



ADVANCING
ANALYTICS

ADVANCING DATABRICKS

NEXT LEVEL ETL



Simon Whiteley
@MrSiWhiteley



DATA:Scotland



Microsoft



Arnold Clark



SentryOne®



<https://github.com/SiWhiteley/AdvancingDatabricks>

**ADVANCING
DATABRICKS**

Agenda

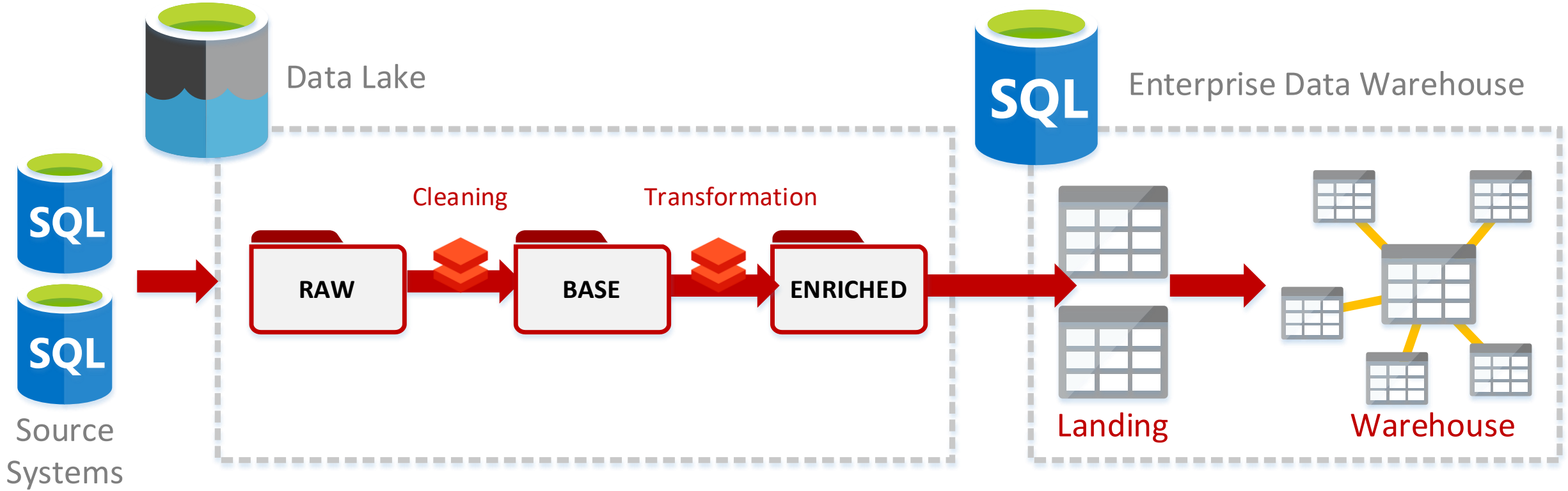
Schemas & Inference

Validating Data

Orchestration

Spark Internals / Deep Questions

**ADVANCING
DATABRICKS**



**ADVANCING
DATABRICKS**

THE DATA FRAME

DataFrame

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

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





APPLYING STRUCTURE

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

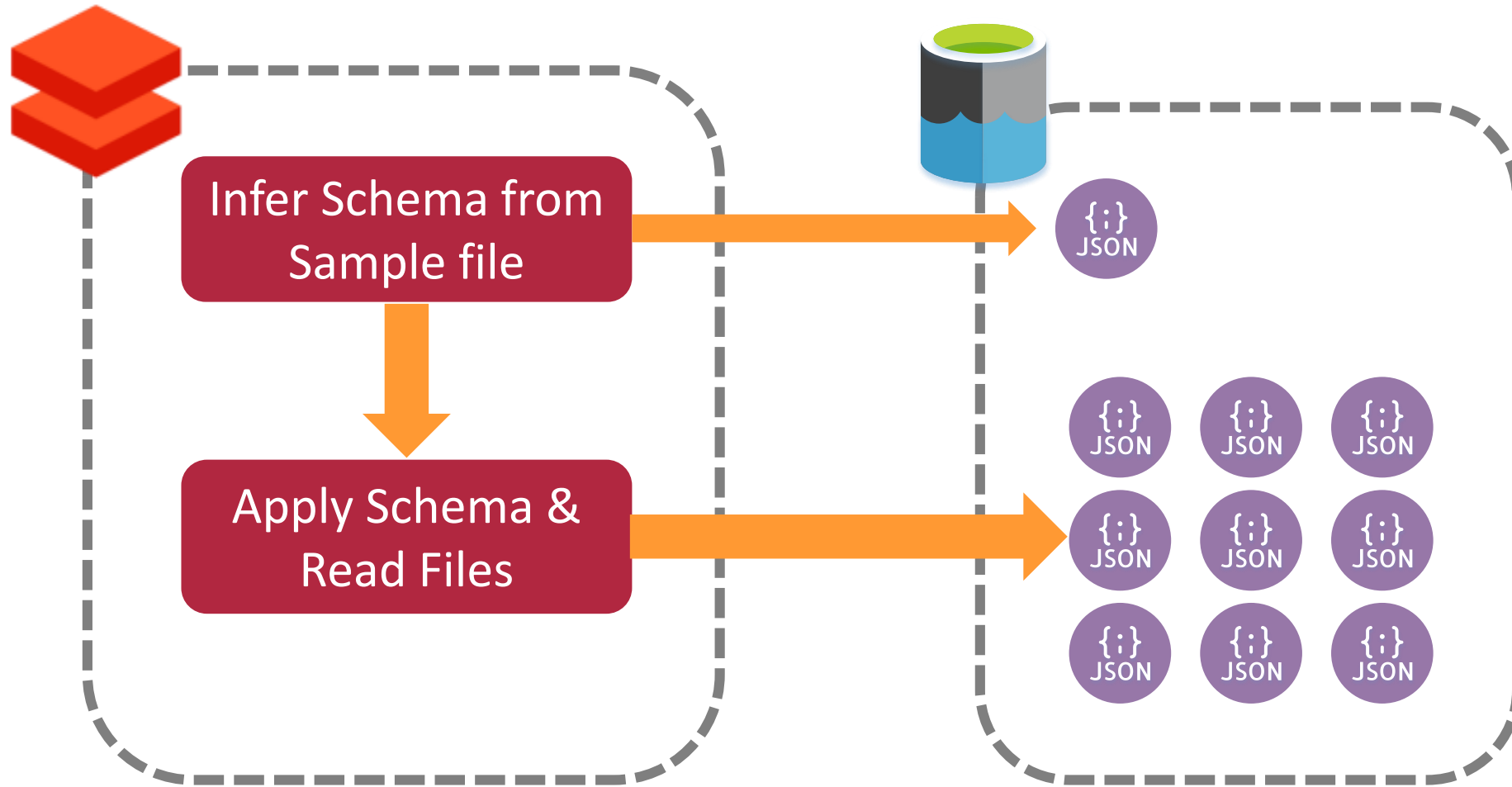
- Dispatching_base_num: string
- Pickup_DateTime: timestamp
- DropOff_datetime: string
- PUlocationID: integer
- DOlocationID: string

SCHEMA ON READ – INFER SCHEMA

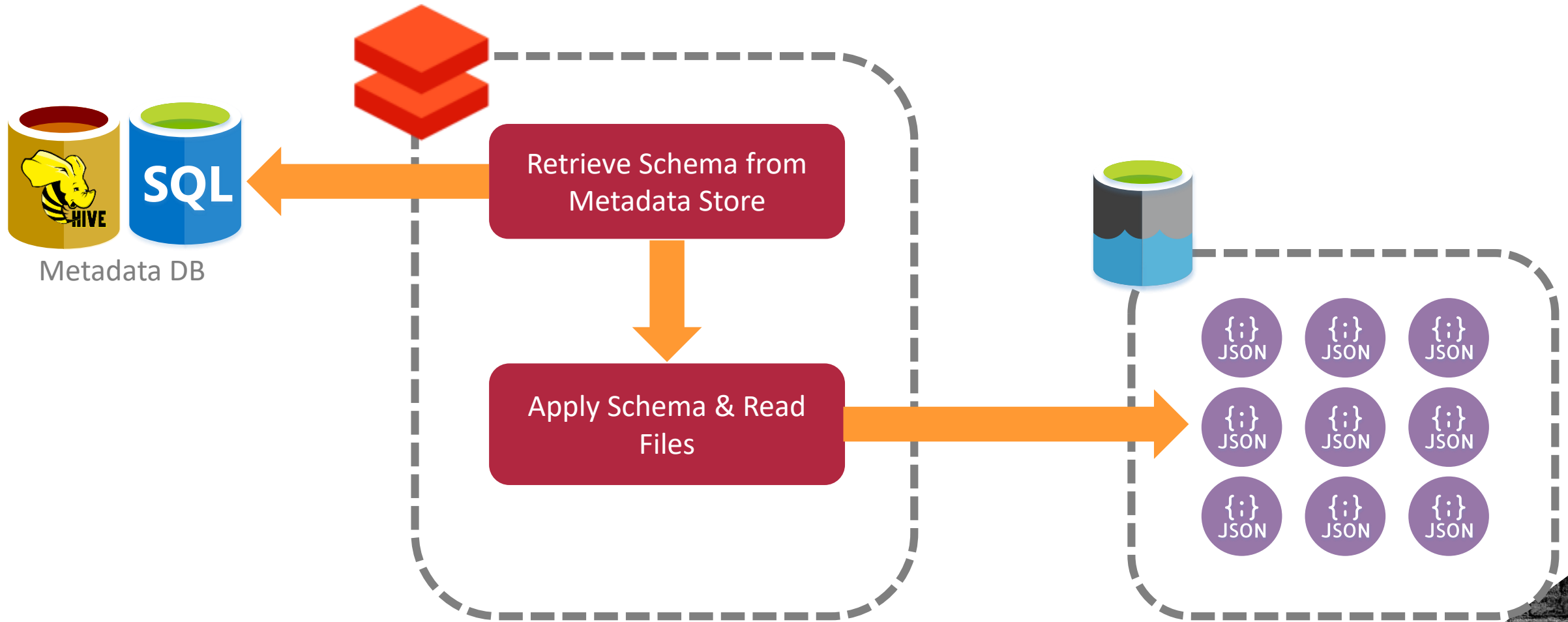


**ADVANCING
DATABRICKS**

SCHEMA ON READ – SAMPLE FILES



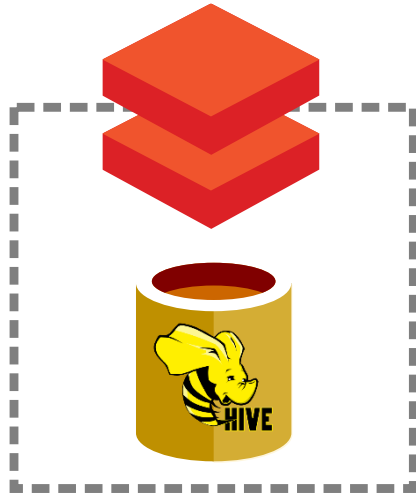
SCHEMA ON READ – METADATA STORE



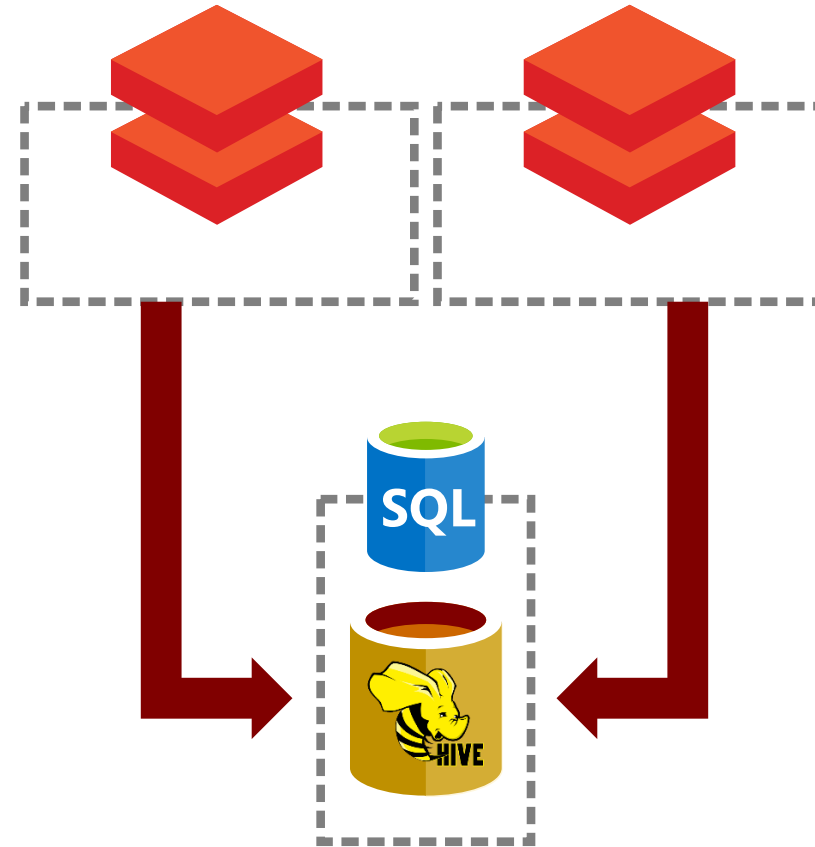
**ADVANCING
DATABRICKS**

USING THE HIVE METASTORE

Internal
(Cross-Cluster)



External
(Cross-Workspace)



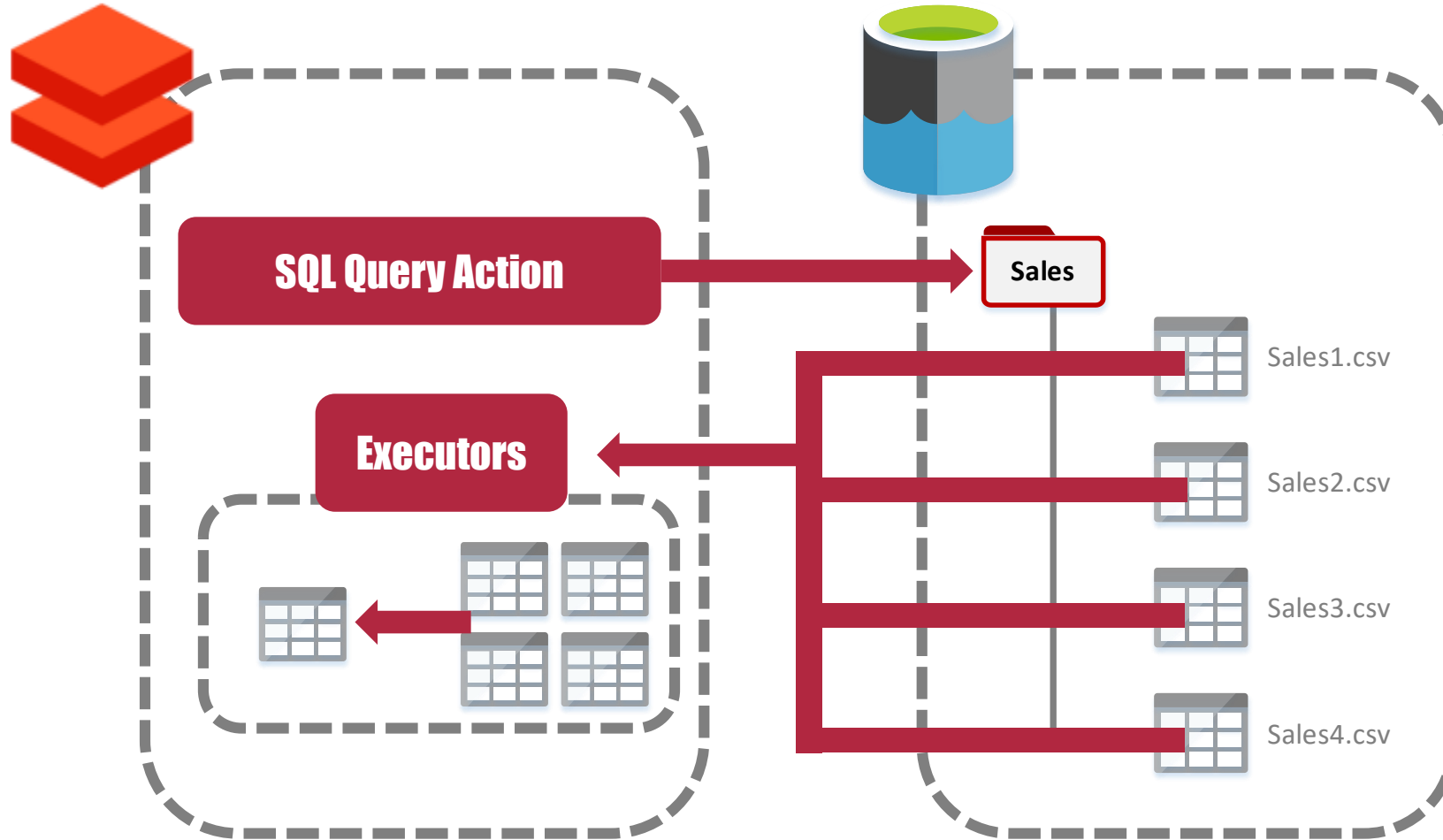
ADVANCING
DATABRICKS



PARTITIONING

READING FILES – NO PARTITIONS

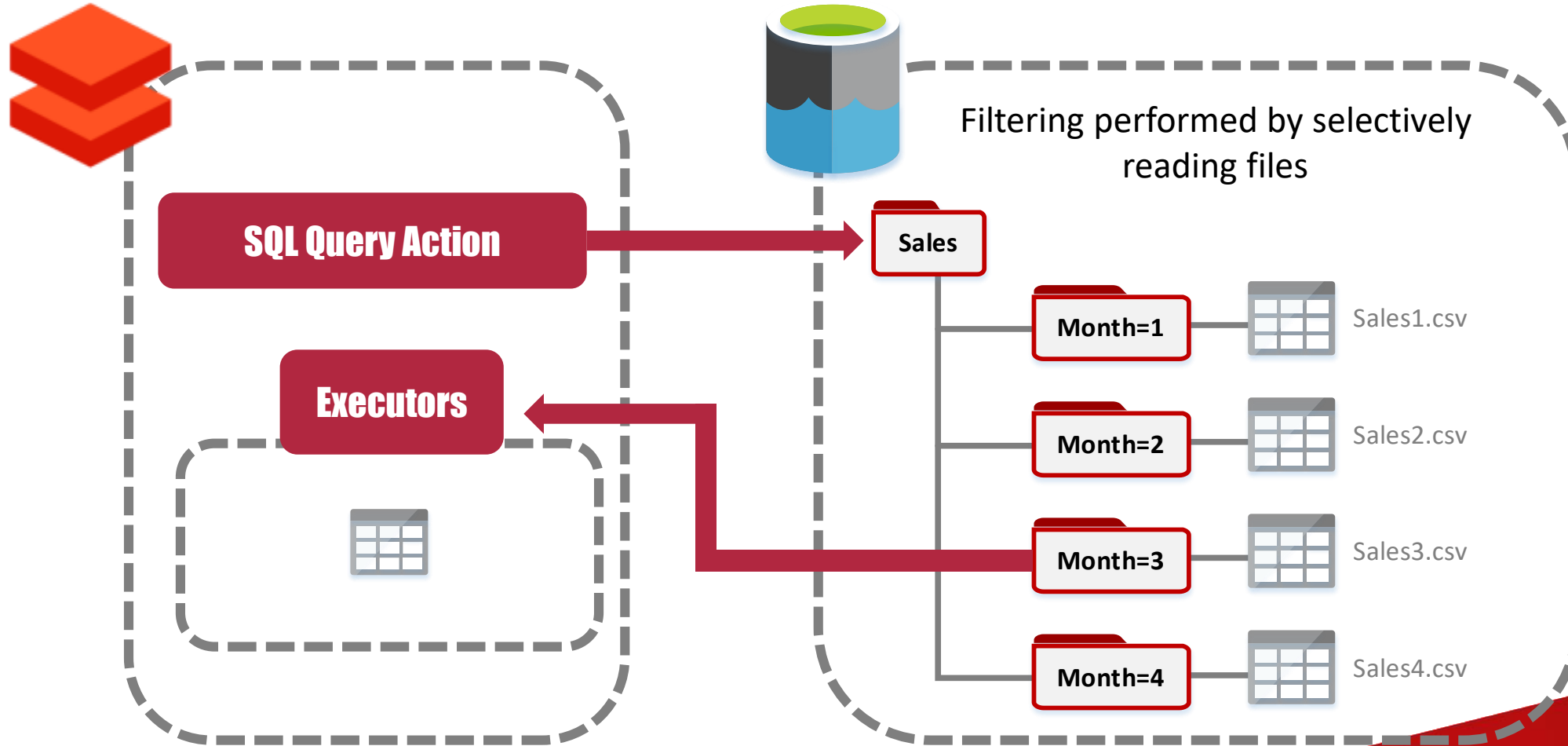
*SELECT * FROM MyFiles WHERE Year = 2019 AND Month = 3*



Filtering performed on executors
after reading all files

READING FILES – PARTITIONED

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



DATA ENGINEERING VS DATA SCIENCE



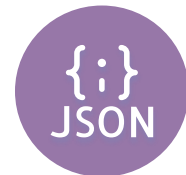
DATA VALIDATION

DATA VALIDATION

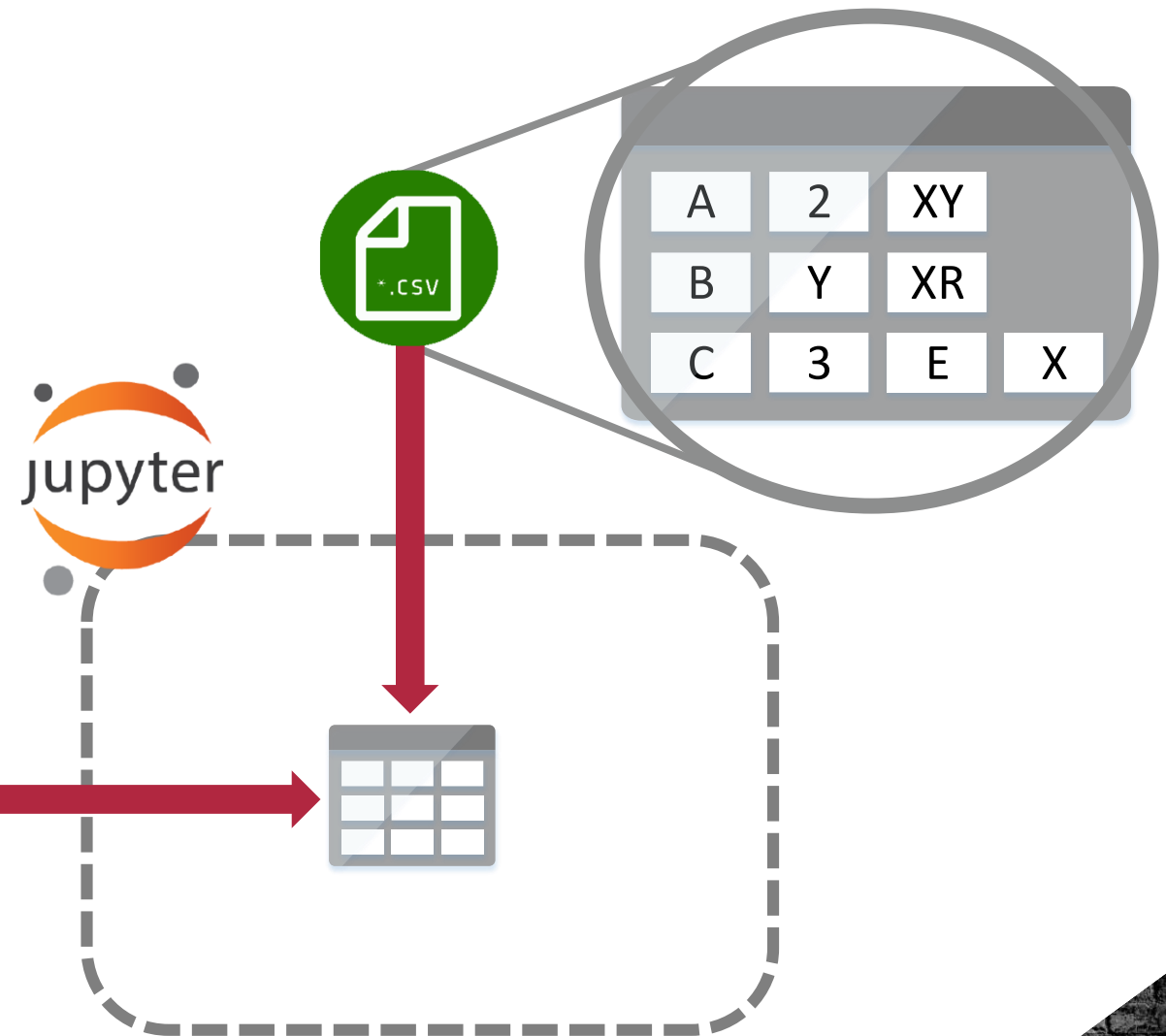
With Schema-On-Read file types such as CSV, we need to make sure that the data is in the format that we expect it to be.

Let's assume that we're specifying the schema ourselves, from an external lookup

Schema



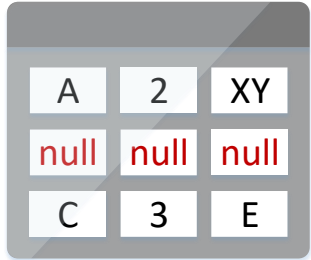
```
{"fields": [
  {"name": "Col1", "nullable": true, "type": "string"},
  {"name": "Col2", "nullable": true, "type": "integer"},
  {"name": "Col3", "nullable": true, "type": "string"}
]}
```



DATA ENGINEERING VS DATA SCIENCE

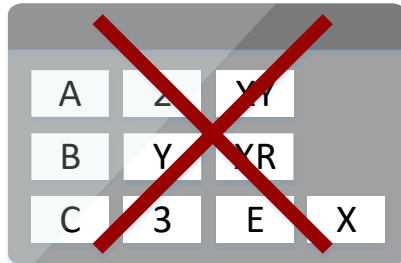
DATA VALIDATION

We have three different methods for handling failed parsing of a DataFrame when accessing text datasets such as csv:



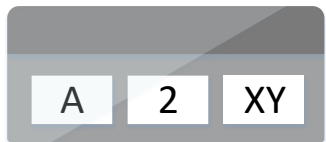
A	2	XY
null	null	null
C	3	E

- 1 **PERMISSIVE** – Extra columns are simply ignored, if any column fails to parse the entire row is nullified



A	2	XY	
B	Y	XR	
C	3	E	X

- 2 **FAILFAST** – If any attributes are different to the specified schema, whether by failing to parse datatypes or attributes added/missing, the entire DataFrame will fail to load



A	2	XY
---	---	----

- 3 **DROPMALFORMED** – Any rows that differ from the schema will be silently dropped from the dataset, also known as the “Nothing to see here” approach to ETL...

string	int	string	
A	2	XY	
B	Y	XR	
C	3	E	X

DEMO:

DATA VALIDATION APPROACHES

- PERMISSIVE
- FAILFAST
- DROP MALFORMED
- A BETTER WAY



ORCHESTRATION

DATABRICKS WIDGETS

Cmd 3

```
1 #dbutils.widgets.removeAll()
2 dbutils.widgets.text("fileName", "Product","AdventureWorks Table")
3 dbutils.widgets.dropdown("entity_name", "Taxi",["Taxi","TaxiZones"] ,"Entity Name")
```

AdventureWorks Table :

Entity Name :



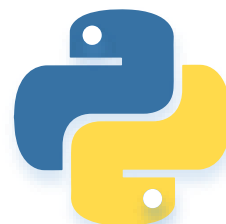
Cmd 8

```
1 fileName = dbutils.widgets.get("fileName")
```

ADVANCING
DATABRICKS



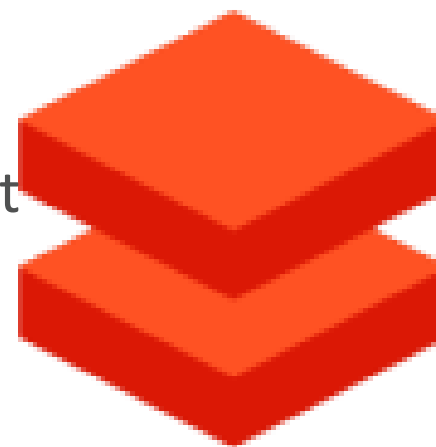
Notebook



Python Script

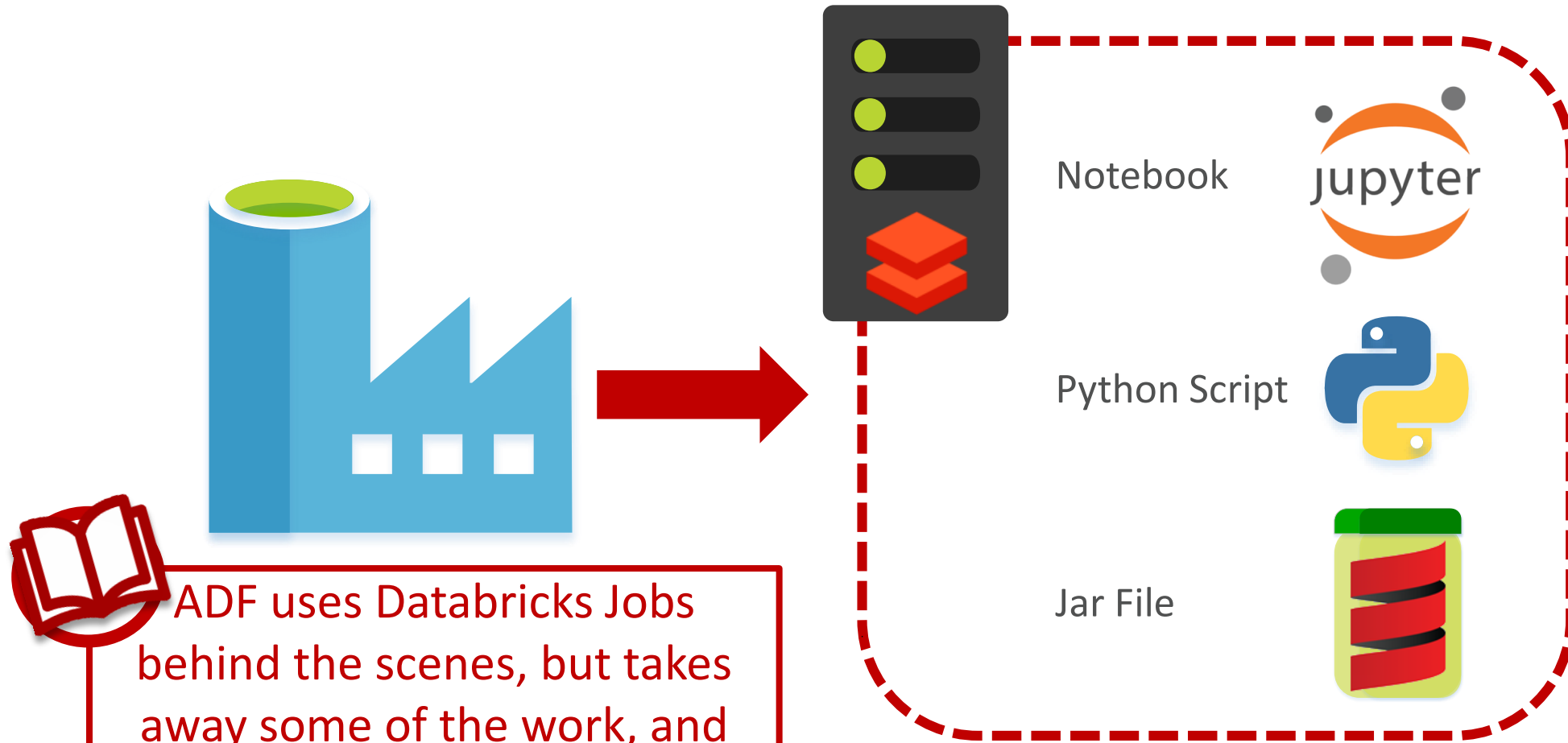


Jar File



ADVANCING
DATABRICKS

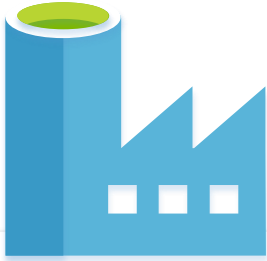
AZURE DATA FACTORY



ADF uses Databricks Jobs behind the scenes, but takes away some of the work, and means you can orchestrate your notebooks with other tasks!

**ADVANCING
DATABRICKS**

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",
```

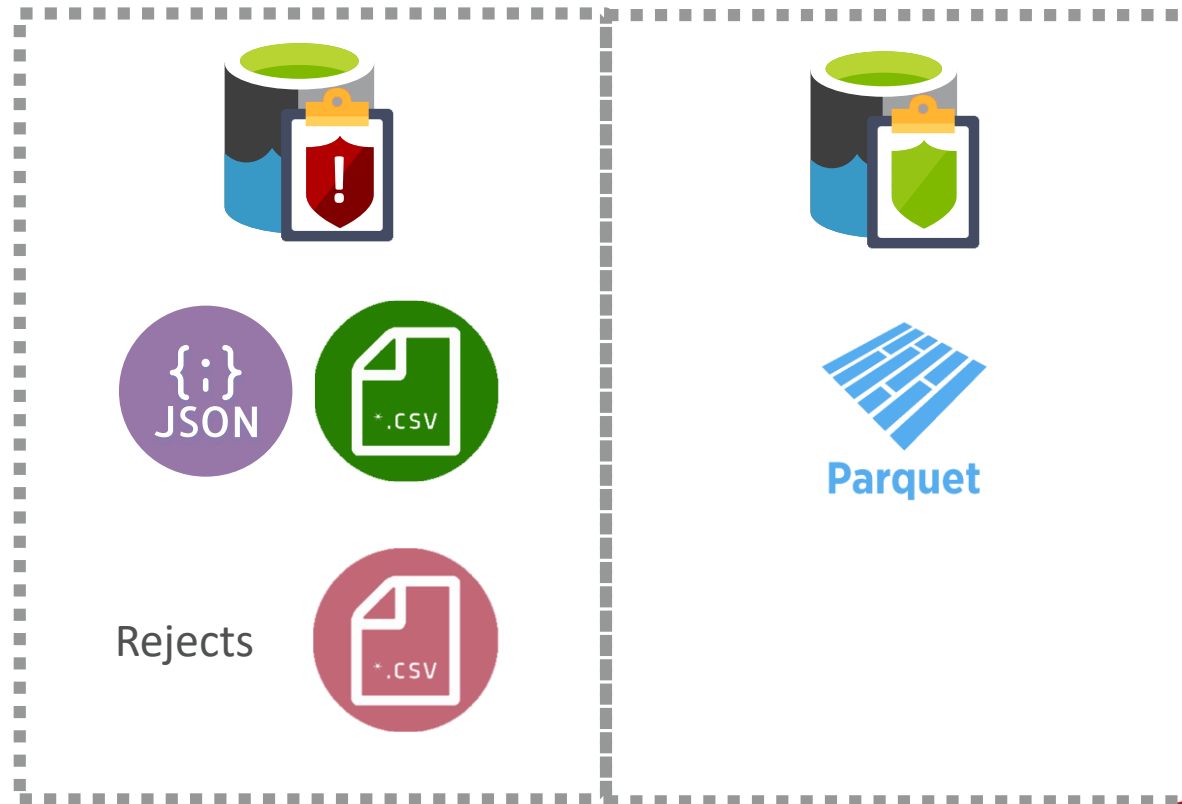
Notebook exited: {"Status": "Succeeded", "ProcessedRows": 66141344, "TransformationsApplie

ADVANCED
DATABRICKS



ONTO THE DEMO...

OUR GOAL:



**ADVANCING
DATABRICKS**

DEMO:

DYNAMIC VALIDATION

- USING WIDGETS
- PARSING SCHEMAS
- DYNAMIC DATAFRAMES
- WRITING DATA
- OUTPUT PARAMETERS

WHY WOULD WE DO THIS?

TIME
(COST)



Old Days

Design

Build SSIS

Build Orchestration

Test Everything

Deploy Everything

Support Everything

Newer Days

Design

Build SSIS

Build Orchestration

DevOps

Support Everything

Now

Design Metadata

Support One Thing

**ADVANCING
DATABRICKS**

Session Feedback
bit.ly/2019sfeedback



Event Feedback
bit.ly/2019efeedback

DATA:Scotland

DATA:Scotland Speaker Scholarship

[BIT.LY/DS-SPEAKER](https://bit.ly/ds-speaker)

DATA:Scotland



**ADVANCING
ANALYTICS**

QUESTIONS?



Simon Whiteley
@MrSiWhiteley