

# ADVANCING DATABRICKS



Understanding Spark to build next-gen ETL systems

Microsoft  
Partner



Gold Data Analytics  
Gold Data Platform  
Silver Cloud Platform



**databricks**



# Sponsors

---

Track



**FORTYTWO**  
— ANALYTICS —



Event





[HTTPS://GITHUB.COM/SIWHITELEY/DATABRICKSETL](https://github.com/siwhiteley/databricksetl)





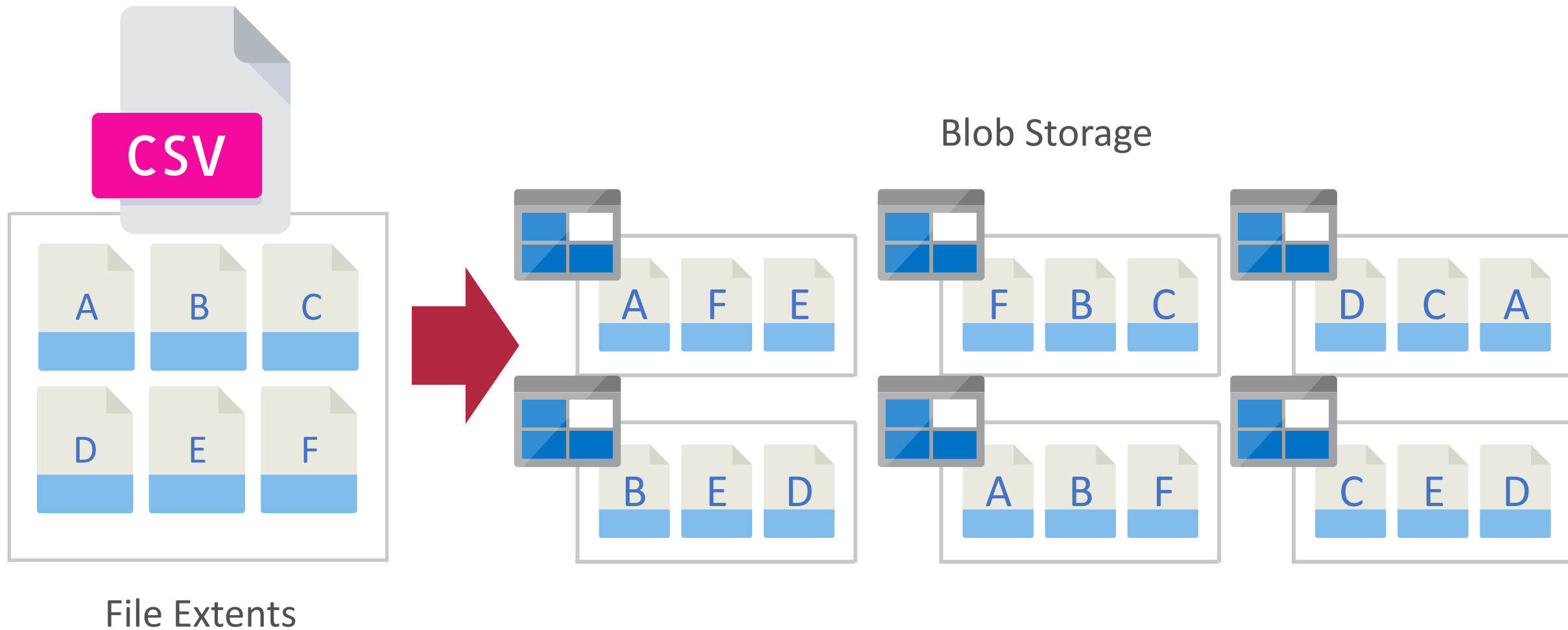
Google File System Papers  
Released

**2003**





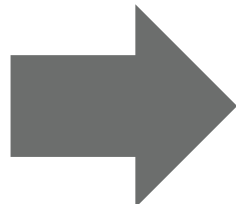
# HDFS IN AZURE





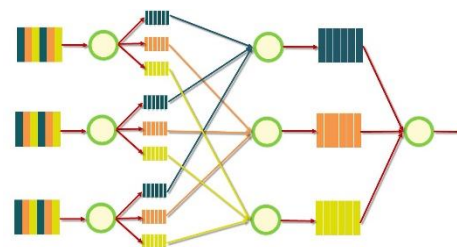
Google File System Papers  
Released

**2003**



Google MapReduce Papers

**2004**



## Input

This is a large document

This might come from a log file

Or it might come from your user input

## Split

This is a large document

This might come from a log file

Or it might come from your user input

## Map

This, 1  
is, 1  
a, 1  
large, 1  
document, 1

This, 1  
might, 1  
come, 1  
from, 1  
a, 1  
log, 1  
file, 1

Or, 1  
it, 1  
might, 1  
come, 1  
from, 1  
your, 1  
user, 1  
input, 1

## Shuffle

This (1,1)

is (1)

a (1,1)

large (1)

document (1)

might (1,1)

come (1,1)

from (1,1)

log (1)

Or (1)

file (1)

it (1)

your (1)

user (1)

input (1)

## Reduce

This, 2  
might, 2  
come, 2  
from, 2  
a, 2  
is, 1  
large, 1  
document, 1  
log, 1  
Or, 1  
file, 1  
it, 1  
your, 1  
user, 1  
input, 1

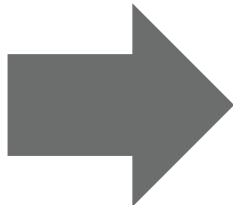






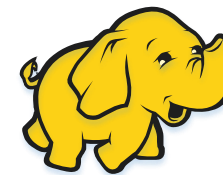
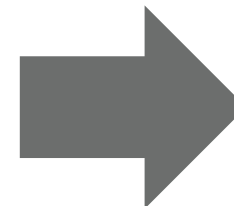
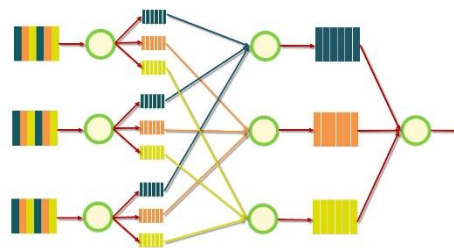
Google File System Papers  
Released

**2003**



Google MapReduce Papers

**2004**



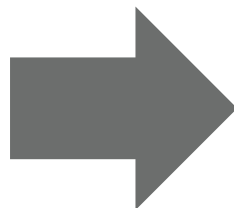
**2006**

Apache Hadoop  
project created



Matei Zaharia starts Spark  
project

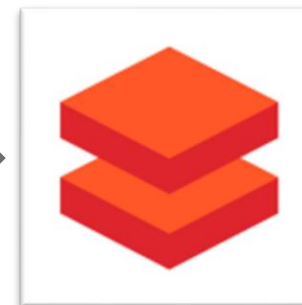
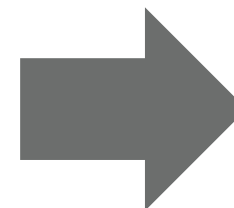
**2009**



THE  
**APACHE**  
SOFTWARE FOUNDATION

Project donated to Apache  
Foundation

**2013**



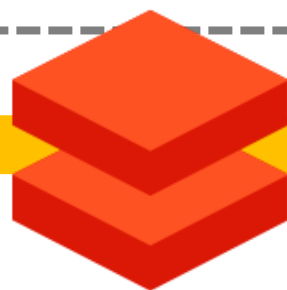
Databricks founded by Matei & UC  
Berkeley Colleagues

**2013**





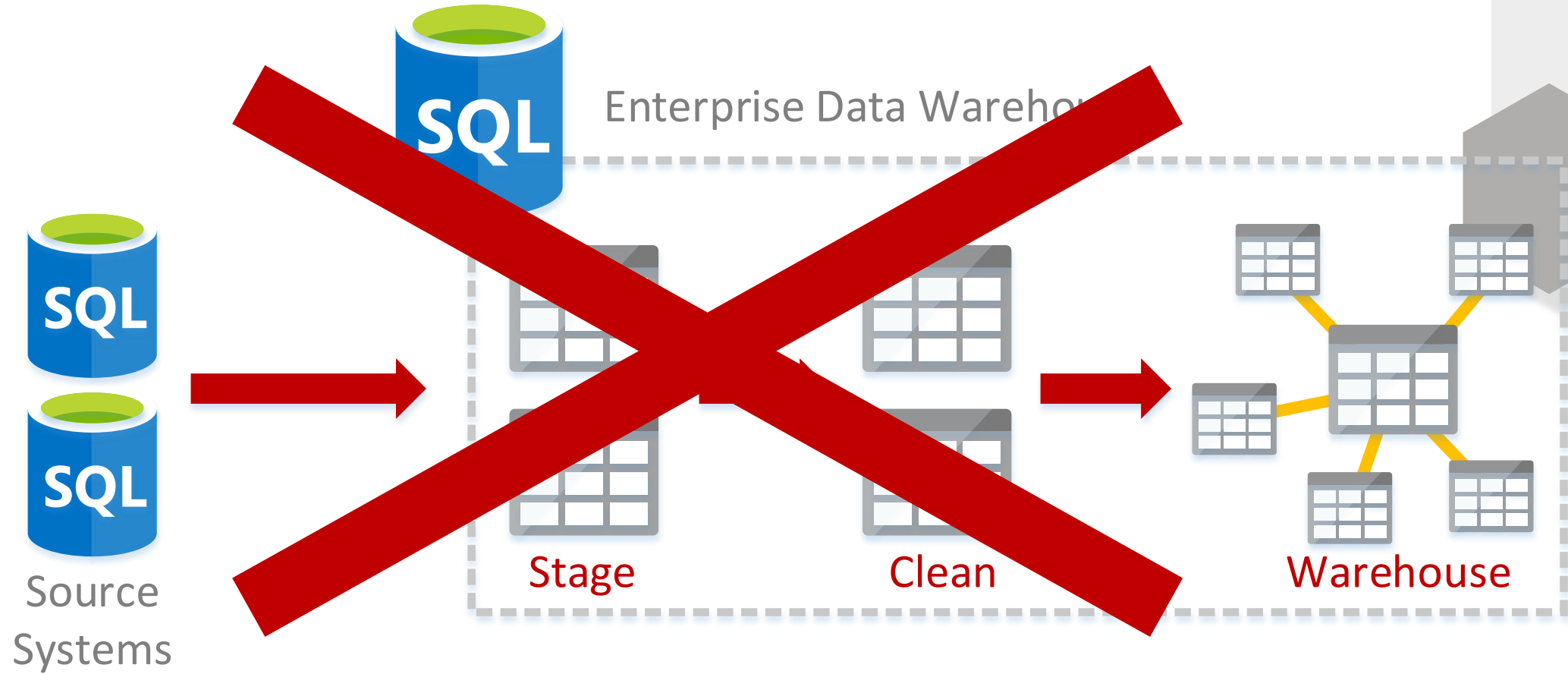
Language Flexibility



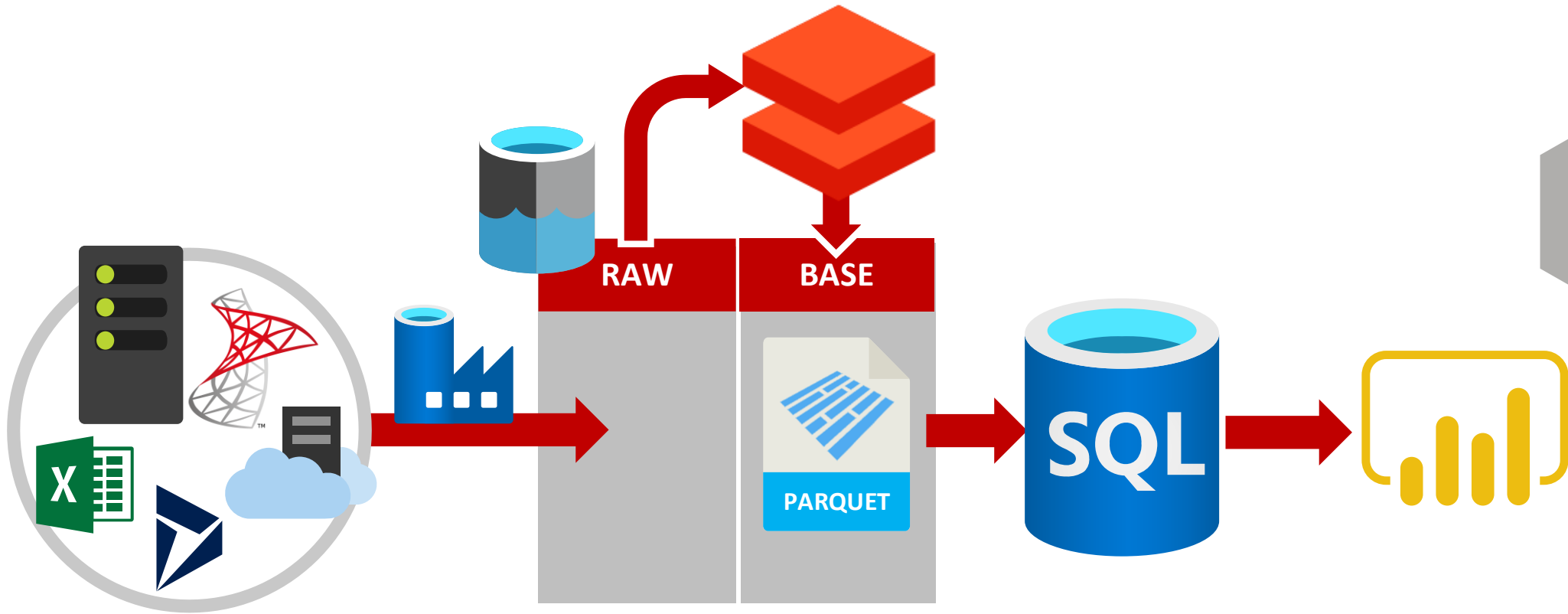
Integrations



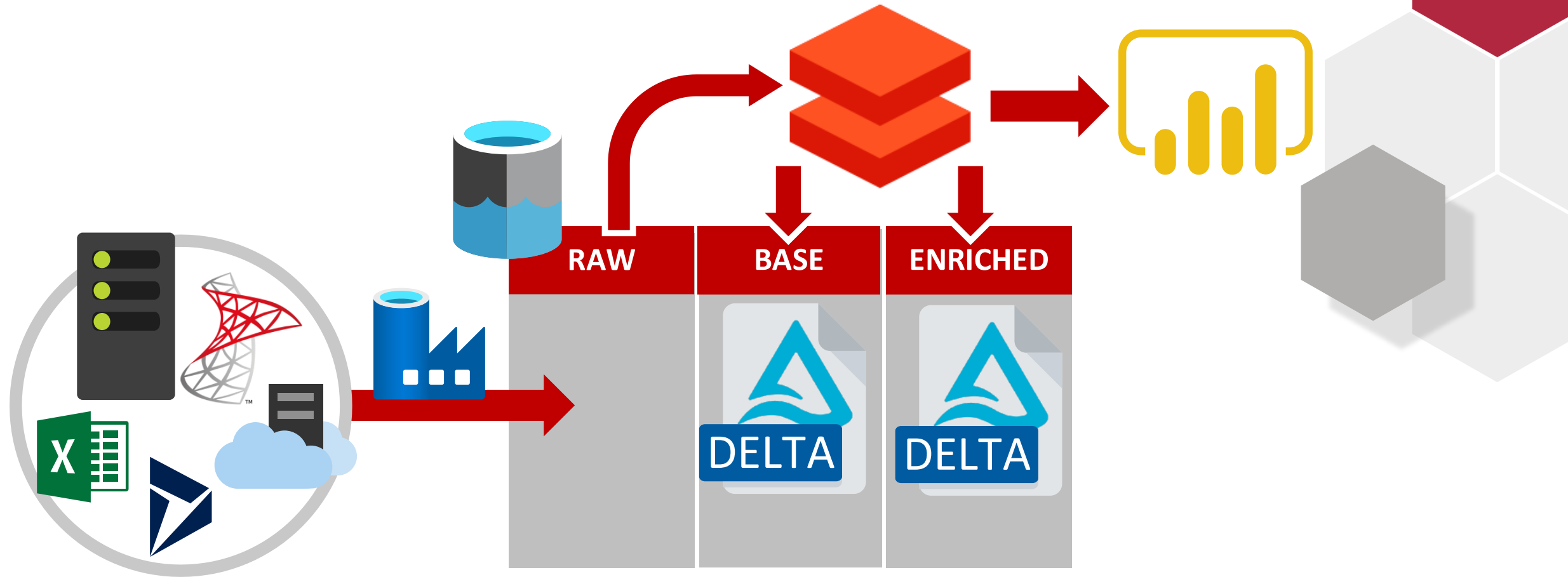
BUT MOST OF ALL...



# MODERN DATA WAREHOUSE



# DATA LAKEHOUSE





[www.advancinganalytics.co.uk](http://www.advancinganalytics.co.uk)

# SPARK INTERNALS



@ADVANCINGANALYTICS

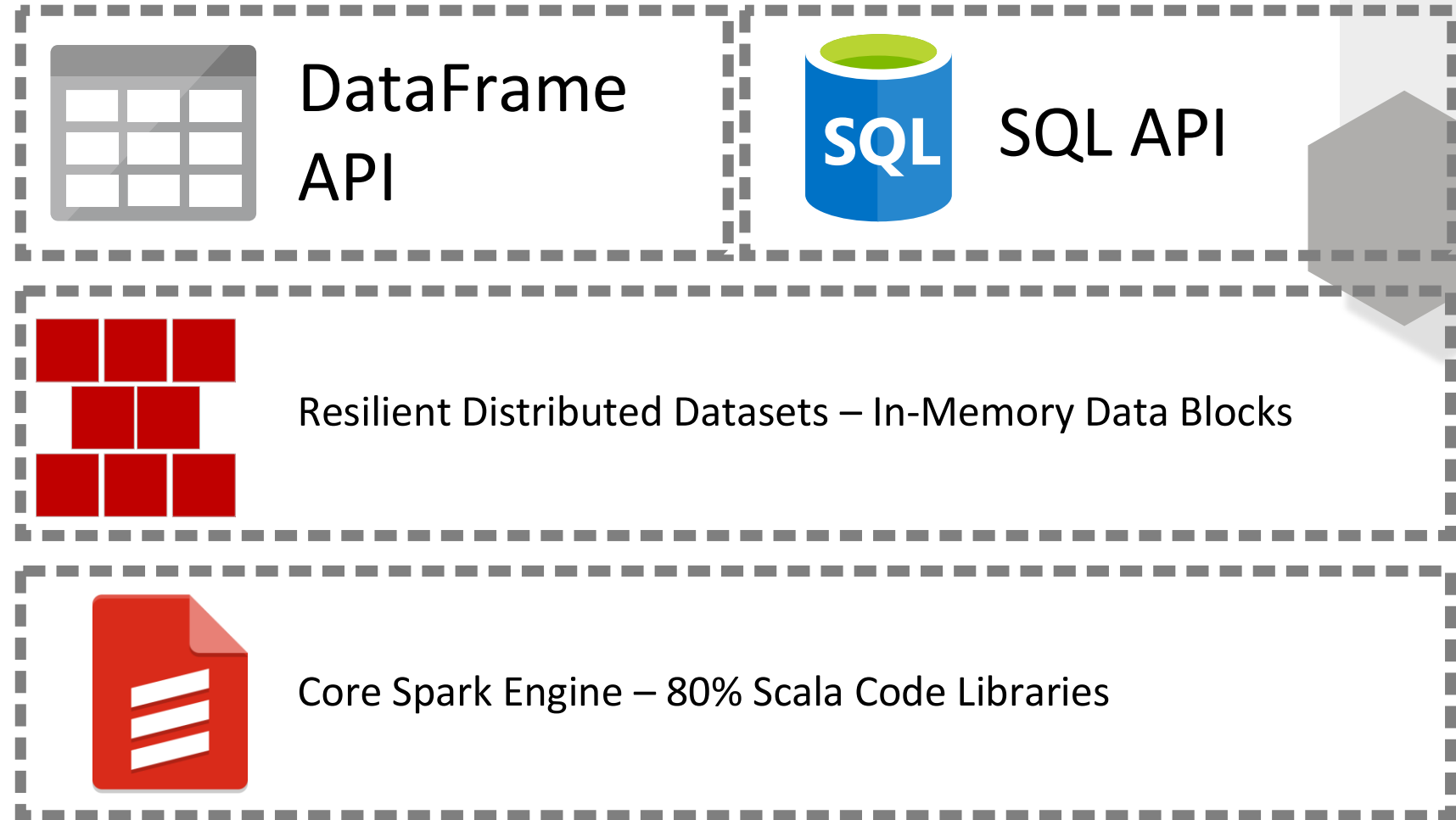


@ADVANALYTICSUK

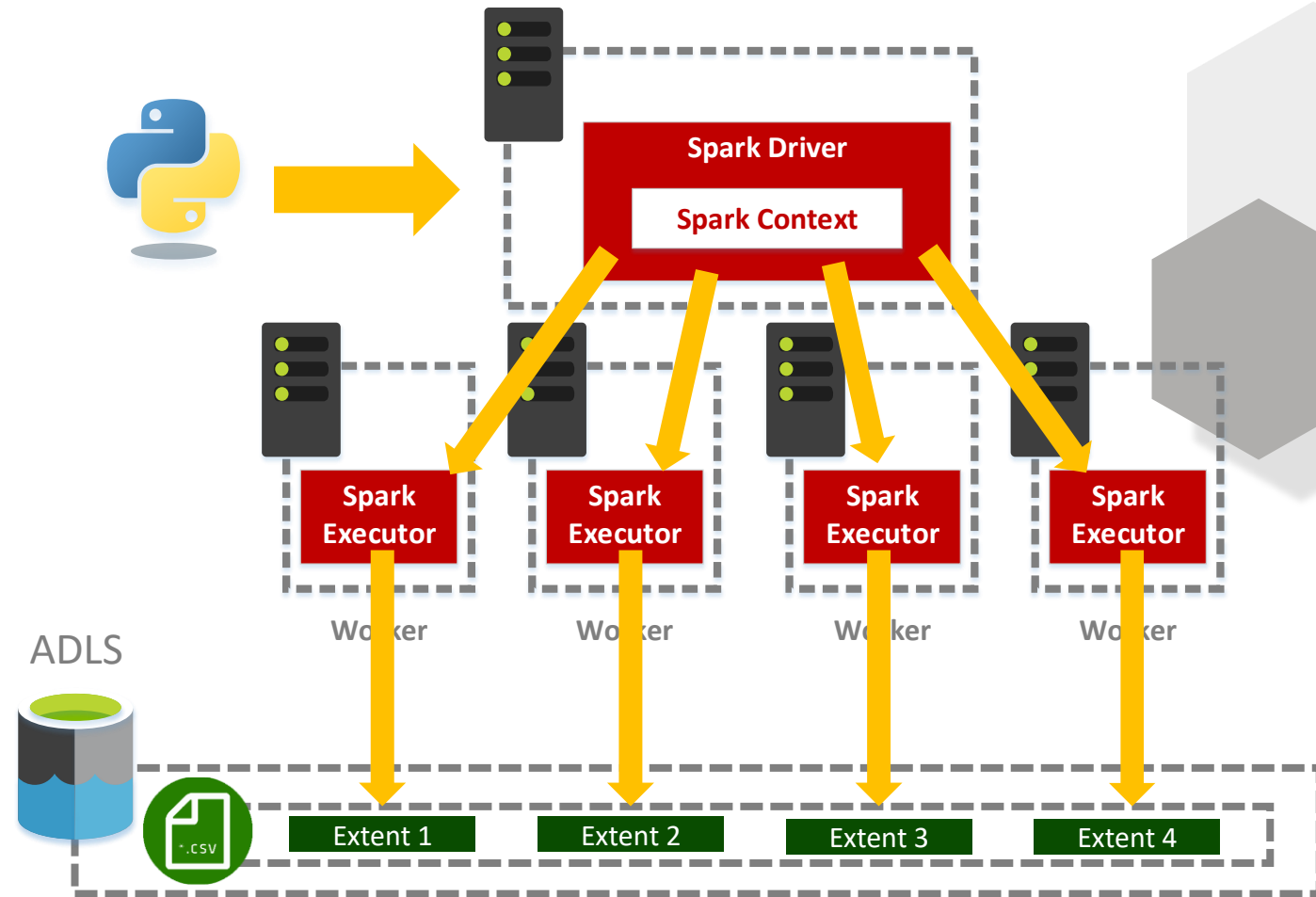


/ADVANCING ANALYTICS

# UNDER THE HOOD

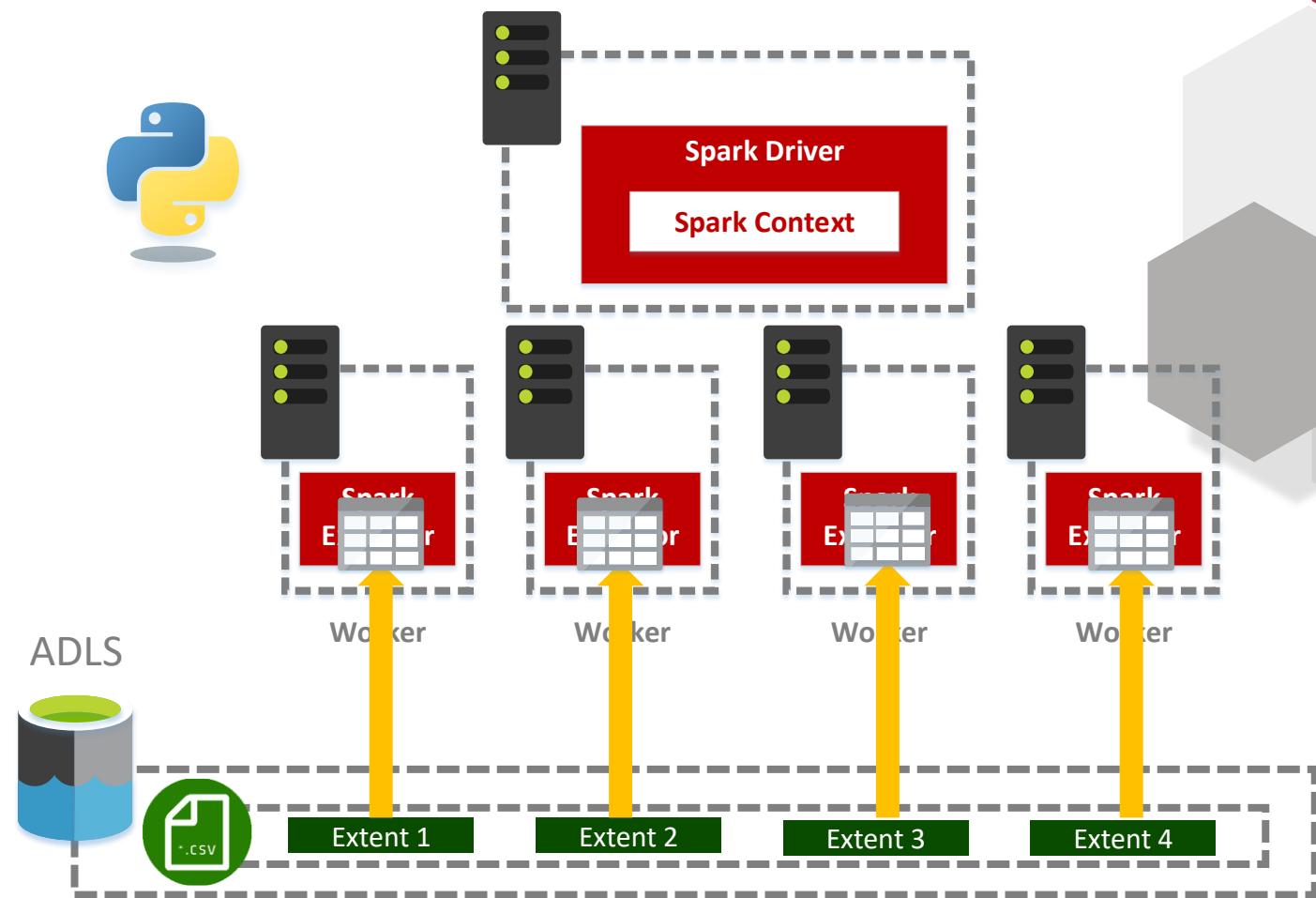


# DISTRIBUTED COMPUTE

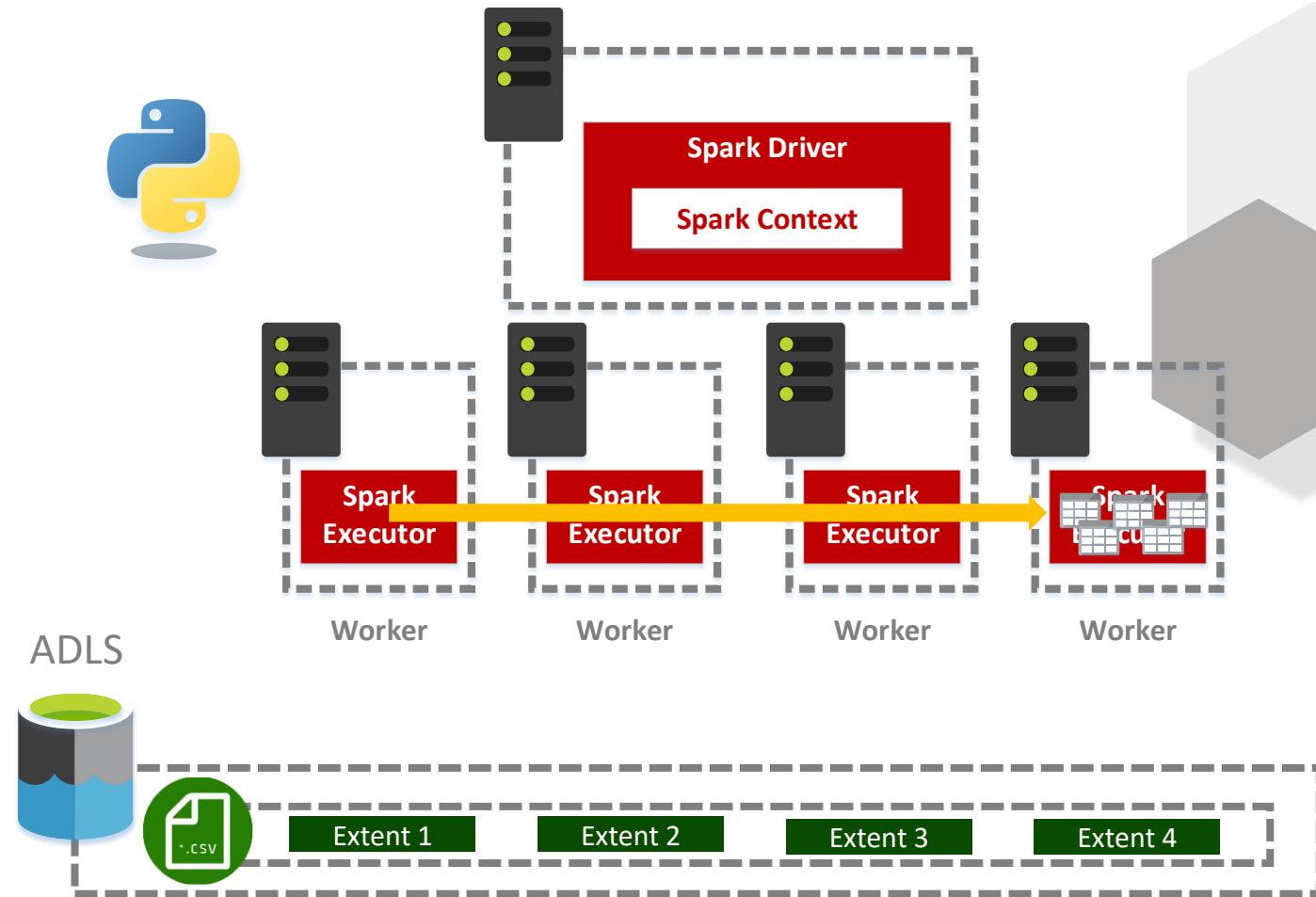




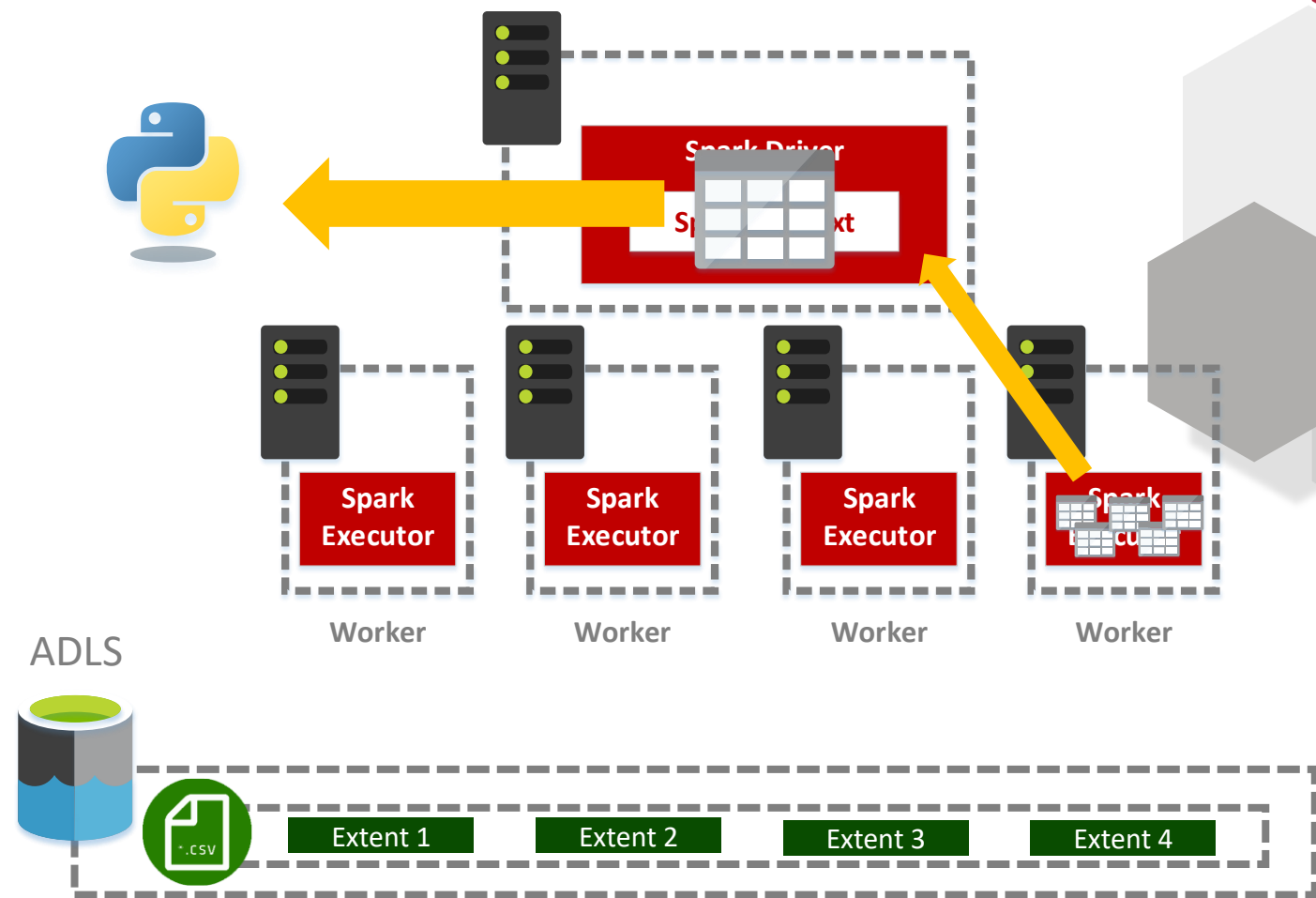
# DISTRIBUTED COMPUTE



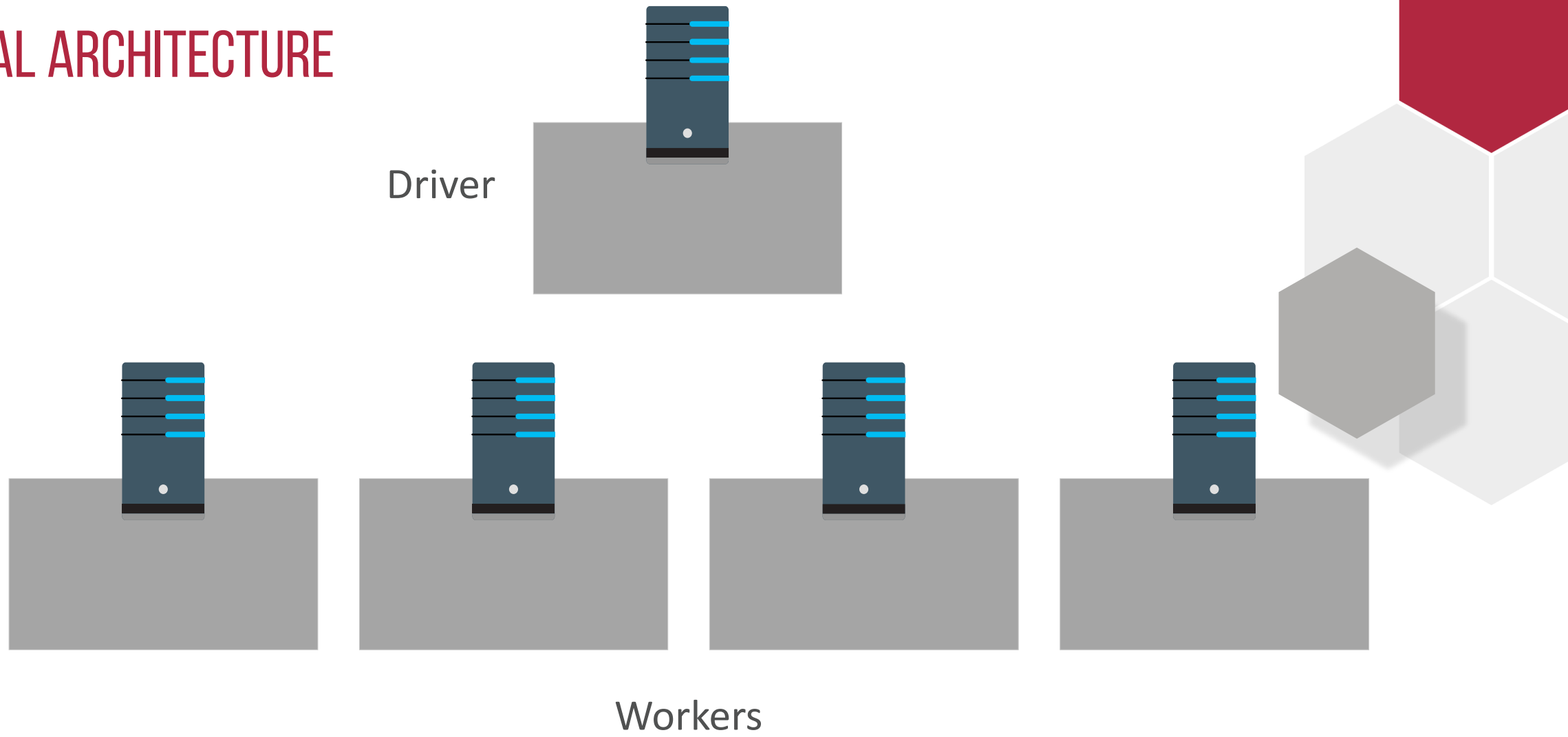
# DISTRIBUTED COMPUTE



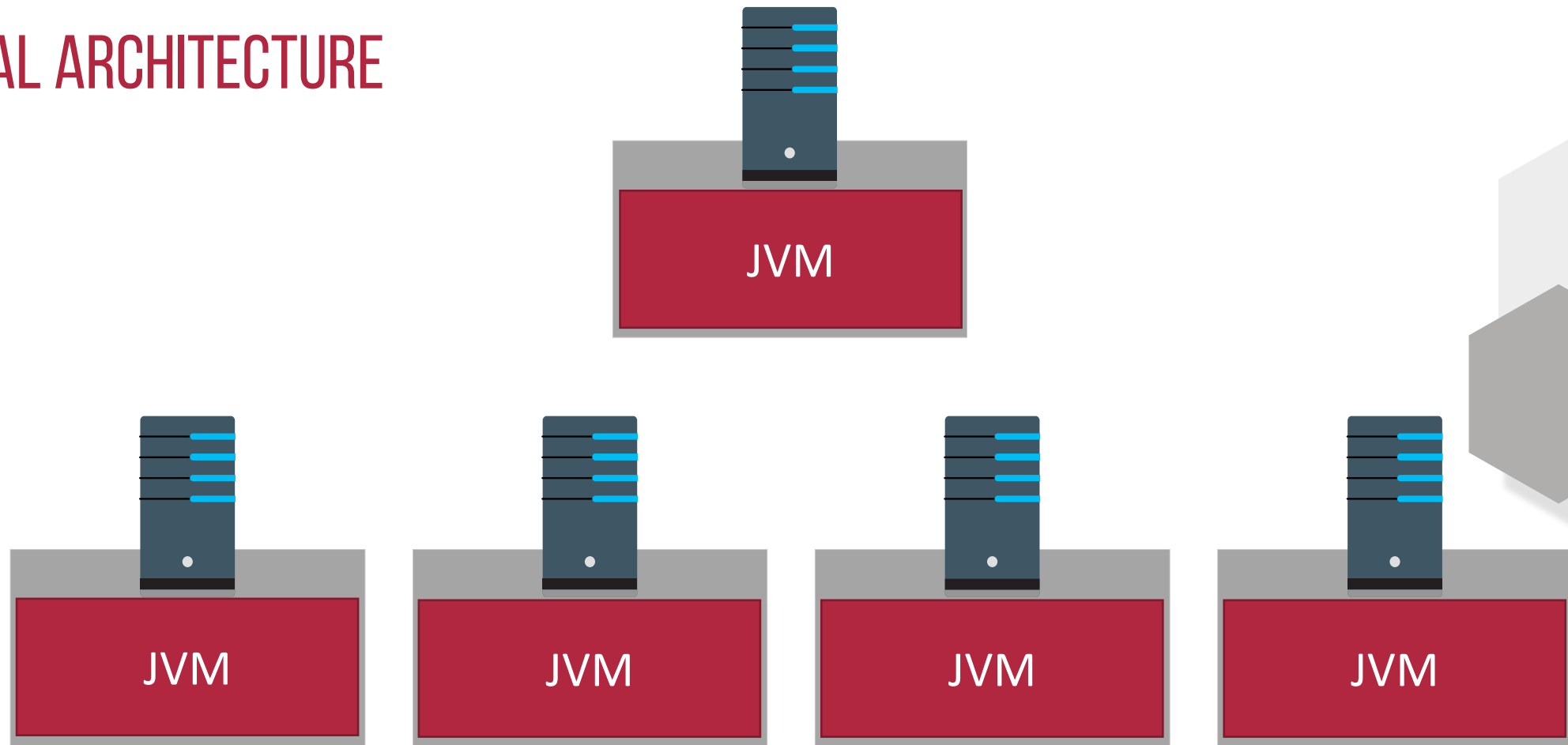
# DISTRIBUTED COMPUTE



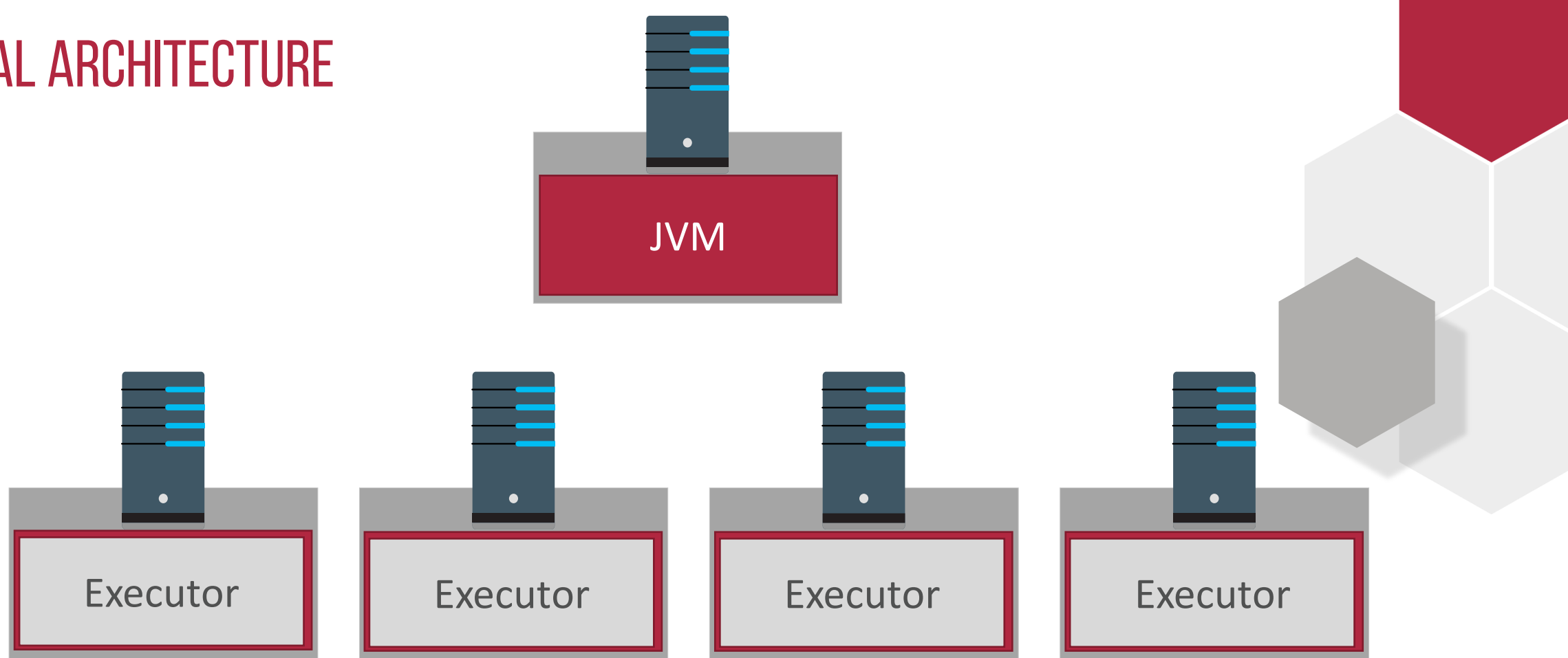
# PHYSICAL ARCHITECTURE



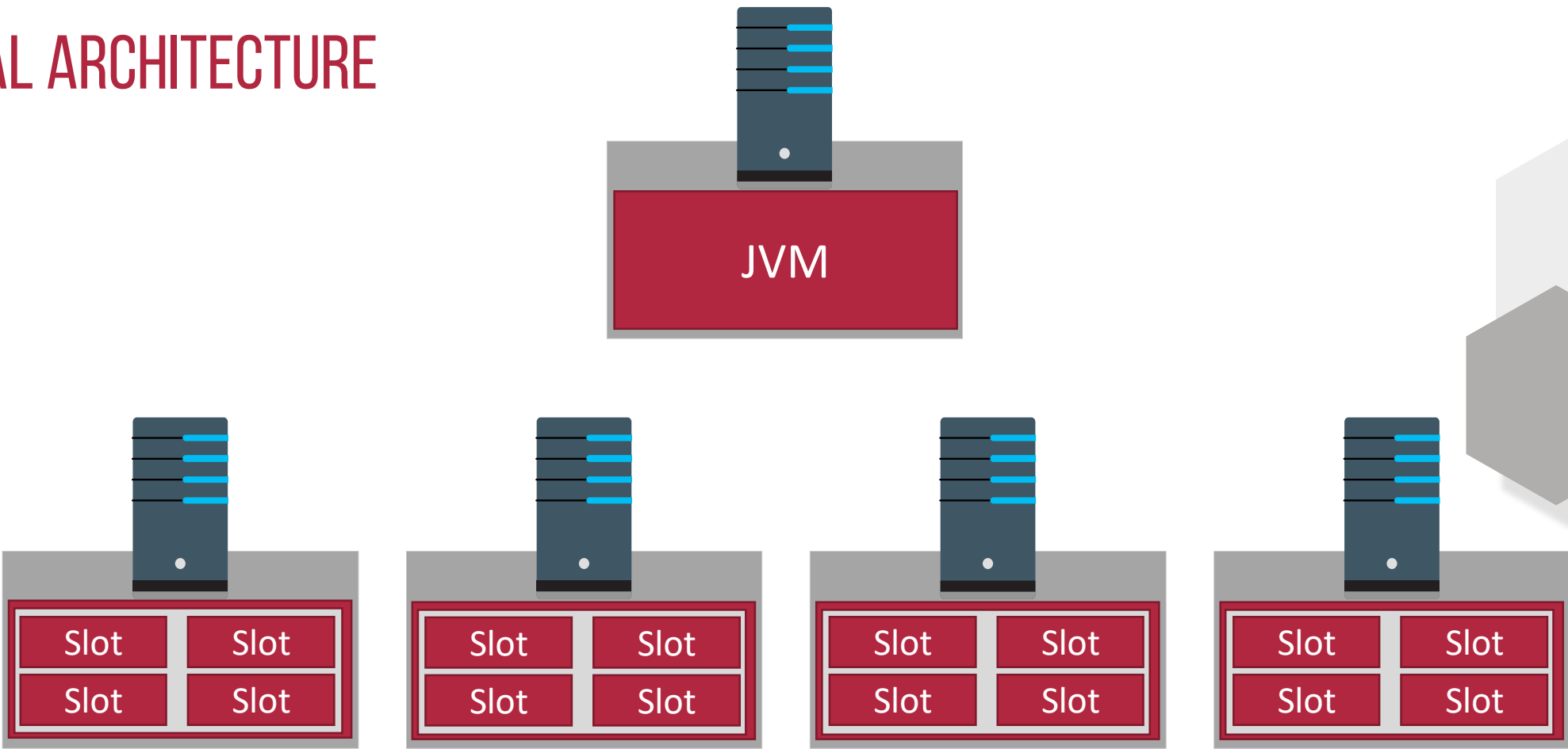
# PHYSICAL ARCHITECTURE



# PHYSICAL ARCHITECTURE

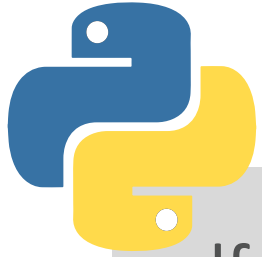


# PHYSICAL ARCHITECTURE





# EXECUTION HIERARCHY



```
df = spark.read...  
    .filter(...)  
    .join(...)  
    .write...
```

Transformations

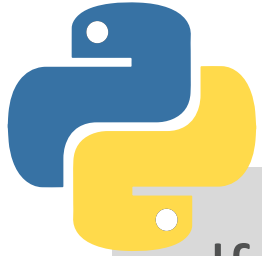
Actions

Narrow

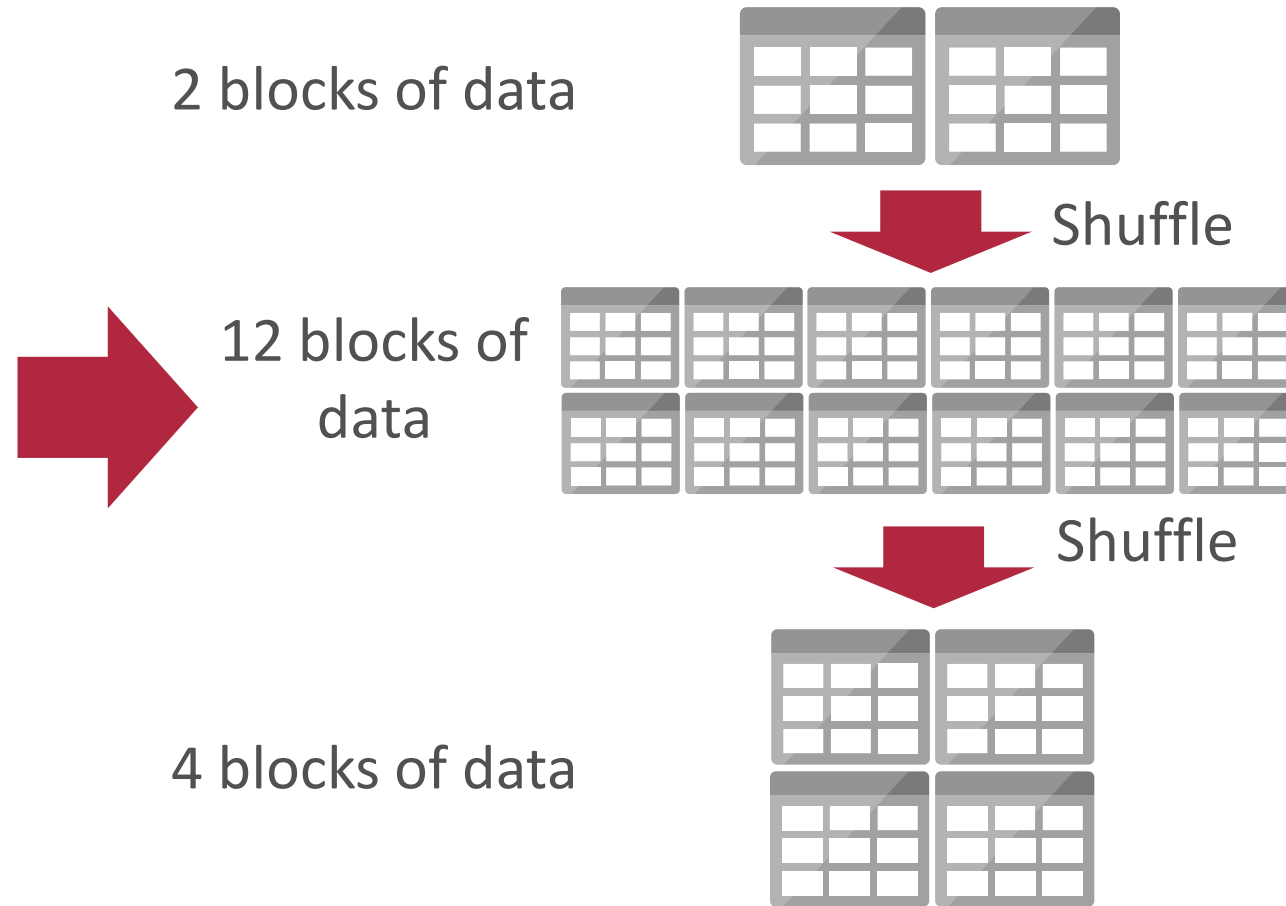
Wide



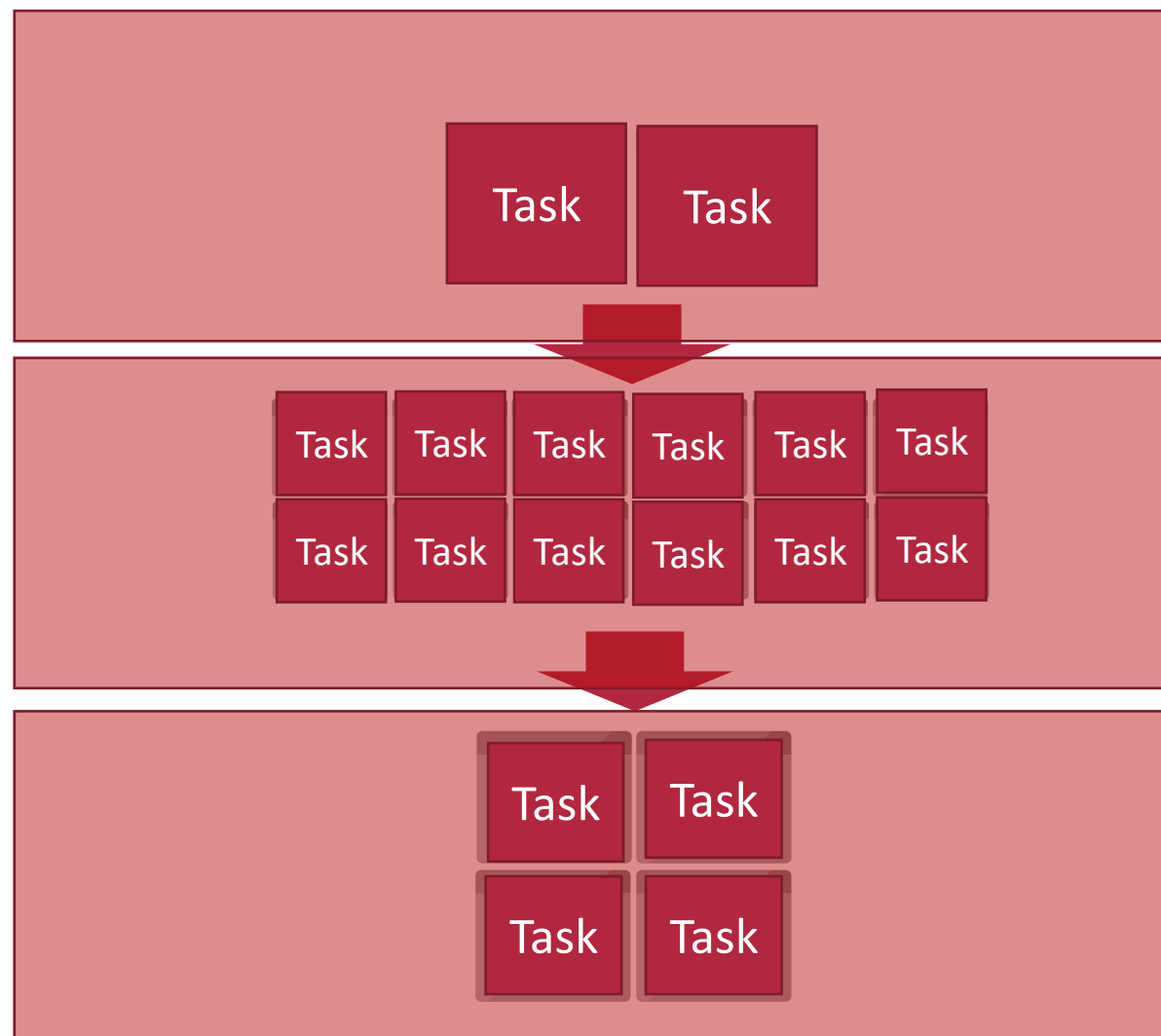
# EXECUTION HIERARCHY



```
df = spark.read...  
    .filter(...)  
    .join(...)  
    .write...
```



JOB



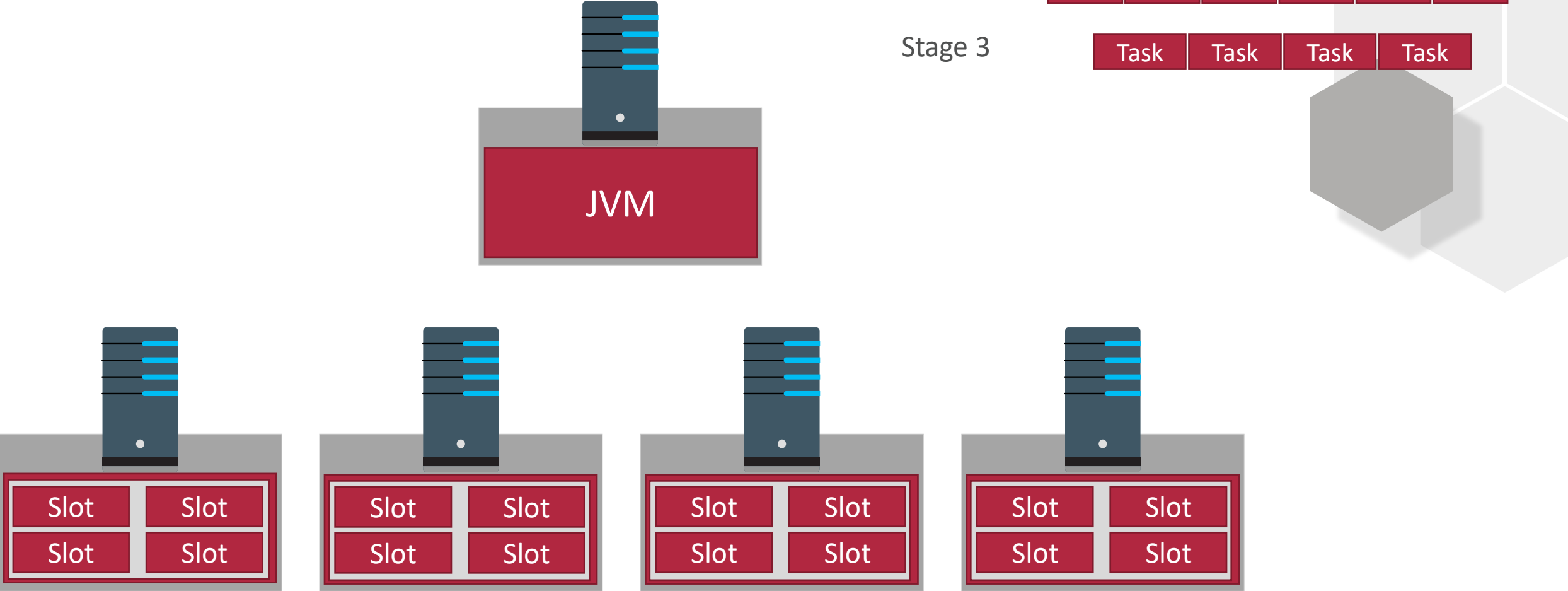
Stage 1

Stage 2

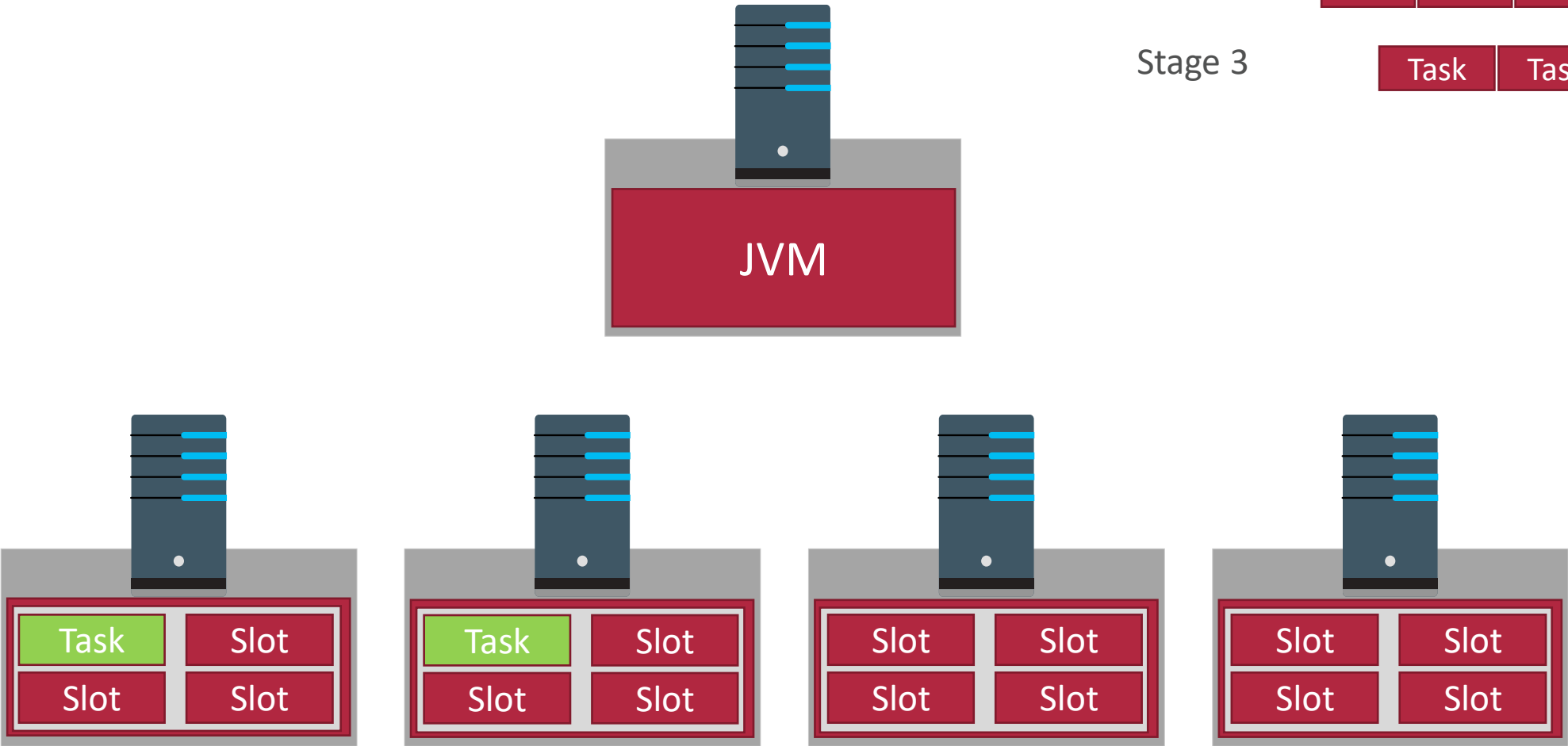
Stage 3



# PHYSICAL ARCHITECTURE



# PHYSICAL ARCHITECTURE



Stage 1



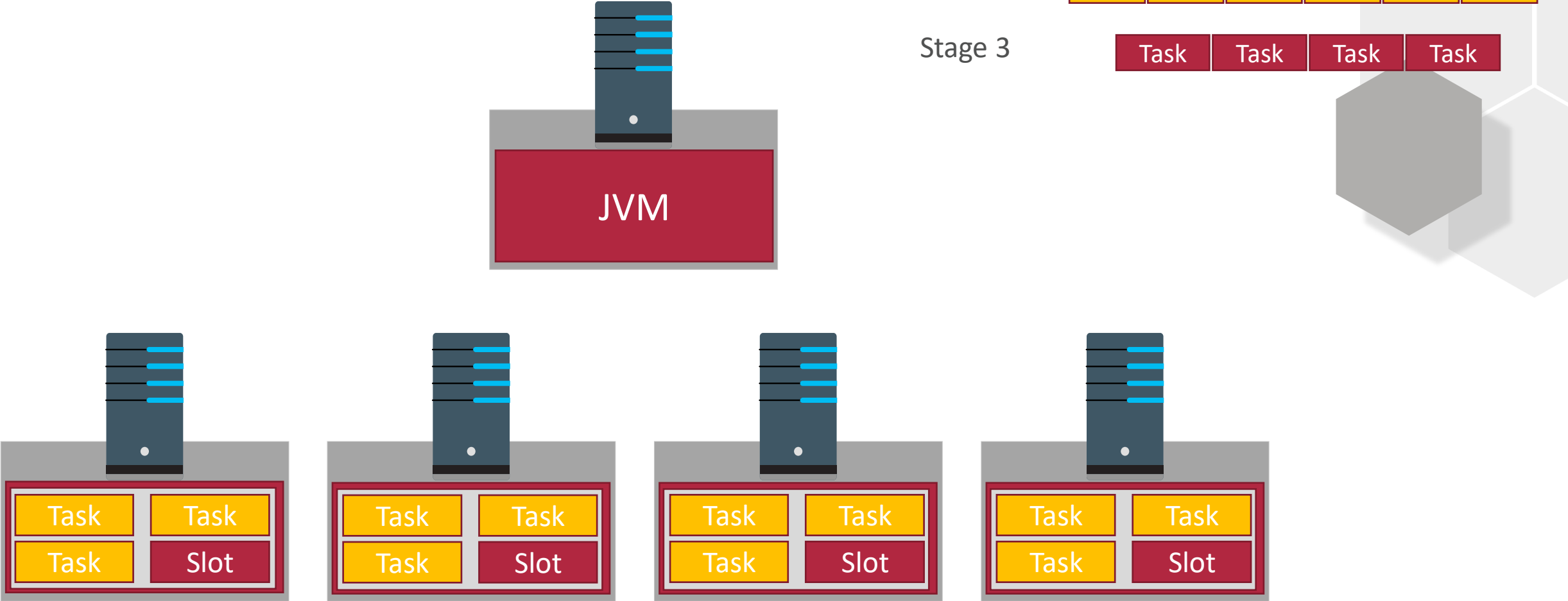
Stage 2



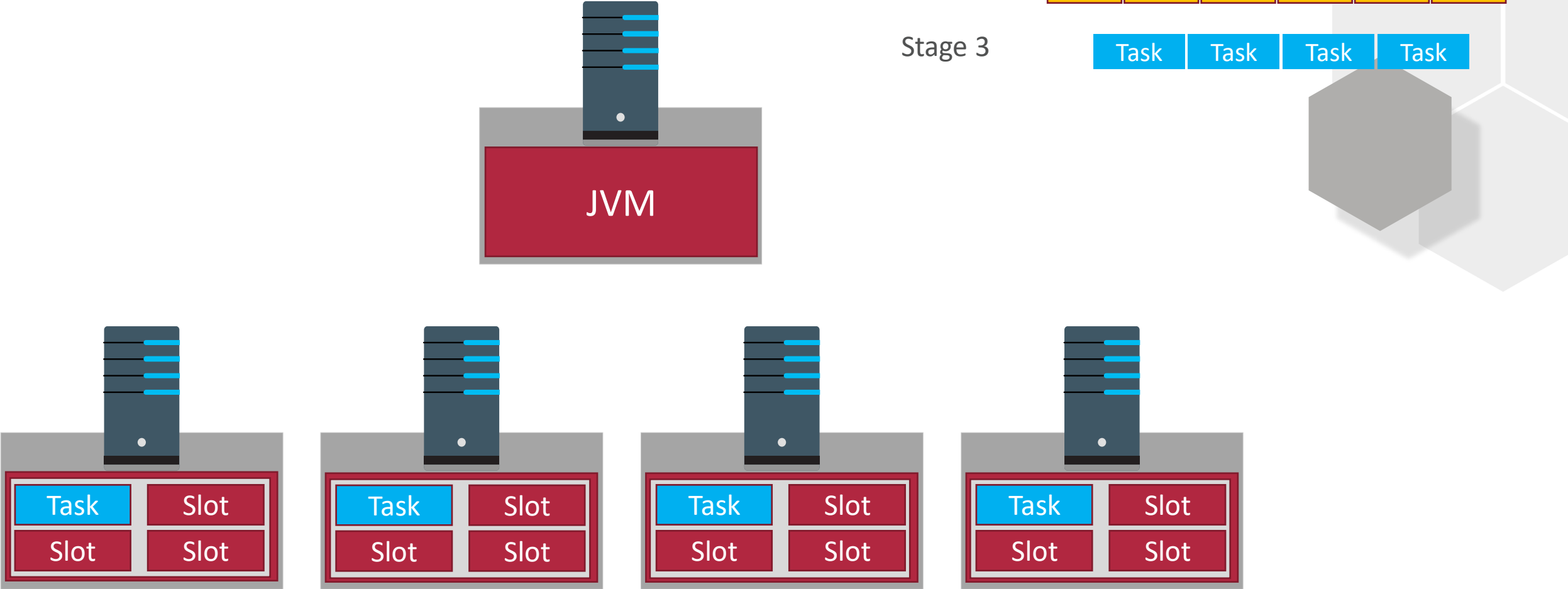
Stage 3



# PHYSICAL ARCHITECTURE



# PHYSICAL ARCHITECTURE







[www.advancinganalytics.co.uk](http://www.advancinganalytics.co.uk)

# ETL IN DATABRICKS



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS

# SCHEMA ON READ – INFER SCHEMA

Cmd 3

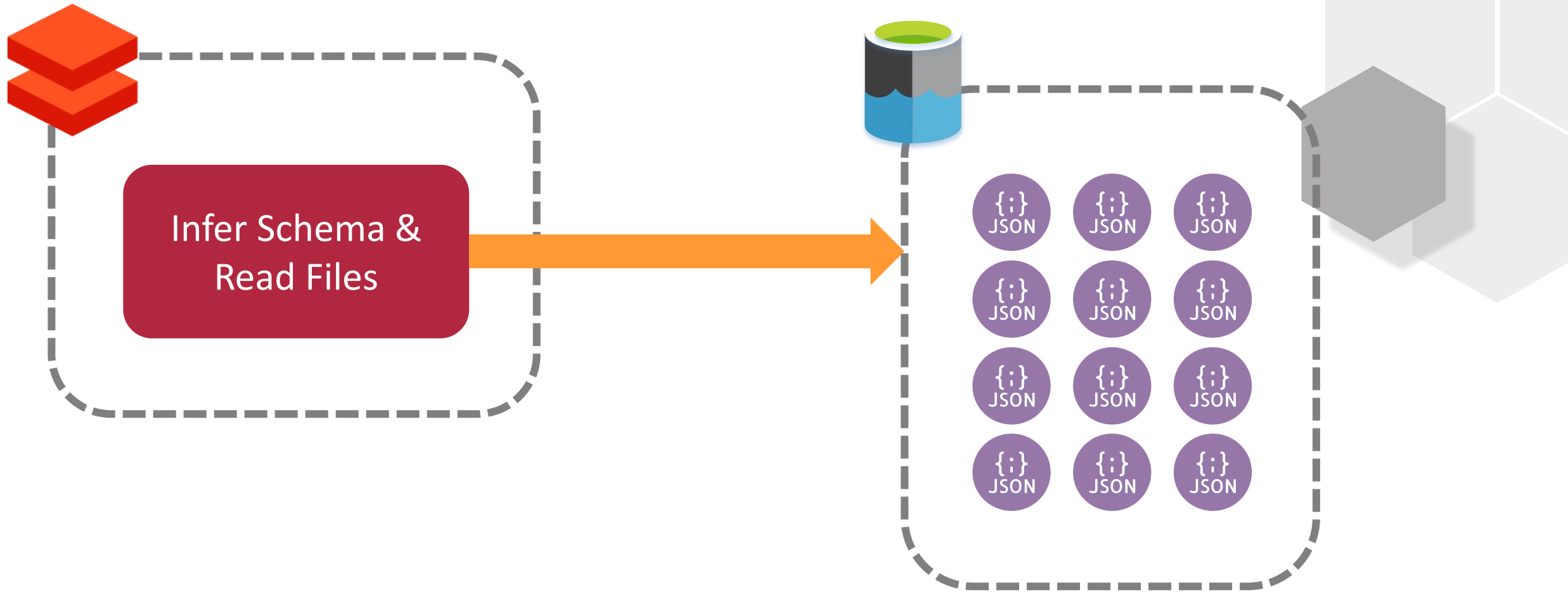
```
1 df = sqlContext.read.format("csv") \  
2   .option("header", "true") \  
3   .option("inferSchema", "true") \  
4   .load("abfss://root@dblake.dfs.core.windows.net/RAW/Public/Taxi/v1/SmallSlice.csv")
```

▼  df: pyspark.sql.dataframe.DataFrame

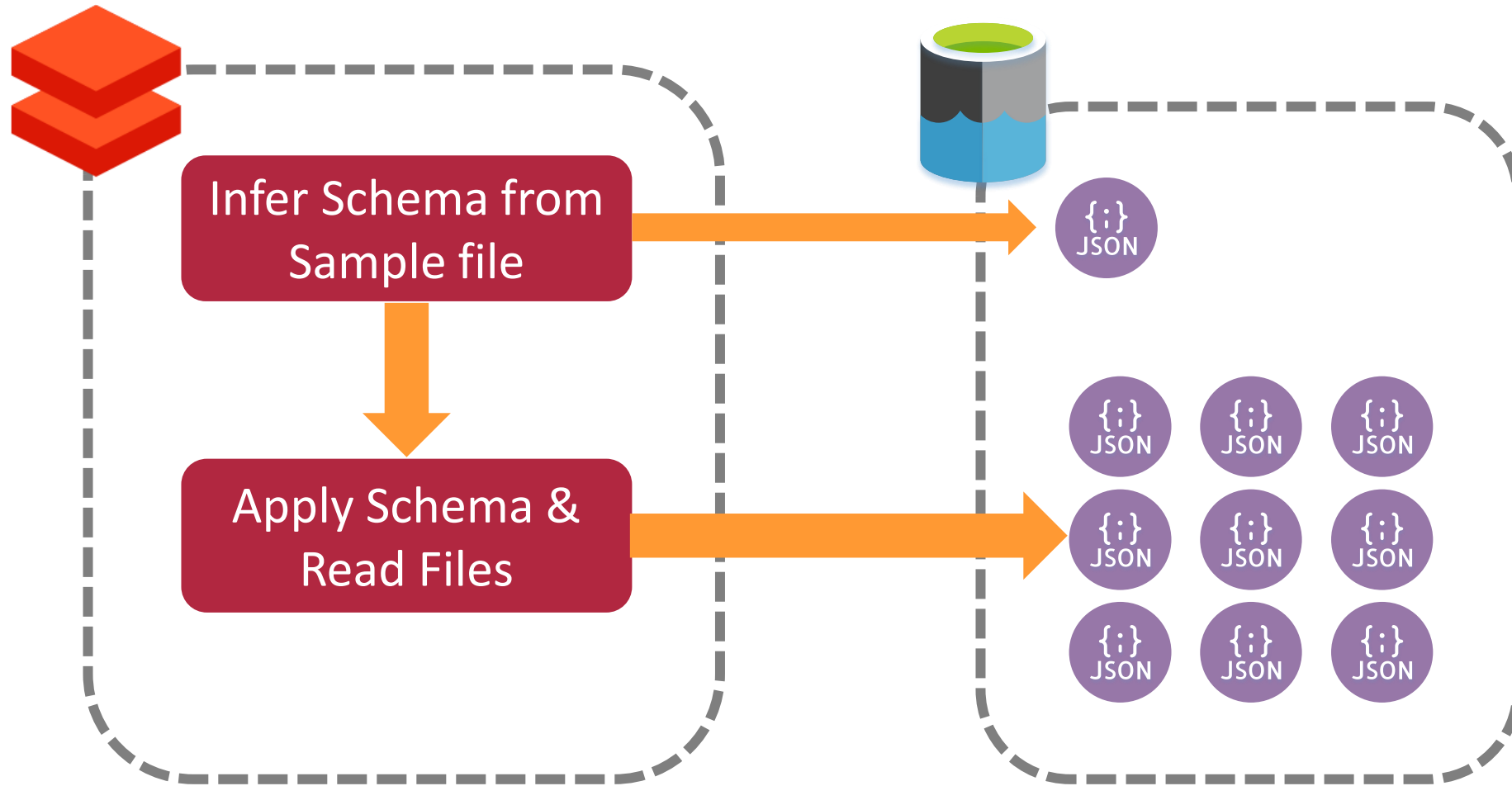
- Dispatching\_base\_num: string
- Pickup\_DateTime: timestamp
- DropOff\_datetime: string
- PUlocationID: integer
- DOlocationID: string



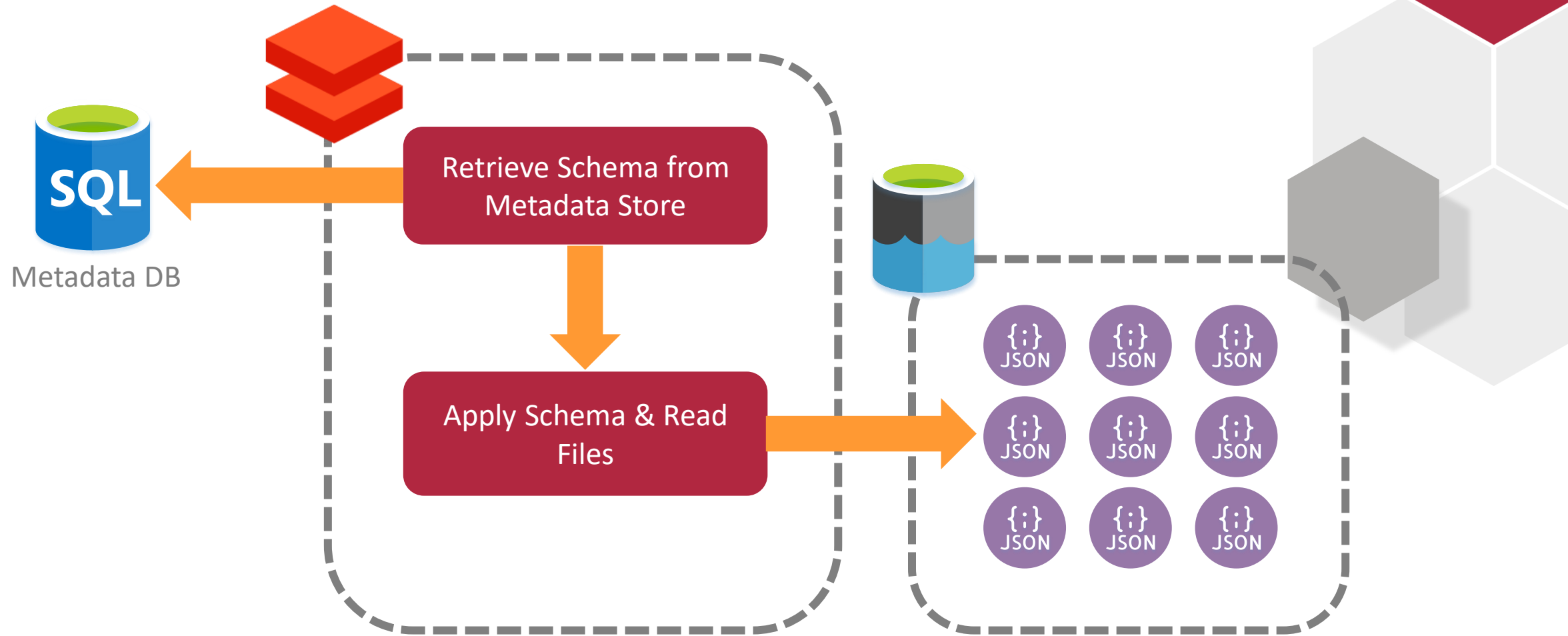
# SCHEMA ON READ – INFER SCHEMA



# SCHEMA ON READ – SAMPLE FILES



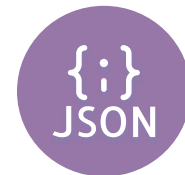
# SCHEMA ON READ – METADATA STORE



# DATA VALIDATION

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

Schema



A	2	XY	
B	Y	XR	
C	3	E	X



# DATA VALIDATION

A	2	XY
null	null	null
C	3	E

- 1 PERMISSIVE**  
NULL DataType Failures  
Drop Extra Columns

A	2	XY	
B	Y	XR	
C	3	E	X

- 2 FAILFAST**  
Fail on any error

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

- 3 DROPMALFORMED**  
Drop rows with DataType Failures  
Drop rows with Extra Columns

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




# DATA VALIDATION

```
#Add a system attribute to our structure
newSchema.add("_corrupt_record", StringType(), True)

#Load the data so we can inspect it
df = (spark
      .read
      .schema(newSchema)
      .option("mode", "PERMISSIVE")
      .csv(dataLocation)
      )
```

```
#Separate bad rows automatically
baddf = (spark
        .read
        .option("badRecordsPath", "/mnt/data/RAW/TaxiZones/_reject/")
        .schema(mySchema)
        .csv(location)
        )
```

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





# DATA PROCESSING WITH SPARK

- Working with Flat Files
- Working with Dodgy Data
- Meta-Driven Systems





[www.advancinganalytics.co.uk](http://www.advancinganalytics.co.uk)

# ORCHESTRATION



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS

# 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")
```

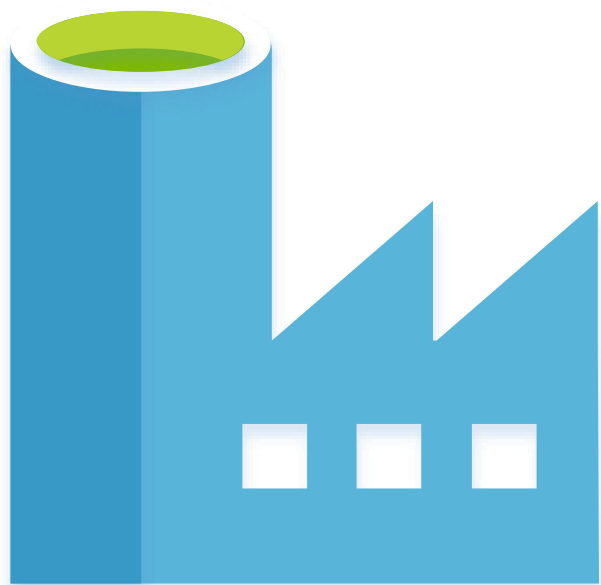
AdventureWorks Table :

Entity Name :

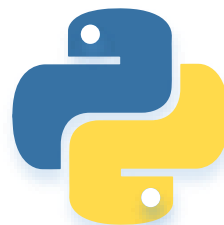


Cmd 8

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



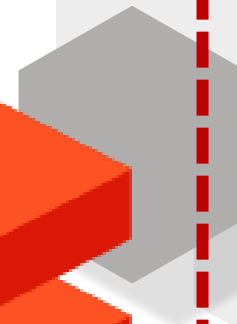
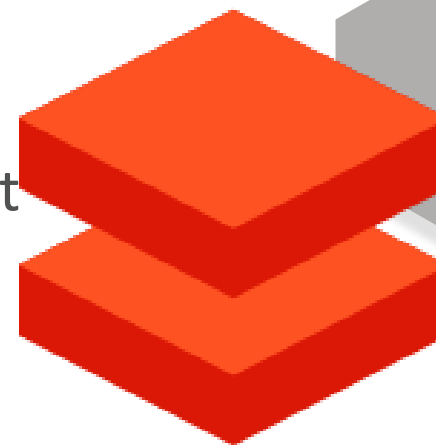
Jupyter Notebook



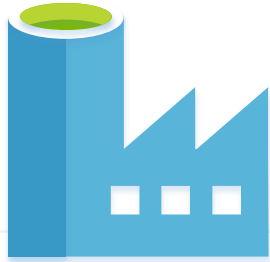
Python Script



Jar File



# AZURE DATA FACTORY



Notebook path \*

[Browse](#)

Base Parameters

[+ New](#) | [Delete](#)

<input type="checkbox"/>	NAME	VALUE
<input type="checkbox"/>	entity_name	Taxi

Entity Name :

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

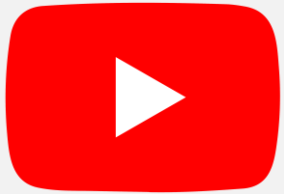


# DATABRICKS ORCHESTRATION

- Adding input & output parameters
- Orchestrating Databricks with ADF



# THANKS!



[youtube.com/c/AdvancingAnalytics](https://youtube.com/c/AdvancingAnalytics)



[@MrSiWhiteley](https://twitter.com/MrSiWhiteley)



[www.advancinganalytics.co.uk](http://www.advancinganalytics.co.uk)





# TEAMS MEETINGS

Azure BI - Join Microsoft Teams Meeting

Azure DB - Join Microsoft Teams Meeting

Cloud Development - Join Microsoft Teams Meeting

Data Platform - Join Microsoft Teams Meeting

Power BI - Join Microsoft Teams Meeting

