Microsoft

databricks

# Azure Advanced Analytics engine for Data Science

**Adrián J. Fernández Zenteno**

*Cloud Solution Architect – Azure Advanced Analytics & Data Science*
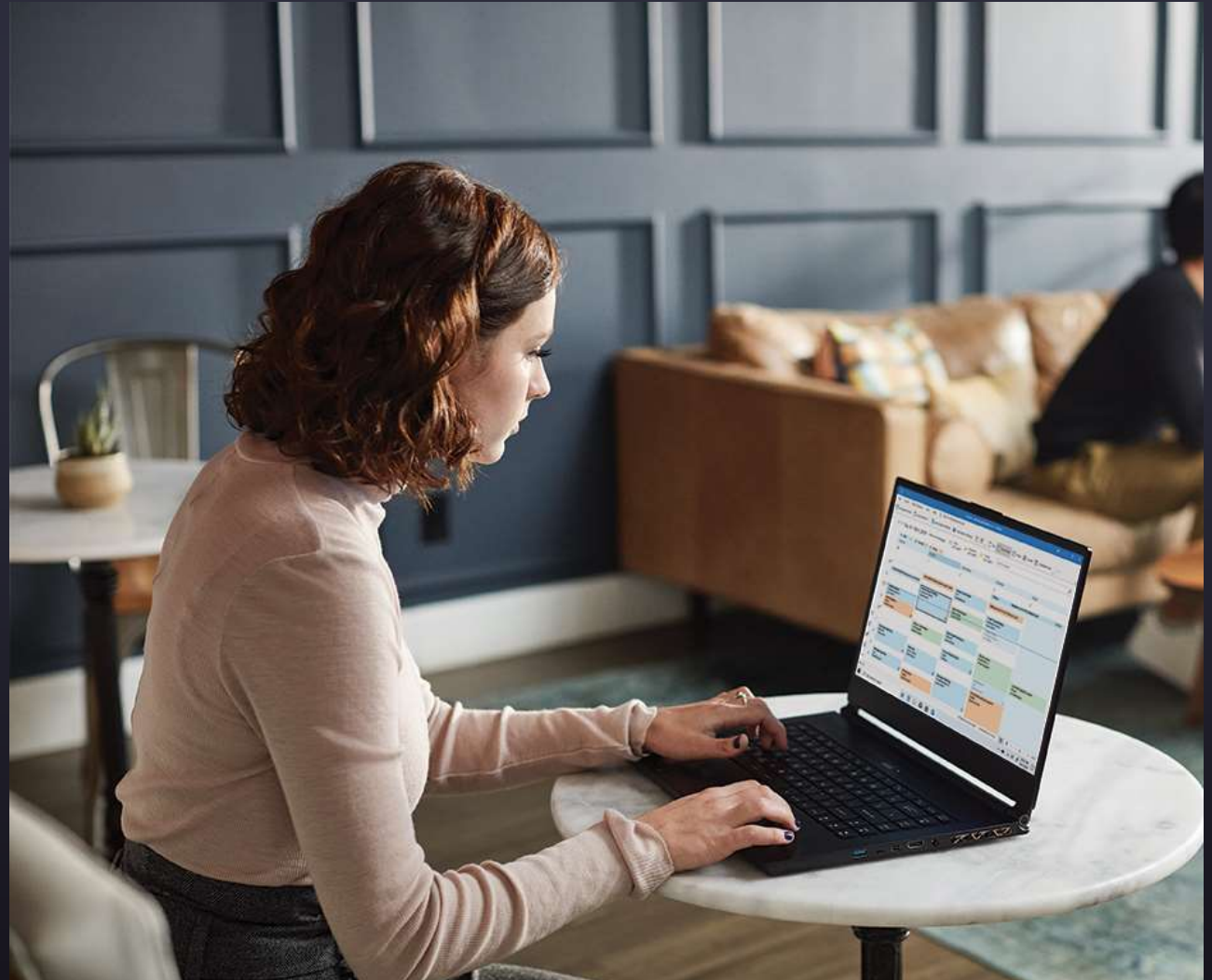
Email: adrian.fernandez@microsoft.com

Twitter: @AdrianFZ10

APACHE Spark

# Agenda

- ➢ Azure Databricks
- ➢ Spark Summit 2020
- ➢ Synapse Analytics Studio
- ➢ Execution Pools engine:
  - ➢ SQL Analytics / On-Demand
  - ➢ Apache Spark
- ➢ Synapse and Azure ML integration

# Machine Learning on Azure

**Domain specific pretrained models**
To simplify solution development

Vision    Speech    Language    Search

**Familiar data science tools**
To simplify model development

Visual Studio Code    Azure Notebooks    Jupyter    Command line

**Popular frameworks**
To build advanced deep learning solutions

PyTorch    TensorFlow    Scikit-Learn    ONNX

**Productive services**
To empower data science and development teams

Azure Machine Learning    Azure Databricks    Synapse Analytics    ML VMs

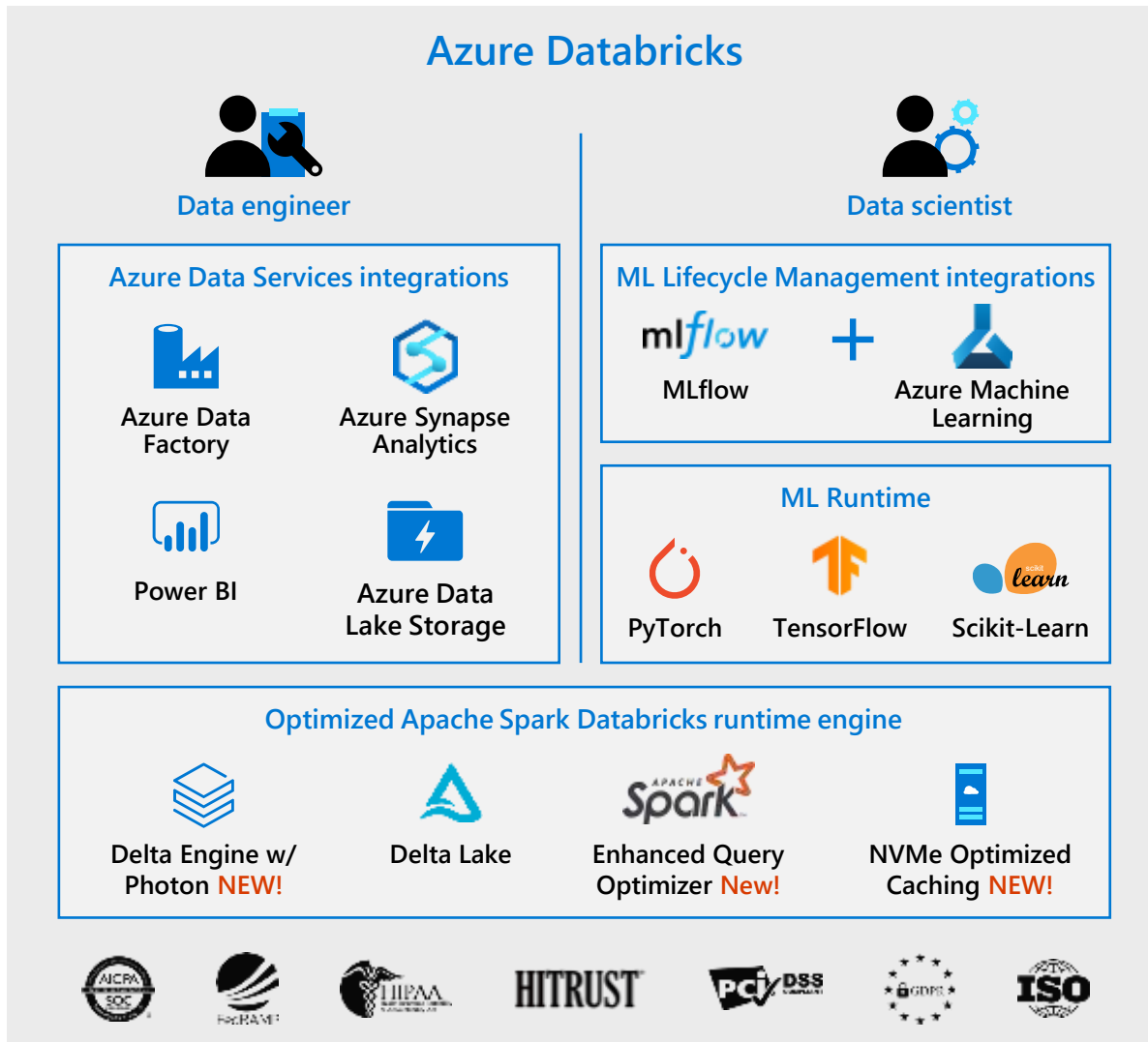**Powerful infrastructure**
To accelerate deep learning

CPU    GPU    FPGA

**From the intelligent Cloud to the intelligent Edge**

# Azure Databricks – Introduction

## Azure Databricks

**Data engineer**

**Data scientist**

### Azure Data Services integrations

Azure Data Factory

Azure Synapse Analytics

Power BI

Azure Data Lake Storage

### ML Lifecycle Management integrations

MLflow + Azure Machine Learning

### ML Runtime

PyTorch

TensorFlow

Scikit-Learn

### Optimized Apache Spark Databricks runtime engine

Delta Engine w/ Photon **NEW!**

Delta Lake

Enhanced Query Optimizer **New!**

NVMe Optimized Caching **NEW!**

---

**Collaborative**
Workspaces for data teams across the full lifecycle

**Connected**
Native integration with the entire Azure Portfolio
Leverage the most popular open source tools

**Fast**
Scalable and reliable data powered by the fastest Spark Engine on the market

**Secure**
Azure Active Directory Single Sign-On

# Azure Databricks – Top Announcements 📣

- Databricks Runtime 7.x with Apache Spark 3.0!

- Delta Engine with Photon!

- Koalas

- Redash

- Workspace 2.0

- Azure US Gov Preview with FedRAMP High Certification
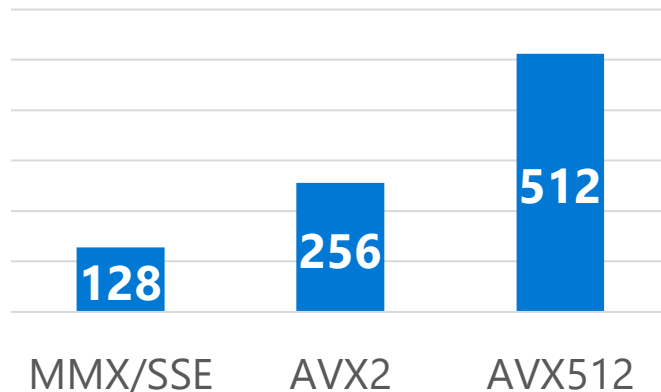
- Spark AI Summit

https://databricks.com/sparkaisummit/north-america-2020

# Azure Databricks - Delta Engine

# Azure Databricks – Delta Engine Motivation

## Hardware Trends

|  | 2010 | 2015 | 2020 |  |
|---|---|---|---|---|
| Storage | 50 MB/s (HDD) | 500 MB/s (SSD) | 16 GB/s (NVMe) | 10X |
| Network | 1 Gbps | 10 Gbps | 100 Gbps | 10X |
| CPU | ~3 GHz | ~3 GHz | ~3 GHz | ☹️ |

## So where are CPU innovations going?

### Data level parallelism (SIMD register width)

| MMX/SSE | AVX2 | AVX512 |
|---|---|---|
| 128 | 256 | 512 |

### Instruction level parallelism:

| Sandy Bridge | Haswell | Skylake |
|---|---|---|
| 168 | 192 | 224 |

# Azure Databricks – Delta Engine

- Builds On Apache Spark 3.0
- 100% Spark Compatibility
- Fully API compatible
- Accelerates SQL and DataFrame workloads with:
  - Improved query optimizer
  - Native vectorized execution engine
  - Caching



| SQL | Spark DataFrame | Koalas |

Query Optimizer — **Today**

Native Execution Engine Photon — **Coming 2021**

Caching — **Today**

**Delta Engine**

# Delta Engine's Improved Query Optimizer

- Extends Spark's cost-based optimizer and adaptive query execution with advanced stats

- Up to 18x performance increase for star schema workloads



SQL  Spark DataFrame  Koalas

Query Optimizer

Native Execution Engine
Photon

Caching

**Delta Engine**

# Delta Engine's Caching

- Automatically caches input

- Transcodes data into a more CPU-efficient format fully leveraging NVMe SSDs

- Up to 5x scan performance increase

SQL

Spark DataFrame

Koalas

Query Optimizer

Native Execution Engine Photon

Caching

**Delta Engine**

# Delta Engine Photon

- New execution engine for Delta Engine to accelerate Spark SQL

- Built from scratch in C++, for performance:
    - Vectorization: exploit data-level parallelism and instruction-level parallelism
    - Optimize for modern structured and semi-structured workloads

SQL      Spark DataFrame      Koalas

Query Optimizer

Native Execution Engine Photon

Caching

**Delta Engine**

# Azure Databricks – Workspace 2.0

# Azure Databricks – Projects API for CI/CD

Development / Experimentation

Production Jobs

dev    Projects Service    main

**Projects API**

Programmatically manage Projects

Perform Git operations on Projects

Git / CI/CD Systems

Version → Review → Test

Specify Git branch, tag, hash for Jobs

Azure DevOps    Bitbucket    GitLab    GitHub

# Repository-level git integration

# Azure Databricks - Redash

# Azure Databricks – Redash: A home for SQL Users

**Self-serve on the Data Lake!**
Collaborative queries, dashboards &
alerts on your data lake

**Simple, SQL oriented UX**
Analysts don't need to understand or
get exposed to notebooks or jobs

**Ready-to-go**
Tightly integrated with Databricks
compute & security

# Spark Interfaces

**Resilient Distributed Dataset (RDD)**

Spark RDD is a resilient, partitioned, distributed and immutable collection of data.

**DataFrame**

Distributed collection of data organized into named columns. It is conceptually equivalent to a table in a relational database or a data frame in R/Python, but with richer optimizations under the hood

**Dataset**

A Dataset is a strongly-typed, immutable collection of objects that are mapped to a relational schema.

An extension of the DataFrame API that provides a *type-safe, object-oriented programming interface.*

# DataFrames

DataFrame is a distributed collection of data organized into named columns.

Conceptually equivalent to a table in a relational database or a data frame in R/Python, but with richer optimizations.

DataFrames can be constructed from a wide array of sources such as: structured data files, tables in Hive, external databases, or existing RDDs.

DataFrames are evaluated lazily, meaning, computation only happens when an action (e.g. display result, save output) is required.

# Loading Data in DataFrames

```python
%python
# Use the Spark CSV datasource with options specifying:
# - First line of file is a header
# - Automatically infer the schema of the data

data = spark.read.format("csv") \
  .option("header", "true") \
  .option("inferSchema", "true") \
  .load("/databricks-datasets/samples/population-vs-price/data_geo.csv")

data.cache() # Cache data for faster reuse
data = data.dropna() # drop rows with missing values
```

# Viewing DataFrames

## Using Spark Command **take()** to view raw records

```python
%python data.take(10) #view 10 records of DataFrame
```

▸ (1) Spark Jobs

Out[3]:
```
[Row(2014 rank=101, City=u'Birmingham', State=u'Alabama', State Code=u'AL', 2014 Population estimate=212247, 2015 median sales price=162.9),
 Row(2014 rank=125, City=u'Huntsville', State=u'Alabama', State Code=u'AL', 2014 Population estimate=188226, 2015 median sales price=157.7),
 Row(2014 rank=122, City=u'Mobile', State=u'Alabama', State Code=u'AL', 2014 Population estimate=194675, 2015 median sales price=122.5),
```

## Using **display()** to view in tabular mode

```python
%python display(data)
```

▸ (2) Spark Jobs

| 2014 rank | City | State | State Code | 2014 Population estimate | 2015 median sales price |
|-----------|------------|---------|------------|--------------------------|-------------------------|
| 101 | Birmingham | Alabama | AL | 212247 | 162.9 |
| 125 | Huntsville | Alabama | AL | 188226 | 157.7 |
| 122 | Mobile | Alabama | AL | 194675 | 122.5 |

# Datasets

The Apache Spark Dataset API provides a type-safe, object-oriented programming interface

**DataFrame** is an alias for an untyped **Dataset [Row]**

Datasets provide compile-time type safety

The Dataset API also offers high-level domain-specific language operations

# Load Sample Data in Dataset

Read a data file from an external data source.

```scala
val df = spark.read.json("/databricks-datasets/samples/people/people.json")
```

At the time of reading the JSON file, Spark does not know the structure of your data.
It doesn't know how you want to organize your data into a typed-specific JVM object.
It attempts to infer the schema from the JSON file
This creates a `DataFrame` = `Dataset[Row]` of generic Row objects.

# Viewing Dataset

Viewing data in tabular mode using display()

```
// display the dataset table just read in from the JSON file display(ds)
```

Using standard Spark commands like take(), foreach() and println() API calls

```
// Using the standard Spark commands, take() and foreach(), print the first
// 10 rows of the Datasets.
ds.take(10).foreach(println(_))
```

# Learn More

Azure Databricks Overview - https://aka.ms/LearnAzureDatabricks

Delta Engine Documentation - https://aka.ms/DeltaEngineDocs

Next Gen Data Science Workspace - https://aka.ms/AzureDatabricks_DataScienceWorkspace

Spark and AI Summit 2020 content - https://aka.ms/SparkAISummit2020

Microsoft Azure

# Azure Analytics

| On-premises data |
| --- |

| Cloud data |
| --- |

| Devices data |
| --- |

| SaaS data |
| --- |

**Ingest**

Azure
Data Factory

**Prep**

Azure Databricks

**Model & Serve**

Azure Synapse
Analytics

**Store** — Azure Data Lake Storage

Power BI

# Modern Data Warehouse

**On-premises data**
Oracle, SQL, Teradata, fileshares, SAP

**Cloud data**
Azure, AWS, GCP

**SaaS data**
Salesforce, Dynamics

**INGEST**

Azure
Data Factory

**PREPARE**

Azure
Data Factory

Azure
Databricks

**TRANSFORM & ENRICH**

Azure
Data Factory

Azure
Databricks

**SERVE**

Azure
SQL Data
Warehouse

**VISUALIZE**

Power BI

**STORE**

Azure Data Lake Storage

# Azure Synapse Analytics - Data Lakehouse

# Azure Synapse Analytics

Integrated data platform for BI, AI and continuous intelligence

**Synapse Analytics**

Artificial Intelligence / Machine Learning / Internet of Things
Intelligent Apps / Business Intelligence

**Experience**

**Synapse Analytics Studio**

**Platform**

| MANAGEMENT |
| SECURITY |
| MONITORING |
| METASTORE |

**Languages**

| SQL | Python | .NET | Java | Scala | R |

**Form Factors**

| PROVISIONED | ON-DEMAND |

**Analytics Runtimes**

| SQL | Spark |

DATA INTEGRATION

**Azure**
**Data Lake Storage**

Common Data Model
Enterprise Security
Optimized for Analytics

Designed for analytics **workloads at any scale**

SaaS **developer experiences** for code free and code first

Multiple **languages** suited to different analytics workloads

Integrated analytics runtimes available provisioned and serverless on-demand

**SQL Analytics** offering T-SQL for batch, streaming and interactive processing

**Spark** for big data processing with Python, Scala, R and .NET

Integrated **platform services** for, management, security, monitoring, and metastore

Data **lake integrated** and Common Data Model aware

# Azure Synapse Analytics

Integrated data platform for BI, AI and continuous intelligence

**Artificial Intelligence / Machine Learning / Internet of Things**
**Intelligent Apps / Business Intelligence**

Designed for analytics **workloads at any scale**

## Synapse Analytics

| Experience | **Synapse Analytics Studio** |
|---|---|

SaaS **developer experiences** for code free and code first

| Platform | | |
|---|---|---|
| MANAGEMENT | **Languages** | |
| | SQL  Python  .NET  Java  Scala  R | |
| SECURITY | **Form Factors** | |
| | PROVISIONED | ON-DEMAND |
| MONITORING | **Analytics Runtimes** | |
| | SQL | Spark |
| METASTORE | DATA INTEGRATION | |

Multiple **languages** suited to different analytics workloads

Integrated analytics runtimes available provisioned and serverless on-demand

**SQL Analytics** offering T-SQL for batch, streaming and interactive processing

**Spark** for big data processing with Python, Scala, R and .NET

Integrated **platform services** for, management, security, monitoring, and metastore

## Azure
**Data Lake Storage**

**Common Data Model**
**Enterprise Security**
**Optimized for Analytics**

Data **lake integrated** and Common Data Model aware

# Synapse Studio

https://web.azuresynapse.net

# Synapse Studio

[https://web.azuresynapse.net](https://web.azuresynapse.net)

- **Data**: Shows the available data sources available to the workspace. These can exist internally in the workspace (such as a SQL compute database or a Spark database), or externally (such as a Data Lake Store Gen2, or Azure Blob Storage account)
- **Develop**: Shows the different objects used to query or operate with the data, such as SQL scripts, notebooks, data flows, Spark job definitions, Power BI, etc
- **Orchestrate**: Shows the objects used to automate analytics processes (such as pipelines, datasets, etc.)
- **Monitor**: Shows metrics for pipeline runs, trigger runs, integration runtimes, and spark applications
- **Manage**: Create linked services, pipeline triggers, integration runtimes, and manage access to Synapse

Microsoft

Azure Synapse Analytics
Spark

# Azure Synapse Apache Spark - Summary

- **Apache Spark 2.4** derivation
  - Linux Foundation Delta Lake 0.6.0 support
  - Apache Spark in Azure Synapse includes .NET Core 3.1
  - Python 3.6.1 + Anacondas support
- Operating System version
  - Apache Spark in Azure Synapse runs on Ubuntu 16.04.
- Tightly coupled to other Azure Synapse services
  - Integrated security and sign on
  - Integrated Metadata
  - Integrated and simplified provisioning
  - Integrated UX including nteract based notebooks
  - Fast load of SQL Analytics pools

- Core scenarios
  - Data Prep/Data Engineering/ETL
  - Machine Learning via Spark ML and Azure ML integration
  - Extensible through library management
- Efficient resource utilization
  - Fast Start
  - Auto scale (up and down)
  - Auto pause
  - Min cluster size of 3 nodes
  - Max cluster size 200 nodes
- Multi Language Support
  - .Net (C#), PySpark, Scala, Spark SQL, Java

https://docs.microsoft.com/en-us/azure/synapse-analytics/spark/apache-spark-version-support

# What is Delta Lake?

- OSS storage layer for Spark
- Provides:
  - ACID transactions
  - History of changed
  - Time travel in data history
  - Schema evolution
  - …

# Azure Synapse Analytics

Integrated data platform for BI, AI and continuous intelligence

**Artificial Intelligence / Machine Learning / Internet of Things**
**Intelligent Apps / Business Intelligence**

## Synapse Analytics

**Experience** — **Synapse Analytics Studio**

**Platform**

| MANAGEMENT | **Languages** | | | | | |
|---|---|---|---|---|---|---|
| | SQL | Python | .NET | Java | Scala | R |

**Form Factors**

| SECURITY | PROVISIONED | ON-DEMAND |
|---|---|---|

**Analytics Runtimes**

| MONITORING | SQL | Spark |
|---|---|---|

| METASTORE | DATA INTEGRATION | |
|---|---|---|

**Azure**
**Data Lake Storage**

Common Data Model
Enterprise Security
Optimized for Analytics

Designed for analytics **workloads at any scale**

SaaS **developer experiences** for code free and code first

Multiple **languages** suited to different analytics workloads

Integrated analytics runtimes available provisioned and serverless on-demand

**SQL Analytics** offering T-SQL for batch, streaming and interactive processing

**Spark** for big data processing with Python, Scala, R and .NET

Integrated **platform services** for, management, security, monitoring, and metastore

Data **lake integrated** and Common Data Model aware

# Languages

## Overview

Supports multiple languages to develop notebook

- pySpark (Python)
- Spark (Scala)
- SparkSQL
- .NET for Apache Spark (C#)

## Benefits

Allows to write multiple languages in one notebook

Using Magic command **%%** <language>

%%pyspark

%%spark

%%sql

%%csharp

use of **temporary tables** across **languages**

1. In Cell 1, read a DataFrame from SQL pool connector using Scala and create a temporary table.

Scala

```scala
%%scala
val scalaDataFrame = spark.read.option("format", "DW connector predefined type")
scalaDataFrame.registerTempTable( "mydataframetable" )
```

2. In Cell 2, query the data using Spark SQL.

SQL

```sql
%%sql
SELECT * FROM mydataframetable
```

3. In Cell 3, use the data in PySpark.

Python

```python
%%pyspark
myNewPythonDataFrame = spark.sql("SELECT * FROM mydataframetable")
```

# Create Notebook on files in storage

# Create Notebook on files in storage

# Languages – PySpark (Python)

# Languages – Spark (Scala)

# Languages – Spark SQL

# Languages – Spark.NET (C#)

# Library Management - Python

## Overview

Customers can add new python libraries at Spark pool level

## Benefits

Input requirements.txt in simple pip freeze format

Add new libraries to your cluster

Update versions of existing libraries on your cluster

Libraries will get installed for your Spark pool during cluster creation

Ability to specify different requirements file for different pools within the same workspace

## Constraints

The library version must exist on PyPI repository

Version downgrade of an existing library not allowed

### In the Portal

Specify the new requirements while creating Spark Pool in Additional Settings blade

# Library Management - Python

Microsoft

Azure Synapse Analytics
Foundation

# Manage – Access Control

## Overview

It provides access control management to workspace resources and artifacts for admin and users

## Benefits

Share workspace with the team

Increases productivity

Manage permissions on code artifacts and Spark pools

# Spark Monitoring

## Overview

Monitor Spark pools, Spark applications for the progress and status of activities

## Benefits

Monitor Spark pools for the status as paused, active, resume, scaling and upgrading

Build a dashboard to monitor performance

Track the usage of resources

# SQL Monitoring

## Overview

Monitor SQL Pool in Azure Portal for overall usage and query activities.

## Benefits

Access SQL Audit Logs for my SQL computes

Monitor status and progress of all/specific activities

Dashboard view to monitor performance

Get to know scale of SQL compute resource

# Azure Synapse Analytics

## Limitless analytics service with unmatched time to insight

**On-premises data**

**Cloud data**

**SaaS data**

### Unified platform and experience

**Synapse Studio**

| Integration | Management | Monitoring | Security |
|---|---|---|---|

### Analytics Runtimes

| SQL | Spark |
|---|---|

**Azure Data Lake Storage**

**Azure Machine Learning**

**Power BI**

# Azure Machine Learning Services

## Overview

Data Scientists can use Azure ML notebooks to do (distributed) data preparation on Synapse Spark compute.

## Benefits

Connect to your existing Azure ML workspace and project

Use the AutoML Classifier for classification or regression problem

Train the model

Access open datasets

API Azure Cognitive Services

# Azure Machine Learning Services (continued)

**Configure AutoML and Train the Models**

Cell 9

```
1  l_config = AutoMLConfig(task = 'regression',debug_log = 'automl_errors.log',
2                          primary_metric = 'normalized_root_mean_squared_error', iteration_timeout_minutes = 10,
3                          iterations = 2, preprocess = True, n_cross_validations = 2,max_concurrent_iterations = 2,
4                          verbosity = logging.INFO,spark_context=sc, enable_onnx_compatible_models=True, cache_store=Tru
```

Cell 10

```
1  local_run = experiment.submit(automl_config, show_output = True)
```

**Best Model**

Cell 12

```
1  best_run, fitted_model = local_run.get_output(return_onnx_model=True)
2  print(fitted_model)
```

**Portal URL for Monitoring Runs**

Cell 14

```
1  more Insights of experiment
2  displayHTML("<a href={} target='_blank'>Your experiment in Azure Portal: {}</a>".format(local_run.get_portal_url(), local_r
```

# Spark to Cosmos DB Connector

## Overview

**Spark to Cosmos DB Connector**

## Benefits

1. Connection is made between Spark master node and Cosmos DB gateway node.

2. Partition map data is transmitted back to Spark master node.

3. Query is submitted from Spark worker nodes to

4. Cosmos DB data nodes and the data is transmitted back to Spark worker nodes for further processing

# Azure Synapse Link for Cosmos DB

⬆ Publish all ❷   ✓ Validate all   ↻ Refresh   🗑 Discard all

**Data**   + ⌄ «

⬜ Notebook 2 ●   ⬜ Notebook 3 ●   ⋯

| Workspace | Linked |
| --- | --- |

+ Cell ⌄   ▷ Run all   ↺ Undo | ⌄   ⬆ Publish | ⋯   🦑 ⚙ ⋯

🔍 Filter resources by name

Cell 1

▷ Storage accounts   4

▷ Cosmos DB   2

▷ Datasets   64

```
[4]    1    # Load a streaming Spark DataFrame from a Cosmos DB container
       2    # To select a preferred list of regions in a multi-region Cosmos
       3
       4    dfStream = spark.readStream\
       5        .format("cosmos.oltp")\
       6        .option("spark.synapse.linkedService", "manufacturing")\
       7        .option("spark.cosmos.container", "mfg-quality")\
       8        .option("spark.cosmos.changeFeed.readEnabled", "true")\
       9        .option("spark.cosmos.changeFeed.startFromTheBeginning", "tr
      10        .option("spark.cosmos.changeFeed.checkpointLocation", "/loca
      11        .option("spark.cosmos.changeFeed.queryName", "streamQuery")\
      12        .load()
```

Command executed in 7mins 27s 621ms by odl_user_209652 on 07-30-2020 17:53:03.362 -04:00

> **Job execution** Succeeded   **Spark** 2 executors 8 cores   View in monitoring   Open Spark UI ☒

⬆ Publish all **2**   ✓ Validate all   ↻ Refresh   🗑 Discard all

**Data**   + ⤓ «

Workspace | **Linked**

🔍 Filter resources by name

▷ **Storage accounts**                    4

◢ **Cosmos DB**                          2

▷ ⬡ CosmosDb1 (industrial)

◢ ⬡ manufacturing (manufacturing-data)

  ▷ 📊 manufacturing

  ◢ 📊 mfg-quality                      ···

▷ **Datasets**                           64

*CosmosDB
Analytics Store*

▭ Notebook 2  ●    ▭ Notebook 3  ●                    ···

+ Cell ⌄   ▷ Run all   ↶ Undo ⌄   ⬆ Publish   ···

```
5      .format( cosmos.olap )\
6      .option("spark.synapse.linkedService", "manufacturing"
7      .option("spark.cosmos.container", "mfg-quality")\
8      .load()
9
10   display(df.limit(10))
```

Command executed in 2mins 55s 537ms by odl_user_209652 on 07-30-2020 17:48:58.727 -04:00

› **Job execution** Succeeded   **Spark** 2 executors 8 cores   View in monitoring   Open Spark UI ⧉

📋 ⬆ ⬒

View   | Table | **Chart** |                                    ⚙

Chart type

pie chart                                    ⌄

Key

_rid                                         ⌄

Values

_ts                                          ⌄

Series Group

dkhAO1ea...        ddkhAO1ea...

✓ Ready (Stop session) | Configure session

# Power BI

## Overview

Power BI is a business analytics service that delivers insights to enable fast, informed decisions

## Benefits

Create Power BI reports in the workspace

Have access to published reports in workspace

Update reports real time from Synapse workspace to get it reflected on Power BI service

Visually explore and analyze data

# Power BI Aggregations and Synapse query performance



**Power BI layers:**
- In Memory
- Dual Table
- DirectQuery

In-Memory storage engine for millisecond latency analytics over aggregated data

Bridge between In-Memory storage engine and underlying raw data

SQL engine to generate queries for non-cache data and pass through to Synapse

**Synapse layers:**
- Resultset Cache
- In-Memory
- SSD Adaptive Cache
- Materialized Views
- IO Optimized Data Access

Compute Pool isolated result cache with resilience to cluster elasticity and response time ~ 200ms

In-Memory cache for sub-second response times

NVMe based SSD cache that acts an extension of In-Memory cache for fast, localized data caching

Pre-Joined + Pre-Aggregated data with guaranteed transactional consistency and automatic query optimizer matching

Analytics optimized data structures on disk including ordered columnar, partitioning, and nonclustered indexing

# Power BI visualization end to end integration with Synapse Analytics

# Augment analytics with Azure Machine Learning
## (coming)

**Data Engineers**

**Azure Synapse**

**Azure ML**

**Data Scientists**

**Core Capabilities**

| Data Warehousing and Data Prep | Build, train and deploy models |
|---|---|

**New Capabilities**

| Discover & Deploy Models to Enrich Data SQL & Spark | Basic training and batch scoring in Spark | Publish Synapse data for ML | Prepare data for ML using Spark |
|---|---|---|---|

**Shared assets**

Shared model registry and data assets

Shared Spark compute pools

Linked Workspaces for asset sharing

Common RBAC and security model

**Collaboration**

# Synapse Spark in Azure ML private preview experience – interactively run data processing on Spark in notebook

## Step2: data exploration and transformation

```
%%synapse pyspark

# Drop columns that are not relevant to ML modeling
columns_to_drop = ['vendorID','pickupLongitude','pickupLatitude','dropoffLongitude','dropoffLatitude','lpepPickupDatetime','l
df = nyc_green_df.drop(*columns_to_drop)

# Transform column tripType
df_t = df.withColumn('tripType', when(df.tripType==2,lit('0')).otherwise(df.tripType))

# Create or replace temp view to prepare for pyspark sql
df_t.createOrReplaceTempView("df_temp")

# Run query by leveraging pyspark sql
sqlDF = spark.sql("""
    SELECT *
    FROM df_temp
    WHERE  (tripDistance>=25 and tripDistance<50)
    AND (passengerCount>0 and totalAmount>0)
""")

# Data exploration and transformation is completed. Print processed data sample.
print("Reading for machine learning")
sqlDF.show(10)
```

## Step1: Start Spark session

```
%synapse start -c synapseMedium
```

```
%%synapse pyspark

import numpy as np
import pyspark
import os
import urllib
import sys
from datetime import datetime
from datetime import datetime
from dateutil import parser
from pyspark.sql.functions import *
from pyspark.ml.classification import *
from pyspark.ml.evaluation import *
from pyspark.ml.feature import *
from pyspark.sql.types import StructType, StructField
from pyspark.sql.types import DoubleType, IntegerType, StringType


# print runtime versions
print('****************')
print('Python version: {}'.format(sys.version))
print('Spark version: {}'.format(spark.version))
print('****************')


# start Spark session
spark = pyspark.sql.SparkSession.builder.appName('NYCGreenTaxi')\
    .config("spark.jars.packages", "io.delta:delta-core_2.12:0.7.0") \
    .config("spark.sql.extensions", "io.delta.sql.DeltaSparkSessionExtension") \
    .config("spark.sql.catalog.spark_catalog", "org.apache.spark.sql.delta.catalog.DeltaCatalog") \
    .getOrCreate()


****************
Python version: 3.6.1 |Continuum Analytics, Inc.| (default, May 11 2017, 13:09:58)
[GCC 4.4.7 20120313 (Red Hat 4.4.7-1)]
Spark version: 2.4.4.2.6.99.201-15911041
****************
```

| Job ID | Job Name | Status | Stages | Tasks | Submission Time | Duration |
|---|---|---|---|---|---|---|
| 2 | showString | COMPLETED | 1/1 | 1/1 | a few seconds ago | 1s |
| 3 | showString | COMPLETED | 1/1 | 4/4 | a few seconds ago | 1s |

Jobs: 2 COMPLETED  Spark: 3 EXECUTORS 8 CORES

```
Reading for machine learning
+-------------+------------+-----------+--------+
|passengerCount|tripDistance|totalAmount|tripType|
+-------------+------------+-----------+--------+
|            1|       28.13|      85.06|       1|
|            1|       32.09|       80.0|       0|
|            1|       25.96|      83.72|       1|
|            1|       30.93|      90.06|       1|
|            1|       28.54|       82.8|       1|
|            1|       25.94|     116.76|       1|
|            1|       29.95|      87.06|       1|
|            1|       25.13|       76.8|       1|
|            1|       30.42|      120.8|       1|
```

### step3: Stop Spark session

When current session reach the status timeout, dead or any failure, you must explicitly stop it before start new one.

```
%synapse stop
```

Session stopped.

**Learn More**
- What is SQL on-demand?: link
- What is Apache Spark in Azure Synapse Analytics?: link
- Best practices for SQL pool in Azure Synapse Analytics: link
- Best practices for SQL on-demand in Azure Synapse Analytics: link
- Azure Synapse Analytics shared metadata: link
- Use maintenance schedules to manage service updates and maintenance: link
- Cheat sheet for Azure Synapse Analytics (formerly SQL DW): link
- Best practices for SQL Analytics in Azure Synapse Analytics (formerly SQL DW): link
- Synapse Analytics documentation is here: aka.ms/SynapseDocs

# Q&A

**Adrián J. Fernández Zenteno**

*Sr. Cloud Solution Architect - Azure Data & AI*

*[Adrian.Fernandez@microsoft.com](mailto:Adrian.Fernandez@microsoft.com)*