# TECHORAMA

https://github.com/SiWhiteley/DatabricksETL

# Agenda

| What is Databricks? | Patterns & Implementation | Orchestration | The Big Picture |

A HISTORY OF SPARK

Google

Google File System Papers
Released

**2003**

# HDFS IN AZURE



Blob Storage

File Extents

Google File System Papers Released

**2003**

Google MapReduce Papers

**2004**

Google File System Papers Released
**2003**

Google MapReduce Papers
**2004**

**2006**
Apache Hadoop project created

Matei Zaharia starts Spark project
**2012**

Project donated to Apache Foundation
**2013**

Databricks founded by Matei
**2013**

PYTHON PIPELINE PRIMER

# Databricks is...

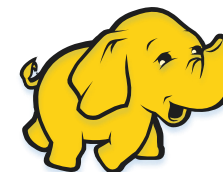*__Apache Spark__, built by the guys who wrote Spark, made **super easy***

**2016**

Microsoft
Azure

It's new to
Azure, not to
everyone else!

**2018**

PYTHON PIPELINE PRIMER

# HDInsight

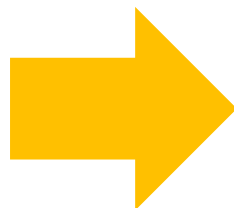| | | | |
|---|---|---|---|
| | **Zeppelin** | | |
| **Ranger** | **Pegasus** | **RHadoop** | **Mahout** |
| **Flume** | **HCatalog** | | |
| **Oozie** | **Hive** | **HBase** | **Storm** **Pig** |
| | **Spark** **MapReduce** | | |
| | **Ambari** | | |
| | **YARN** | | |

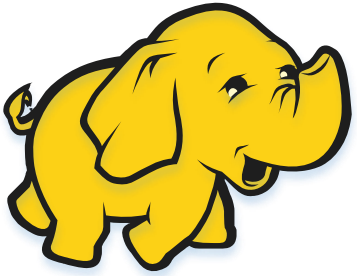| Open Source | Open Source | Proprietary |
|---|---|---|
| 20 min provisioning | 5 min provisioning | 1 min provisioning |
| Integrates Well | Integrates Well | Integrates Poorly |
| Secure | Secure | Secure |
| Hadoop, Spark, Kafka, Hbase, HIVE, Storm… | Spark (Python/Scala/R) | U-SQL |
| Slow Release Cycle | Fast Release Cycle | Slow Release Cycle |

PYTHON PIPELINE PRIMER

1 Scale

2 Flexibility

Microsoft® SQL Server®
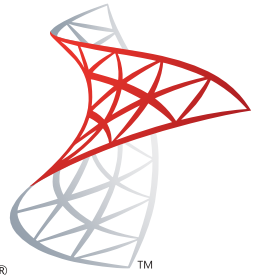
PYTHON PIPELINE PRIMER

WHAT PROBLEMS DOES IT SOLVE?

Language Flexibility

Scale

Integrations

File Types

PYTHON PIPELINE PRIMER

# MACHINE LEARNING AT SCALE

# GEOSPATIAL MAPPING



Geospatial
Shape Files

# IMAGE PROCESSING

# BUT MOST OF ALL...



Enterprise Data Warehouse

Source Systems

Stage

Clean

Warehouse

PYTHON PIPELINE PRIMER

Data Lake

Enterprise Data Warehouse

Cleaning

Transformation

**RAW**

**BASE**

**ENRICHED**

Source
Systems

Landing

Warehouse

PYTHON PIPELINE PRIMER

# UNDER THE HOOD

DataFrame API

SQL API

Resilient Distributed Datasets – In-Memory Data Blocks

Core Spark Engine – 80% Scala Code Libraries

PYTHON PIPELINE PRIMER

# UNDER THE HOOD

# LANGUAGE OPTIONS

# THE CATALYST OPTIMISER

# DEMO:

## AZURE DATABRICKS

- Databricks Workspace

- Clusters

PATTERNS & IMPLEMENTATION

Workload Isolation

Processing Cluster

Streaming Cluster

Interactive Cluster

PYTHON PIPELINE PRIMER

Workload Isolation

Fact Cluster

Interactive Cluster

Dim Cluster

PYTHON PIPELINE PRIMER

Workload Isolation

Fact Cluster

Interactive Cluster

Dim Cluster

(High Concurrency)

PYTHON PIPELINE PRIMER

EXECUTIONS

# DISTRIBUTED COMPUTE



PYTHON PIPELINE PRIMER

# DISTRIBUTED COMPUTE



Spark Driver

Spark Context

Spark Executor

Spark Executor

Spark Executor

Spark Executor

Worker    Worker    Worker    Worker

ADLS

*.csv

Extent 1    Extent 2    Extent 3    Extent 4

# DISTRIBUTED COMPUTE



**Spark Driver**

**Spark Context**

**Spark Executor**

**Spark Executor**

**Spark Executor**

**Spark Executor**

Worker

Worker

Worker

Worker

ADLS

*.csv

| Extent 1 | Extent 2 | Extent 3 | Extent 4 |

# DISTRIBUTED COMPUTE



Spark Driver

Spark Context

Spark Executor — Worker

Spark Executor — Worker

Spark Executor — Worker

Spark Executor — Worker

ADLS

*.csv

| Extent 1 | Extent 2 | Extent 3 | Extent 4 |

# SO HOW DO WE USE IT?

# THE DATA FRAME

**DataFrame**

- Schema ← Parameter
- Format ← Parameter
- Location ← Parameter

```python
df = (spark
        .read
        .schema(newSchema)
        .format(fileFormat)
        .load(dataLocation)
     )
```

# SCHEMA ON READ – INFER SCHEMA

```
Cmd 3

1   df = sqlContext.read.format("csv") \
2     .option("header", "true") \
3     .option("inferSchema", "true") \
4     .load("abfss://root@dblake.dfs.core.windows.net/RAW/Public/Taxi/v1/SmallSlice.csv")
```

```
▼ ▦  df:  pyspark.sql.dataframe.DataFrame

        Dispatching_base_num: string
        Pickup_DateTime: timestamp
        DropOff_datetime: string
        PUlocationID: integer
        DOlocationID: string
```

# SCHEMA ON READ – INFER SCHEMA

Infer Schema &
Read Files

Samples Every File!

# SCHEMA ON READ – SAMPLE FILES



Infer Schema from Sample file

Apply Schema & Read Files

Inferring schema from a small file sample before reading large datasets?

PYTHON PIPELINE PRIMER

# SCHEMA ON READ – METADATA STORE



**Metadata DB**

Retrieve Schema from Metadata Store

Apply Schema & Read Files

Files such as Parquet don't need a schema to be supplied!

PYTHON PIPELINE PRIMER

# READING FILES – NO PARTITIONS

*SELECT * FROM MyFiles WHERE Year = 2019 AND Month = 3*



SQL Query Action

Sales

Sales1.csv

Executors

Sales2.csv

Sales3.csv

Sales4.csv

Filtering performed on executors
after reading all files

PYTHON PIPELINE PRIMER

# READING FILES – PARTITIONED

*SELECT * FROM MyFiles WHERE Year = 2019 AND Month = 3*



Filtering performed by selectively reading files

SQL Query Action

Executors

Sales

Month=1 — Sales1.csv

Month=2 — Sales2.csv

Month=3 — Sales3.csv

Month=4 — Sales4.csv

# ORCHESTRATION

# AZURE DATA FACTORY



Notebook

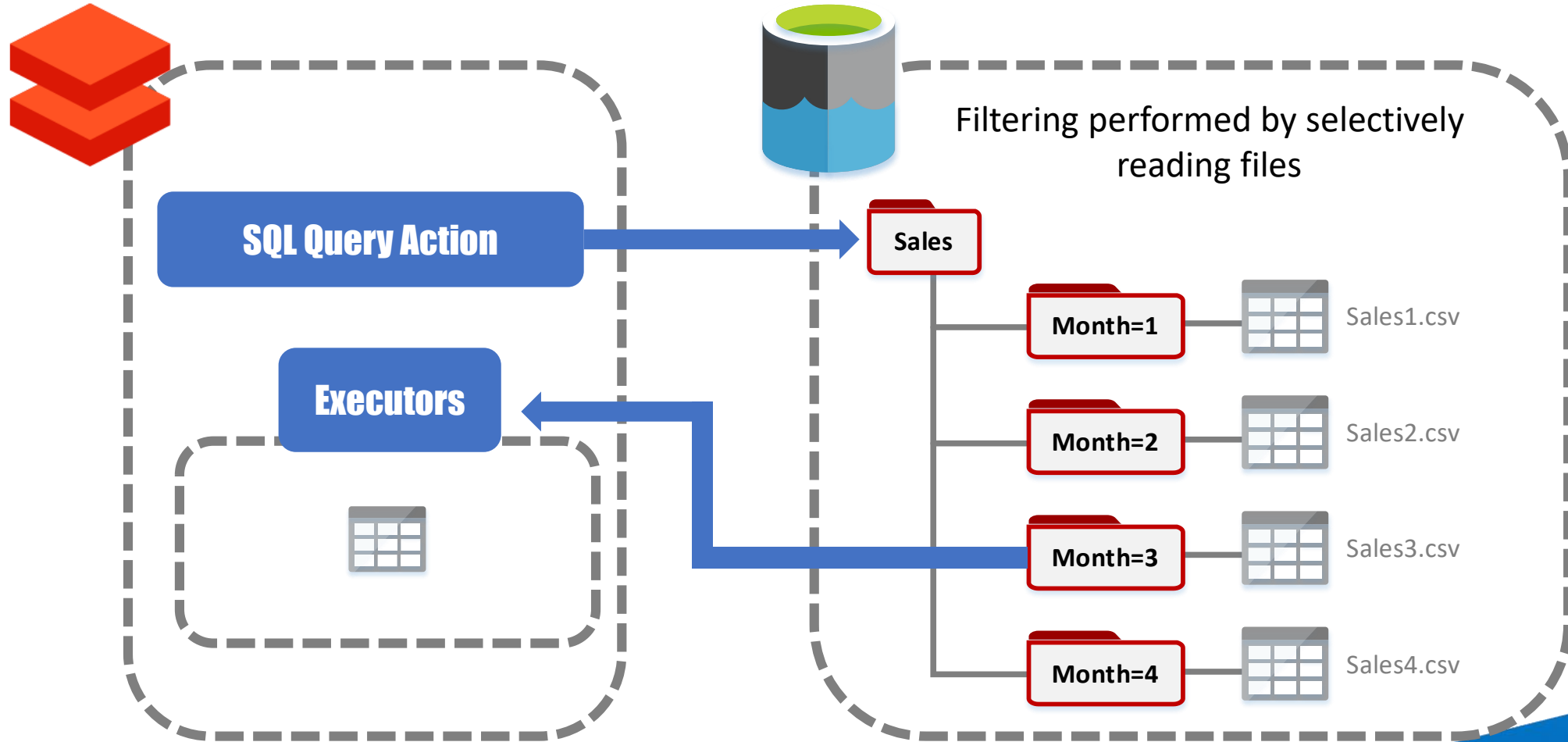Python Script

Jar File

# AZURE DATA FACTORY

Notebook path * `/Demonstrations/Dynamic Transformation/Dy` [Browse]

▲ Base Parameters

+ New | 🗑 Delete

| NAME | VALUE |
|---|---|
| entity_name | Taxi |

Entity Name : `Taxi` ⌄

### Cmd 1

## Configure Widgets

```
1  #dbutils.widgets.removeAll()
2  dbutils.widgets.dropdown("entity_name", "Taxi",["Taxi","TaxiZones"] ,"Entity Name")
```

```
"runOutput": {
    "TransformationsApplied": 2,
    "Status": "Succeeded",
    "ProcessedRows": 66141344
},
```
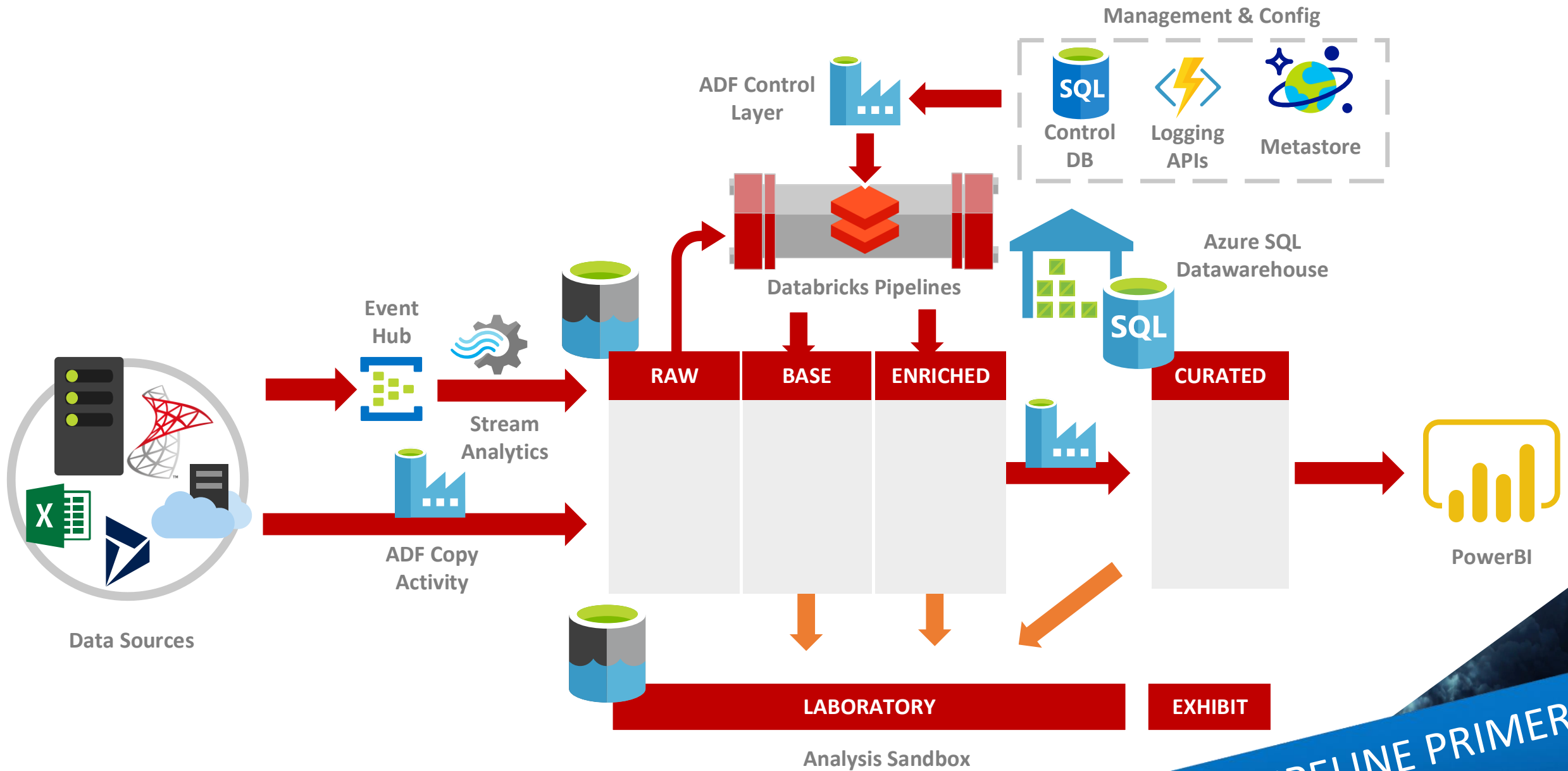
### Cmd 10

## Return Results to Caller

```
1  import json
2  dbutils.notebook.exit(json.dumps({"processedRows":processedRows, "status":"Succeeded",
```

Notebook exited: {"Status": "Succeeded", "ProcessedRows": 66141344, "TransformationsApplie

PYTHON PIPELINE

# THE BIG PICTURE

SIMON WHITELEY

@MRSIWHITELEY

WWW.ADVANCINGANALYTICS.CO.UK

Questions?

TECHORAMA