation
Data Cleansing
IoT Analytics

What can **We** use Data Factory for?

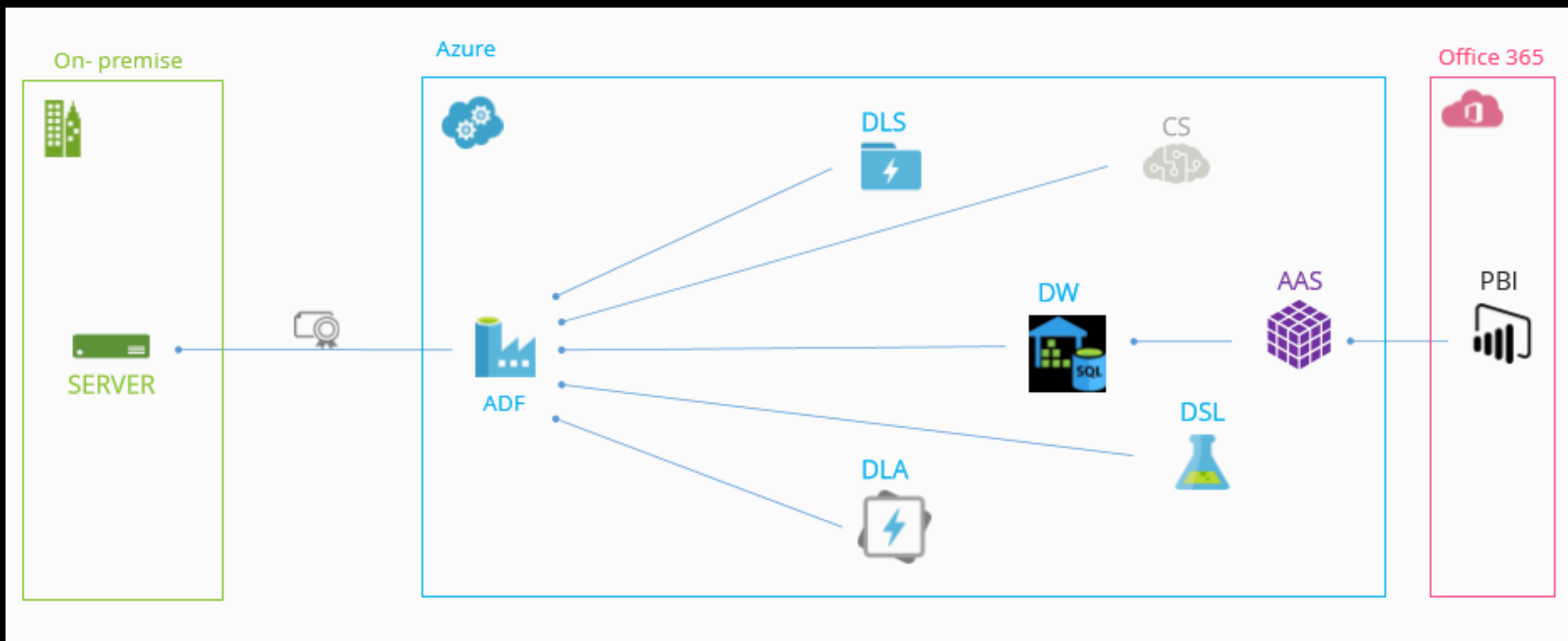
Use it to **ingest data** from multiple on-premises and cloud sources.

Schedule, orchestrate, and manage the **data transformation** and analysis process.

Transform raw data into finished or shaped data that's **ready for consumption** by BI tools or by your on-premises or cloud applications and services.

Manage your entire network of **data pipelines** at a glance to identify issues and take action

Data Sources Usage



What is **Pipeline**




A pipeline is a logical grouping of activities. They are used to group activities into a unit that together perform a task.

To understand pipelines better, you need to understand an activity first.

What is **Activity**



Activities define the actions to perform on your data. For example, you may use a Copy activity to copy data from one data store to another data store. Similarly, you may use a Hive activity, which runs a Hive query on an Azure HDInsight cluster to transform or analyze your data. You may also choose to create a custom .NET activity to run your own code.



PIPELINE





Activity Types

Data Factory supports two types of activities: and.

Data movement activities



Copy Activity in Data Factory copies data from a source data store to a sink data store. Data from any source can be written to any sink.

Data transformation activities



Data Transformation Activity transforms data to desired format and shape. Transformation activities that can be added to pipelines either individually or chained with another activity.

Linked **Services**



Linked services define the information needed for Data Factory to connect to external resources (Examples: Azure Storage, on-premises SQL Server, Azure HDInsight).



Linked **Services**

Linked services are used for two purposes in Data Factory:



To represent a data store

including, but not limited to, an on-premises SQL Server, Oracle database, file share, or an Azure Blob Storage account.



To represent a compute resource

that can host the execution of an activity. For example, the HDInsight Hive activity runs on an HDInsight Hadoop cluster.





Connecting To The Cloud **Data Source**

Linked Services



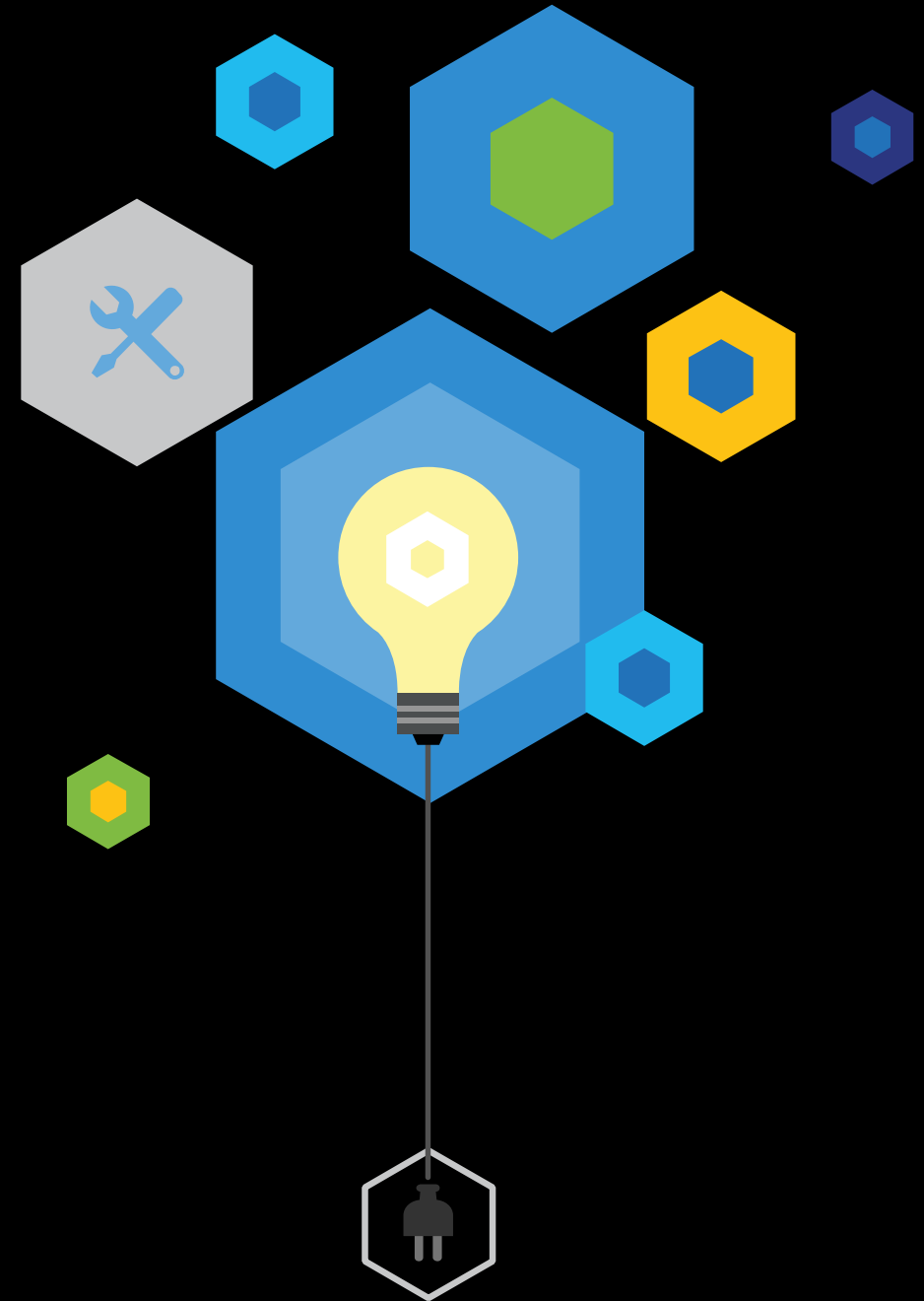
Datasets



Linked services link data stores to an Azure data factory. Datasets represent data structures within the data stores. For example, an Azure Storage linked service provides connection information for Data Factory to connect to an Azure Storage account. An Azure Blob dataset specifies the blob container and folder in the Azure Blob Storage from which the pipeline should read the data. Similarly, an Azure SQL linked service provides connection information for an Azure SQL database and an Azure SQL dataset specifies the table that contains the data.

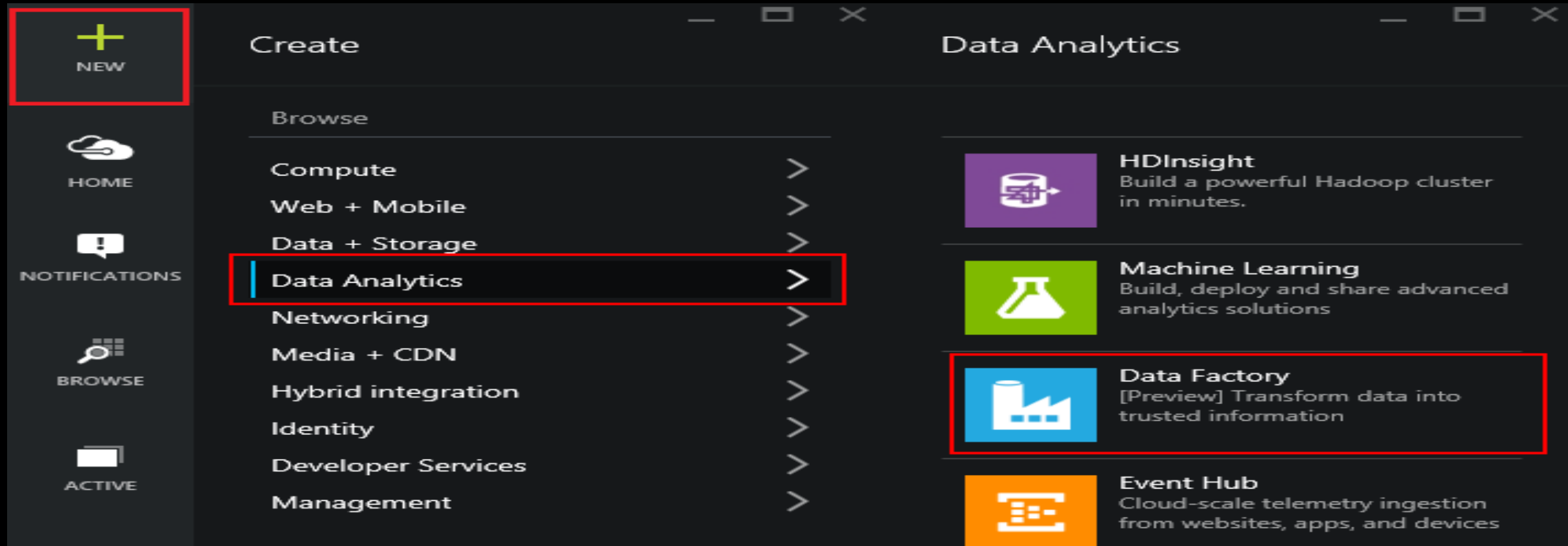
2. Create the Data Factory

Portal, PowerShell and
Visual Studio



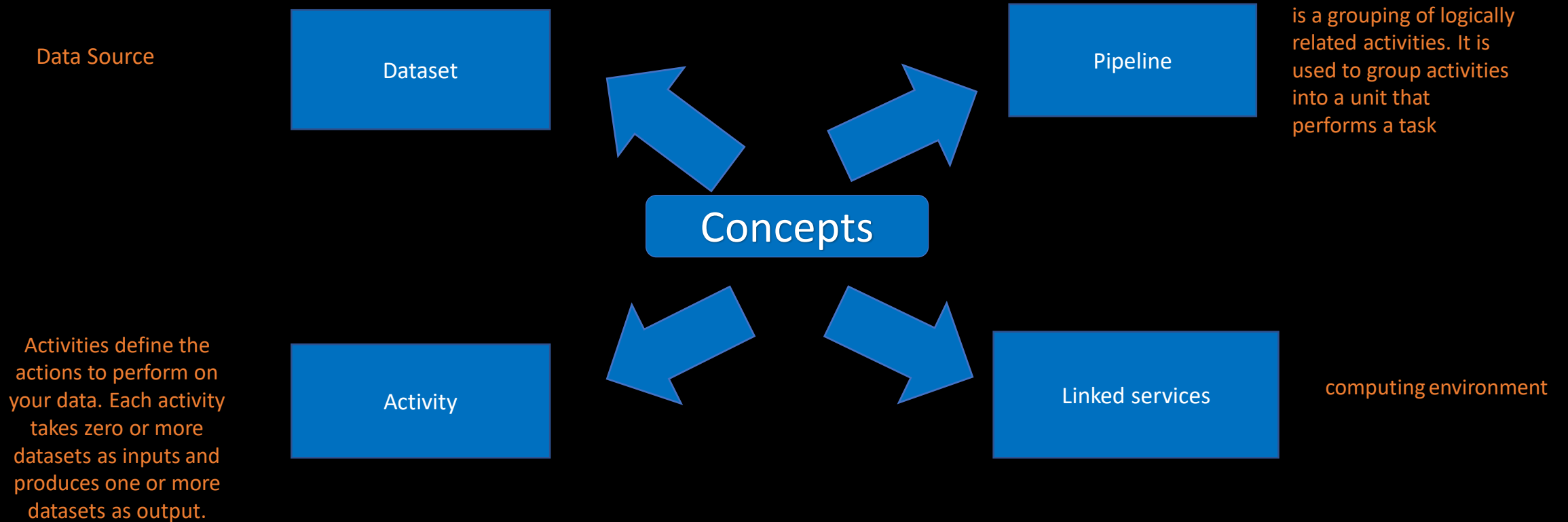
Using the Portal

- Use in Non-MS Clients
- Use for Exploration
- Use when teaching or in a Demo



ADF V2

Concepts



Expressions & Parameters

String functions –
concat, substring,
replace, indexof etc.

**Collection
functions** – length,
union, first, last etc.

Logic functions –
equals, less than,
greater than, and, or,
not etc.

**Conversation
functions** – coalesce,
xpath, array, int,
string, json etc.

Math functions –
add, sub, div, mod,
min, max etc.

Date functions –
utcnow, addminutes,
addhours, format
etc.

System variables

Pipeline scope

Variable Name
@pipeline().DataFactory
@pipeline().Pipeline
@pipeline().RunId
@pipeline().TriggerType
@pipeline().TriggerId
@pipeline().TriggerName
@pipeline().TriggerTime

Trigger scope

Variable Name
trigger().scheduledTime
trigger().startTime



Trigger

Type of triggers

- Manual execution
- Schedule trigger: A trigger that invokes a pipeline on a wall-clock schedule.
- Tumbling window trigger: A trigger that operates on a periodic interval, while also retaining state.

Schedule trigger

A schedule trigger runs pipelines on a wall-clock schedule. This trigger supports periodic and advanced calendar options. For example, the trigger supports intervals like "weekly" or "Monday at 5:00 PM and Thursday at 9:00 PM."

```
{
  "properties": {
    "name": "MyTrigger",
    "type": "ScheduleTrigger",
    "typeProperties": {
      "recurrence": {
        "frequency": "Hour",
        "interval": 1,
        "startTime": "2017-11-01T09:00:00-08:00",
        "endTime": "2017-11-02T22:00:00-08:00"
      }
    },
    "pipelines": [{
      "pipelineReference": {
        "type": "PipelineReference",
        "referenceName": "SQLServerToBlobPipeline"
      },
      "parameters": {}
    },
    {
      "pipelineReference": {
        "type": "PipelineReference",
        "referenceName": "SQLServerToAzureSQLPipeline"
      },
      "parameters": {}
    }
  ]
}
```

Tumbling window trigger

Tumbling window triggers are a type of trigger that fires at a periodic time interval from a specified start time, while retaining state. Tumbling windows are a series of fixed-sized, non-overlapping, and contiguous time intervals.

```
{
  "name": "PerfTWTrigger",
  "properties": {
    "type": "TumblingWindowTrigger",
    "typeProperties": {
      "frequency": "Minute",
      "interval": "15",
      "startTime": "2017-09-08T05:30:00Z",
      "delay": "00:00:01",
      "retryPolicy": {
        "count": 2,
        "intervalInSeconds": 30
      },
      "maxConcurrency": 50
    },
    "pipeline": {
      "pipelineReference": {
        "type": "PipelineReference",
        "referenceName": "DynamicsToBlobPerfPipeline"
      },
      "parameters": {
        "windowStart": "@trigger().outputs.windowStartTime",
        "windowEnd": "@trigger().outputs.windowEndTime"
      }
    },
    "runtimeState": "Started"
  }
}
```




Control flow

Control flow

ForEach activity

Filter activity

Get metadata

Execute Pipeline

If Condition activity

Web activity

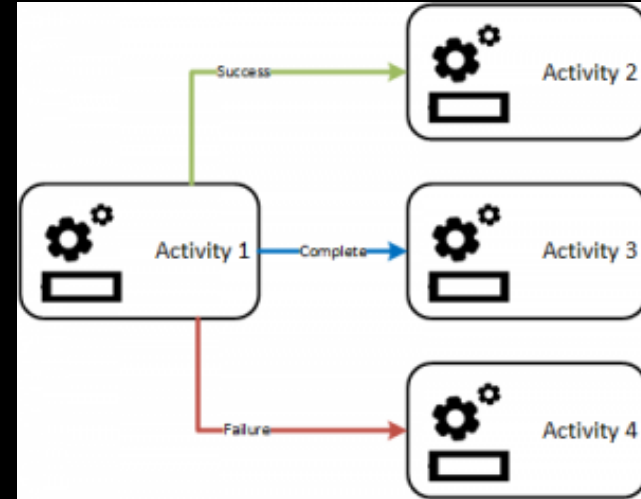
Lookup activity

Wait activity

Until activity

Branching

- On success
- On failure
- On completion
- On skip



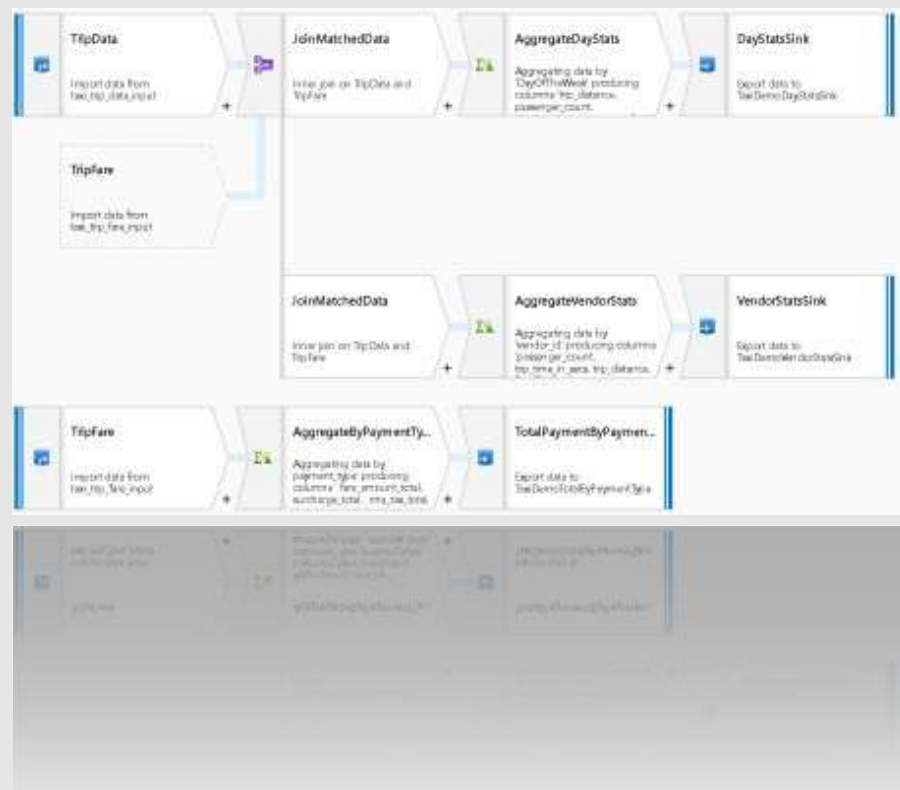
Hands on Usecase

- One of the Manufacturing client needs data coming from multiple sources to be automated to generate reports and even send the email at the end.
- This should be a one time creation of pipeline but any new source integration should not impact the pipeline process.

Usecase

Build 2019

Whats New !!!



```

// Create a DataFrame from the input data
val tripData = spark.read.parquet("taxi_trip_data.parquet")
val tripfare = spark.read.parquet("taxi_fare_tripdata.parquet")

// Inner join on TripData and Tripfare
val joinedData = tripData.join(tripfare, "trip_id", "inner")

// Aggregating data by DayOfWeek producing columns trip_distance, passenger_count
val dayStats = joinedData.groupBy("DayOfWeek").agg(
  sum("trip_distance").as("total_trip_distance"),
  sum("passenger_count").as("total_passenger_count")
)

// Export data to TaxiDemoDayStatsSink
dayStats.write.parquet("taxi_demo_day_stats.parquet")

// Inner join on TripData and Tripfare
val joinedData2 = tripData.join(tripfare, "trip_id", "inner")

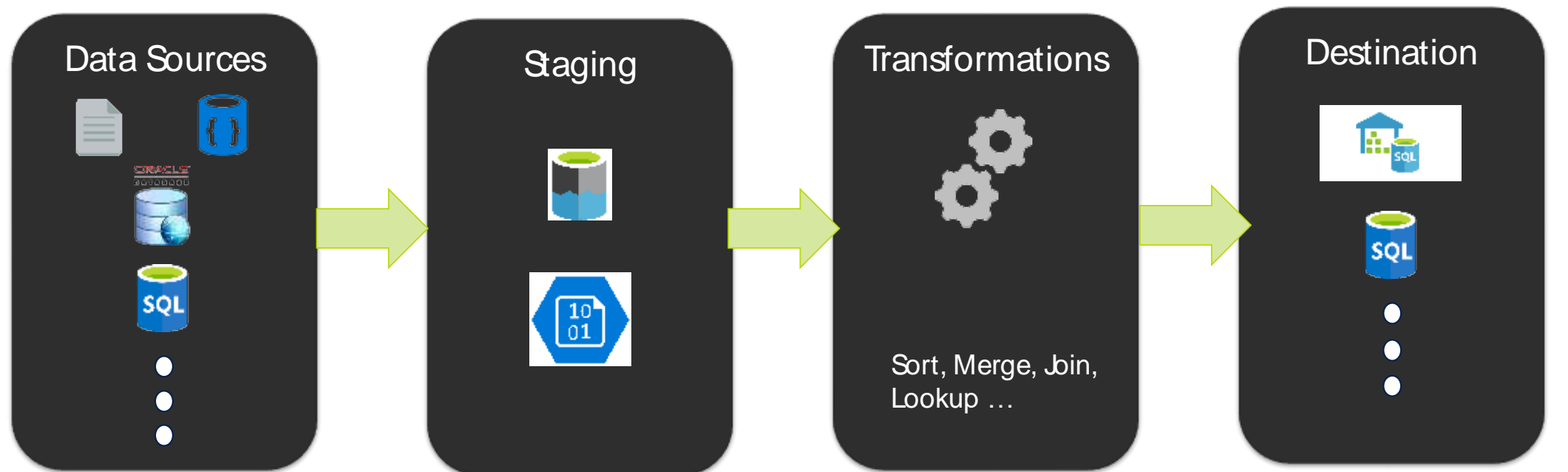
// Aggregating data by Vendor id producing columns passenger_count, trip_time_in_secs, trip_distance
val vendorStats = joinedData2.groupBy("Vendor id").agg(
  sum("passenger_count").as("total_passenger_count"),
  sum("trip_time_in_secs").as("total_trip_time_in_secs"),
  sum("trip_distance").as("total_trip_distance")
)

// Export data to TaxiDemoVendorStatsSink
vendorStats.write.parquet("taxi_demo_vendor_stats.parquet")

// Aggregating data by payment_type producing columns fare_amount,total_amount,total_mta_tax,total
val paymentTypeStats = tripfare.groupBy("payment_type").agg(
  sum("fare_amount").as("total_fare_amount"),
  sum("total_amount").as("total_amount"),
  sum("total_mta_tax").as("total_mta_tax"),
  sum("total").as("total")
)

// Export data to TaxiDemoTotalPaymentTypeSink
paymentTypeStats.write.parquet("taxi_demo_total_payment_type_stats.parquet")
  
```


ADF Data Flow Workstream



- Explicit user action
- User places data source(s) on design surface, from toolbox
- Select explicit sources

- Implicit/Explicit
- Data Lake staging area as default
- User does not need to configure this manually
- Advanced feature to set staging area options
- File Formats / Types (Parquet, JSON, txt, CSV ...)

- Explicit user action
- User places transformations on design surface, from toolbox
- User must set properties for transformation steps and step connectors

- Explicit user action
- User chooses destination connector(s)
- User sets connector property options

TaxiDemo X

CopyFlow * X

Save

Debug

DSL

source1

Add source dataset

+

sink1

Add sink dataset

Add Source

+

—

🔒

100%

🔍

🖱️

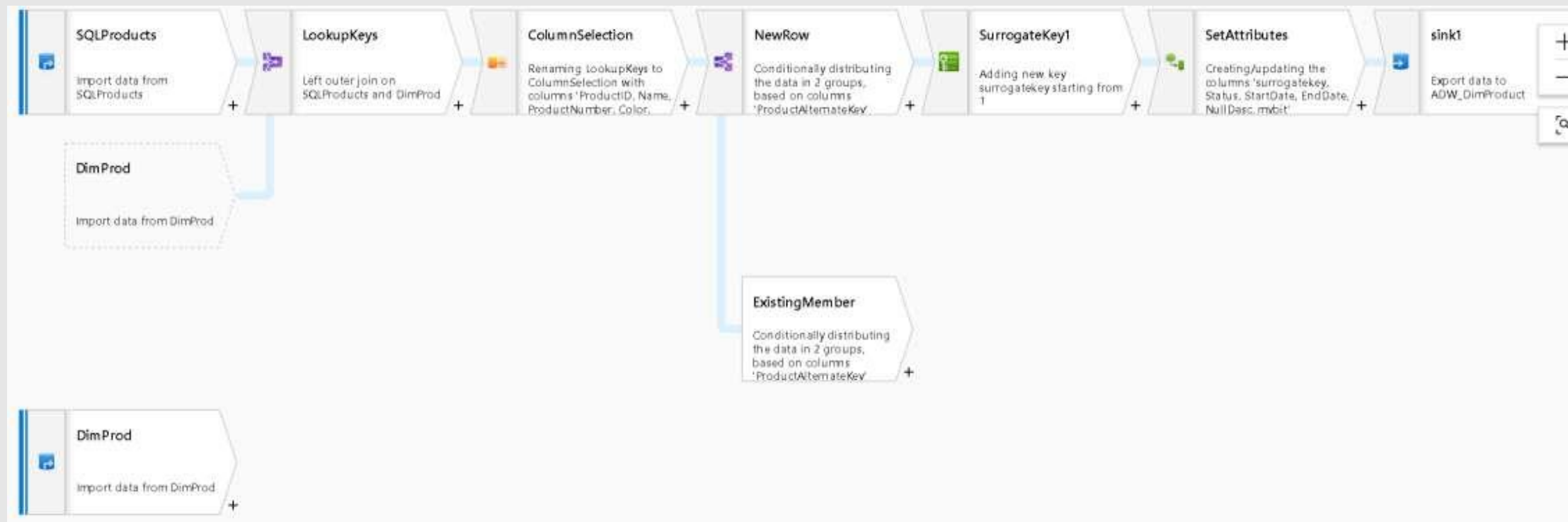
📏

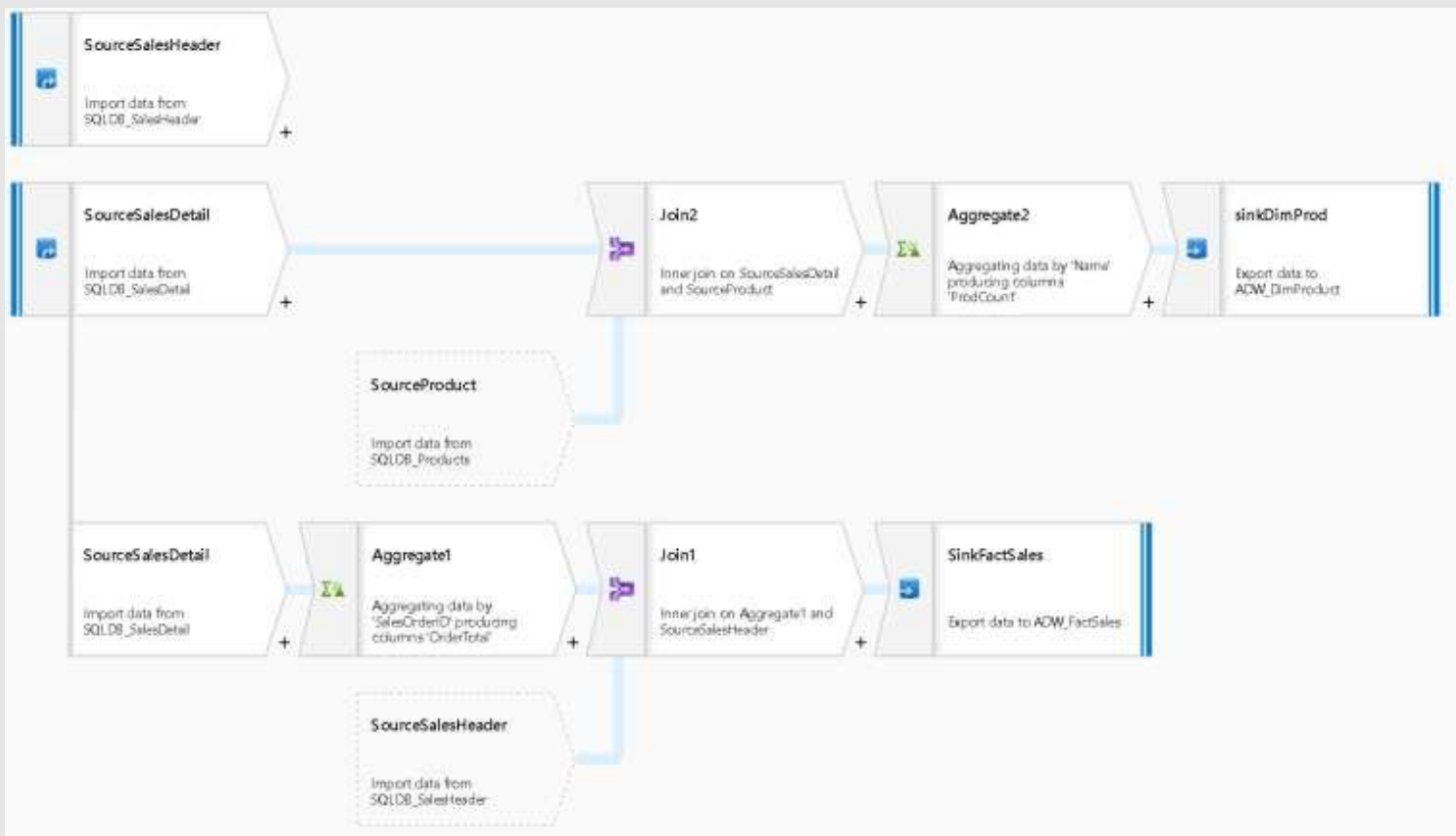
🔲

General

Name *

CopyFlow





CurrencyCleanup * Pipeline 134* Pipeline 76*

On Debug Save DSL

OriginalCurrencyData

Import data from Original Currency Data Azure Data Lake

Join A

Join on (Original Currency : Data.CurrencyName == Daily Currency #file.CurrencyID) & ...

USD

Split on (CurrencyName == 'USD')

Original Currency Data

Columns: 13 total
Rows: 28,024 total

+ - 🔍 📄 🔒 🔄 🗑️

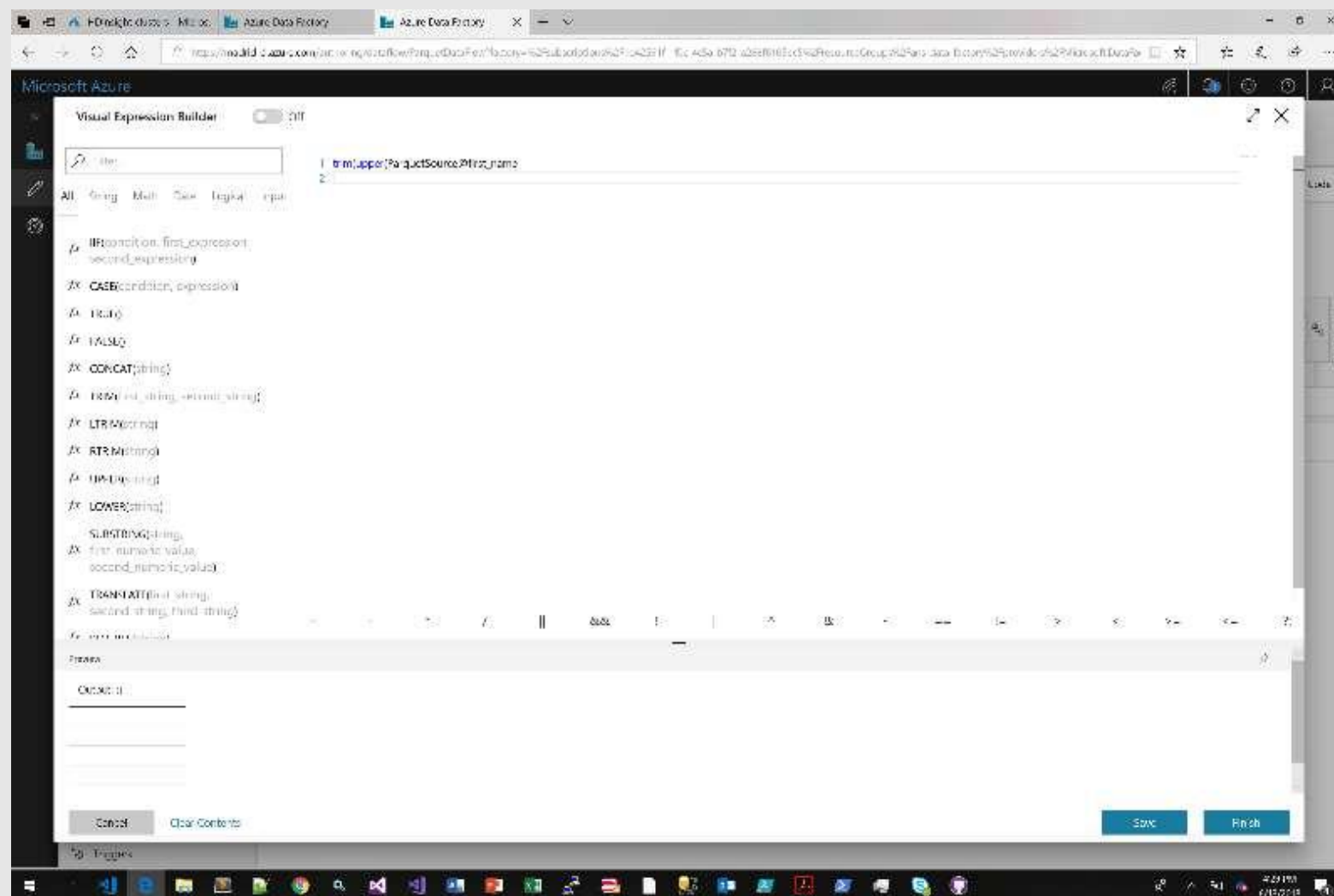
Setting Inspect Error log

	Update *	New *	Unchanged	Total
Nuber of columns	2	1	15	18
Nuber of rows	30	0	2,483	2,234

Output schema Data Preview

Column ▾

	Date * 📅	InUSA * 🇺🇸	Profit * 1.2	Column 123	Column abc	Column abc	Column 📅
1	12/03/2018	True	2455.45	12345	Cell Contents	Cell Contents	09/23/2017, 23:00
2	12/03/2018	False	2455.45	12345	Cell Contents	Cell Contents	09/23/2017, 23:00
3	12/03/2018	True	2455.45	12345	Cell Contents	Cell Contents	09/23/2017, 23:00
4	12/03/2018	False	2455.45	12345	Cell Contents	Cell Contents	09/23/2017, 23:00
5	12/03/2018	False	2455.45	12345	Cell Contents	Cell Contents	09/23/2017, 23:00
6	12/03/2018	False	2455.45	12345	Cell Contents	Cell Contents	09/23/2017, 23:00



DataFlowName-1

Data Flow

30%

Number of transforms: 32



MatchedData

Aggregate

Columns Partitions

Total columns	50
New columns*	2
Updated column*	1
Dropped column	50

Column information

Order	Column	Type	Create by
1	Date	date	IS MatchedData
2	InUSA	boolean	VS InUSA
3	Profit	double	PP Profit
4	CustomerID	int	IS CustomerID
5	CustomerName	string	PP CustomerName
6	CustomerID_sum	int	IS CustomerID_sum

VendorStatsSink

TRANSFORM	ROWS	TIME(ms)
• TripData	3m	
• TripFare	3,420	
• MatchedData	54k	3,420
• VendorStats	14k	14k
• VendorStatsSink	1,000	1,000

DayStatsSink

TRANSFORM	ROWS	TIME(ms)
• TripData		
• TripFare		
• MatchedData	54k	3,420
• USD		
• DayStatsSink		

TotalByPaymentType

TRANSFORM	ROWS	TIME(ms)
• TripData	-	-
• AggregateByPay...	-	-
• TotalByPayment...	-	-

Sink4

TRANSFORM	ROWS	ROWS
• Source3	-	-
• Transform1	-	-
• TripData	-	-
• TripFare	-	-
• MatchedData	-	-
• Transform2	-	-
• Transform3	-	-
• Transform4	-	-

	Update ⁺	New ⁺	Unchanged	Total
Nuber of columns	2	1	15	18
Nuber of rows	30	0	2,483	2,234

Output schema Data Preview

Order	Date ⁺	Column ⁺	Profit ⁺ 1,2	Column 123	Continent_data abc	Column abc
1	12/03/2018	✓	2455.45	Cell Contents	North America	Cell Contents
2	12/03/2018	✓	2455.45	Cell Contents	Europe	Cell Contents
3	12/03/2018		2455.45	Cell Contents	North America	Cell Contents
4	12/03/2018		2455.45	Cell Contents	North America	Cell Contents
5	12/03/2018	✓	2455.45	Cell Contents	Africa	Cell Contents
6	12/03/2018	✓	2455.45	Cell Contents	Asia	Cell Contents
7	12/03/2018	✓	null	Cell Contents	Europe	Cell Contents
8	12/03/2018		null	Cell Contents	North America	Cell Contents
9	12/03/2018		49582.23	Cell Contents	Asia	Cell Contents
10	12/03/2018		49582.23	Cell Contents	Asia	Cell Contents
11	12/03/2018	✓	null	Cell Contents	Europe	Cell Contents
12	12/03/2018		49582.23	Cell Contents	Europe	Cell Contents
13	12/03/2018		49582.23	Cell Contents	North America	Cell Contents

Continent_data abc



Content Type	String
Missing Value	0 (0%)
Unique Values	5 (5%)
Most Common	North America (52.8%)
Skewness	-0.043210436

Future is most promising for ADF

Data Flows are in Public Preview since Build 2019. (Since 1 Week)

Hands on Usecase

- One of the Manufacturing client needs data coming from multiple sources to be automated to generate reports and even send the email at the end.
- This should be a one time creation of pipeline but any new source integration should not impact the pipeline process.

Usecase