

# THE SPARK SHOWDOWN



Simon Whiteley  
@MrSiWhiteley





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

# HERE'S HOW IT WORKS



@ADVANCINGANALYTICS



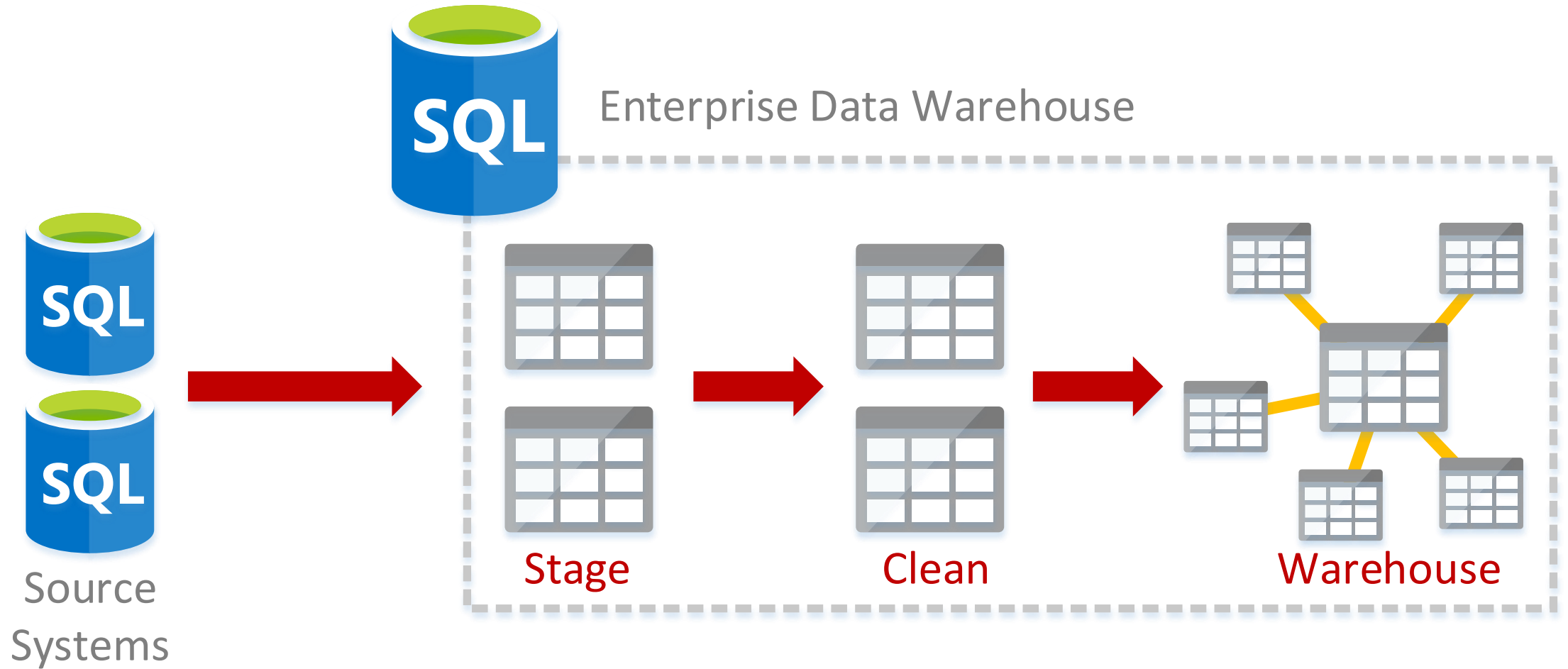
@ADVANALYTICSUK



/ADVANCING ANALYTICS

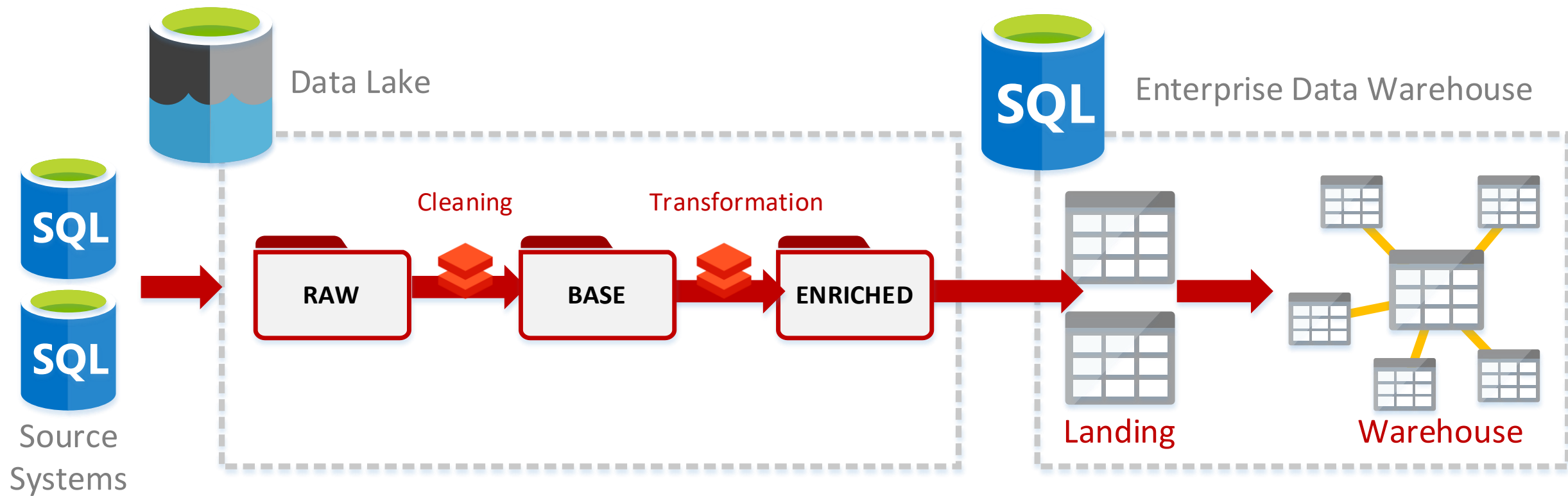


# THE RETIRED CHAMPION

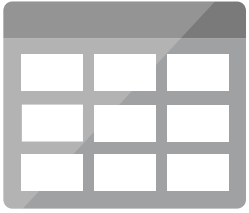




# THE NEW GENERATION



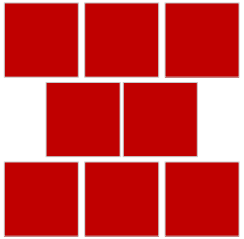
# THE SPARK API ABSTRACTIONS



DataFrame  
API



SQL API



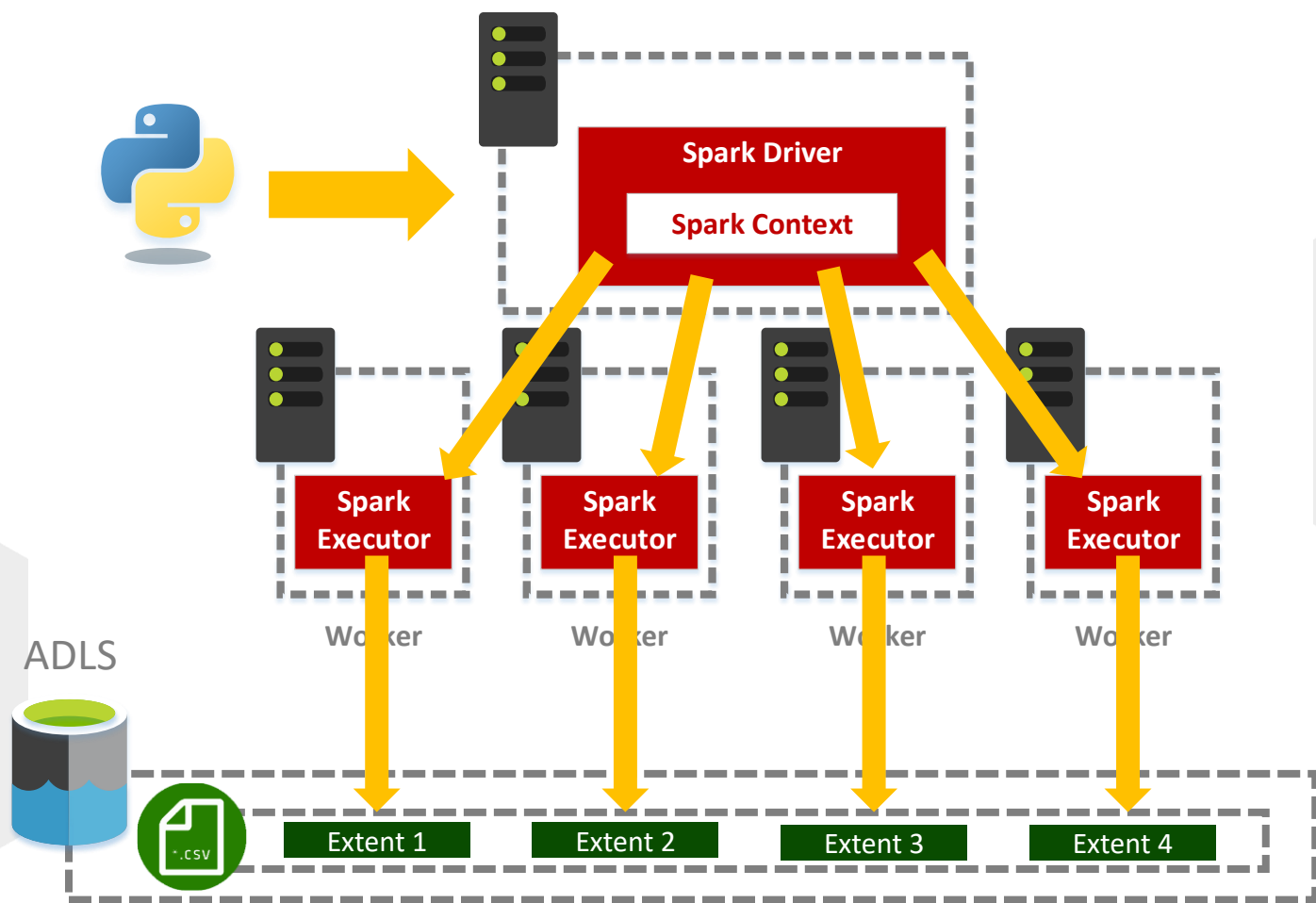
Resilient Distributed Datasets – In-Memory Data Blocks



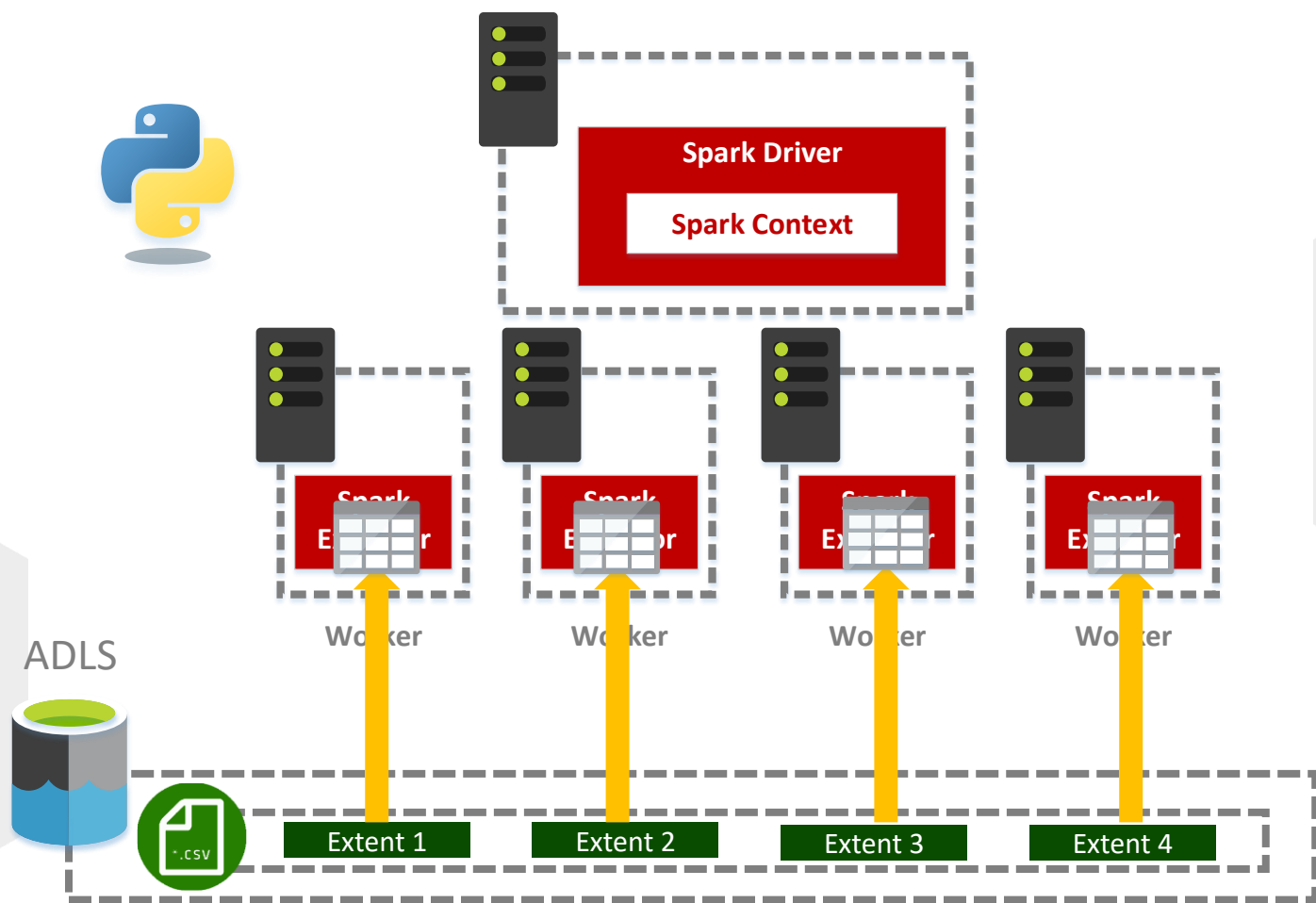
Core Spark Engine – 80% Scala Code Libraries



# DISTRIBUTED COMPUTE

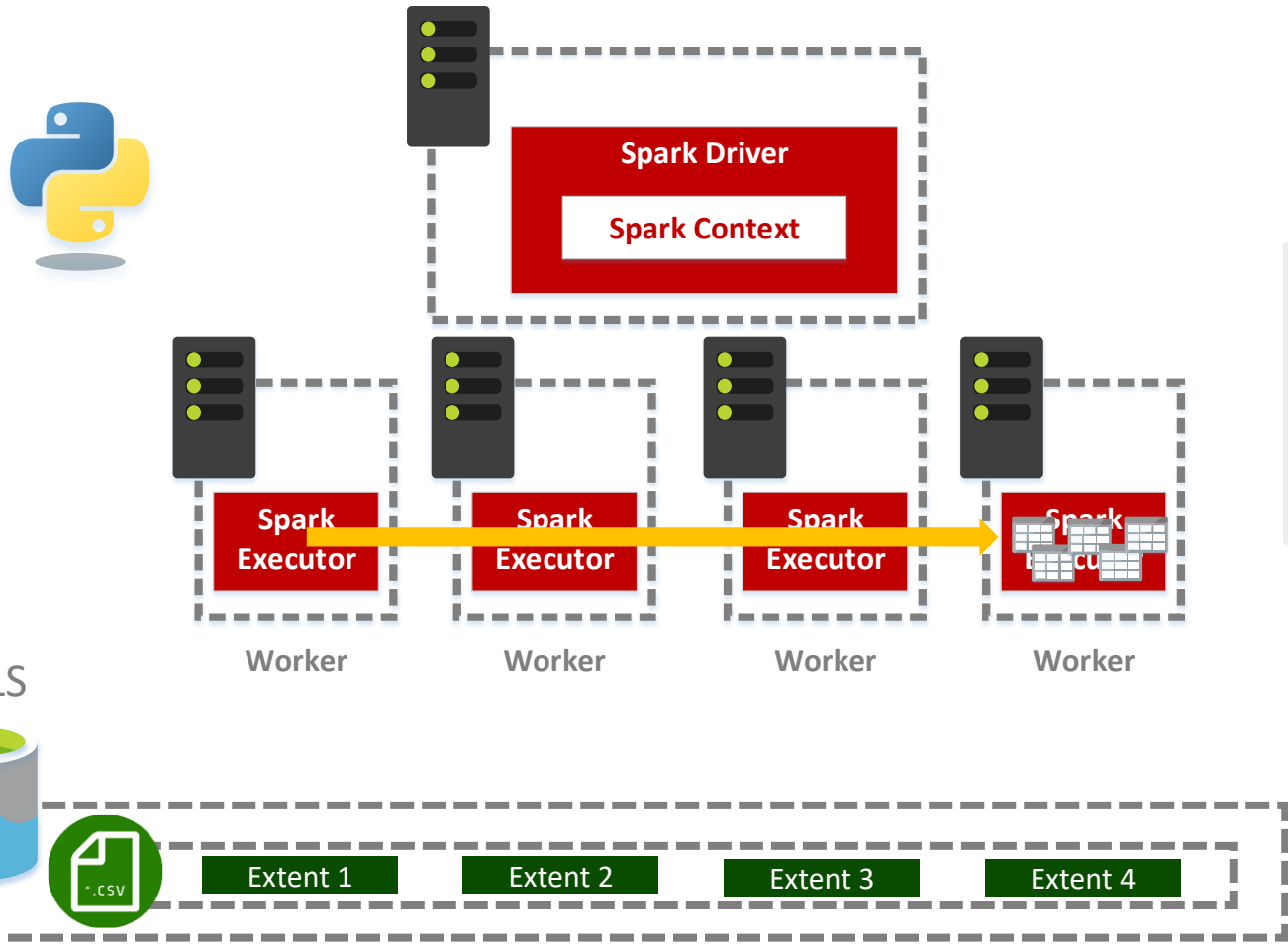


# DISTRIBUTED COMPUTE



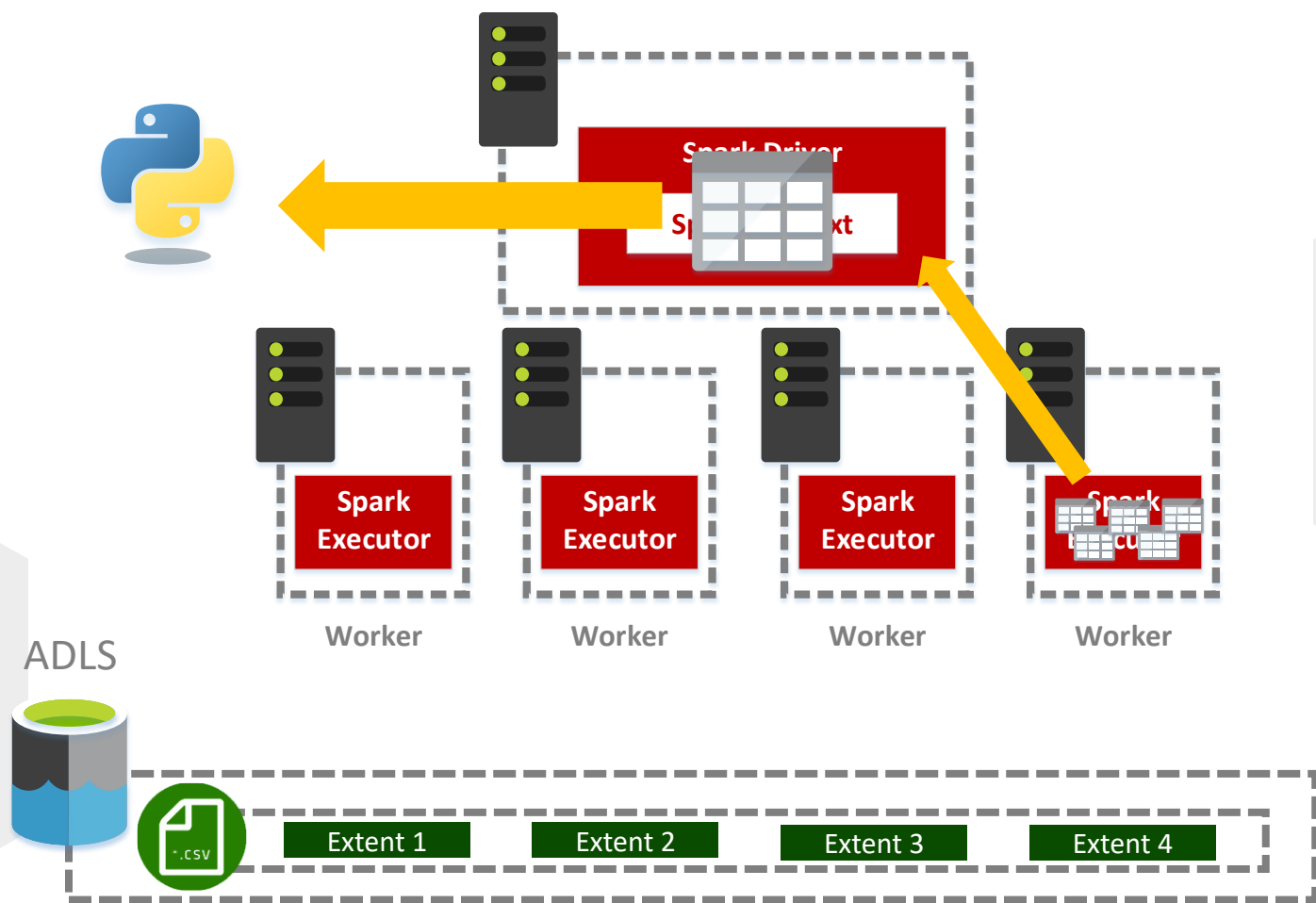


# DISTRIBUTED COMPUTE





# DISTRIBUTED COMPUTE





# TODAY'S COMPETITORS



VS



# Databricks

- Released 2016 (AWS)
- AWS/Azure Cross Platform
- Databricks Proprietary Runtime
- Built by the inventors of Spark

## Special Skills:

- Workspace Features
- Delta Engine

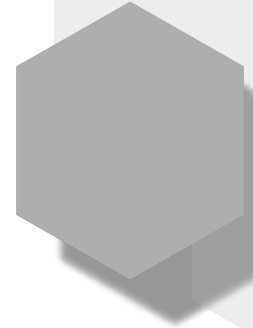


# Synapse Analytics

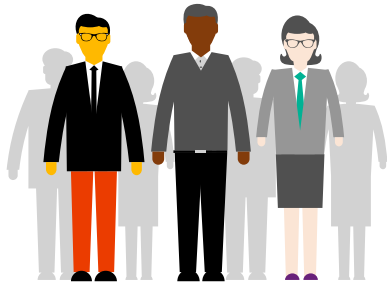
- Still in Preview
- Azure Only
- Vanilla Spark Runtime
- Fresh look at how Spark can work

## Special Skills:

- Integrations
- Spark.NET



# AZURE DATABRICKS WORKSPACE



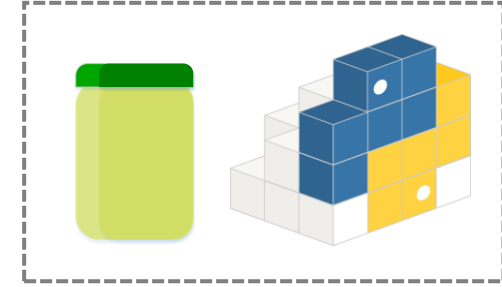
User  
Management



Notebooks



Jobs

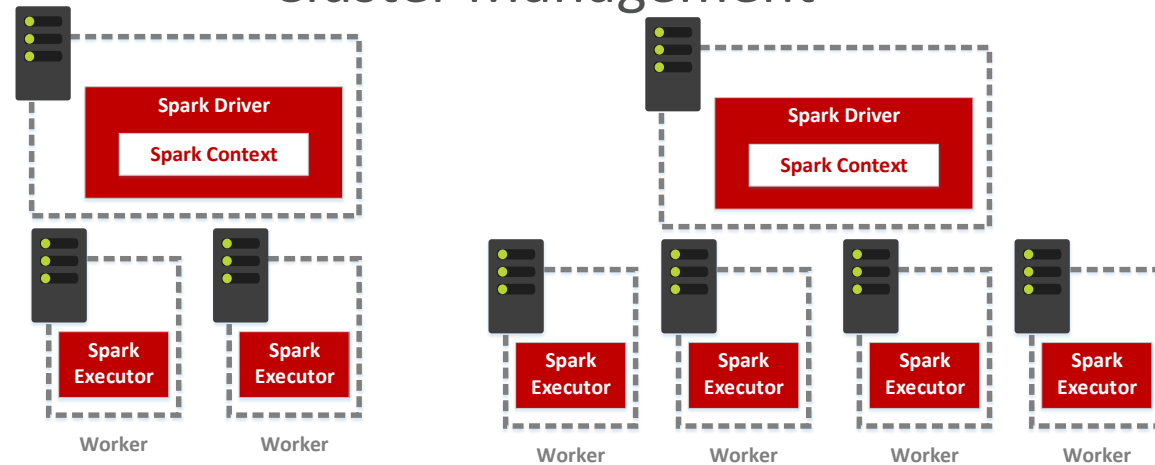


Library

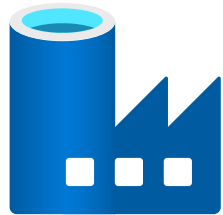
DBFS



Cluster Management







Azure Data  
Factory V2

Generally  
Available



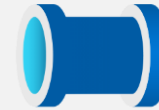
SQL Data  
Warehouse  
(SQLDW)



Azure Synapse  
Analytics  
(SQLDW)



Azure Synapse  
Analytics  
(Workspaces)



Pipelines



Provisioned Spark  
Pools

NEW



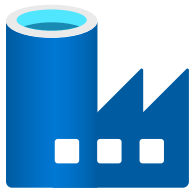
On Demand SQL  
Pools

NEW



Provisioned SQL  
Pools  
(aka SQLDW)

Public Preview



Azure  
Data Factory



ADF Mapping  
Data Flows



Azure Synapse  
Studio



Monitoring



Management



Provisioned  
SQL Pools (*SQLDW*)



On Demand  
SQL Pools



Provisioned  
Spark Pools



Data Lake  
Store Gen 2



Metadata  
Store







Azure  
Data Factory



ADF Mapping  
Data Flows



Azure Synapse  
Studio



Monitoring



Management



Provisioned  
SQL Pools (*SQLDW*)



On Demand  
SQL Pools



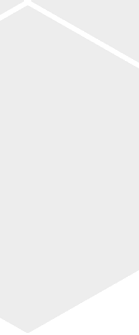
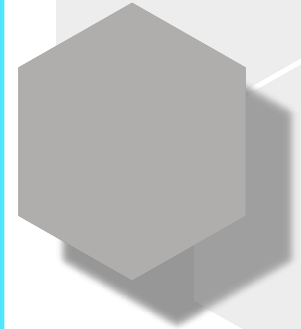
Provisioned  
Spark Pools



Data Lake  
Store Gen 2



Metadata  
Store





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

# ROUND ONE - POWER



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS

## Synapse Workspace



Synapse Runtime



## Databricks Workspace

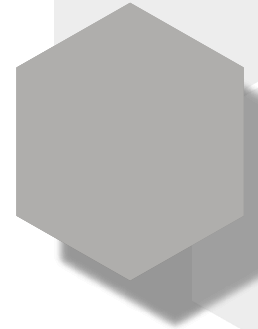


DBX Runtime



**Why do we care?**

- Functions
- Optimisations
- Spark 3.0



Delta Engine  
& Photon

3.0 Optimisation

Spark 3.0

Databricks

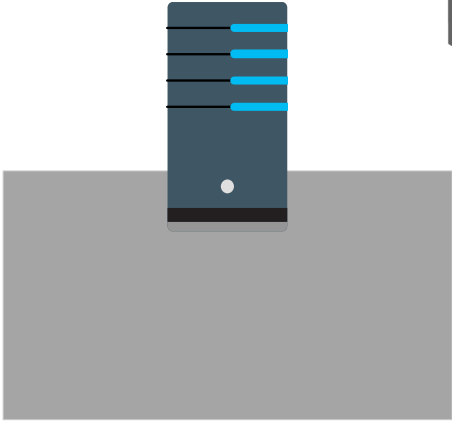
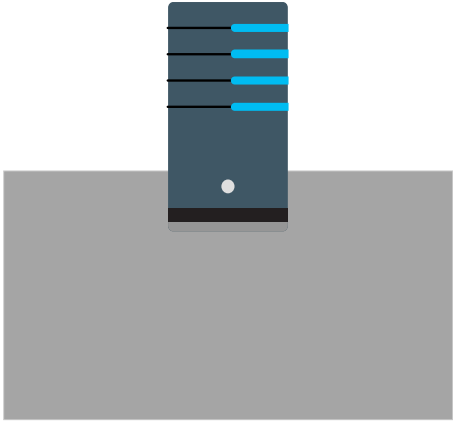
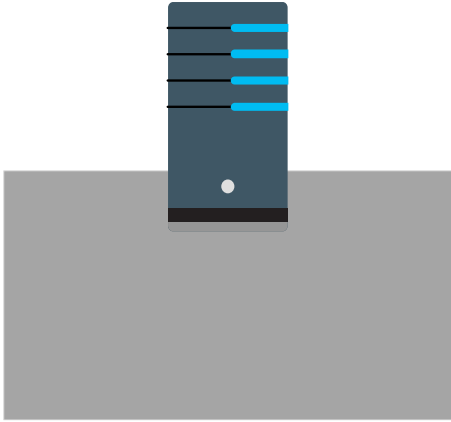
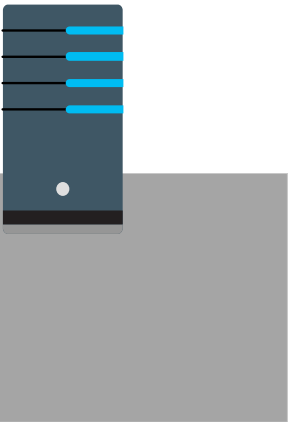
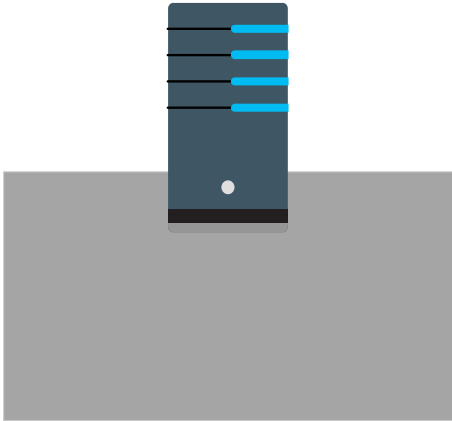
Spark 3.0

Synapse



# PHYSICAL ARCHITECTURE

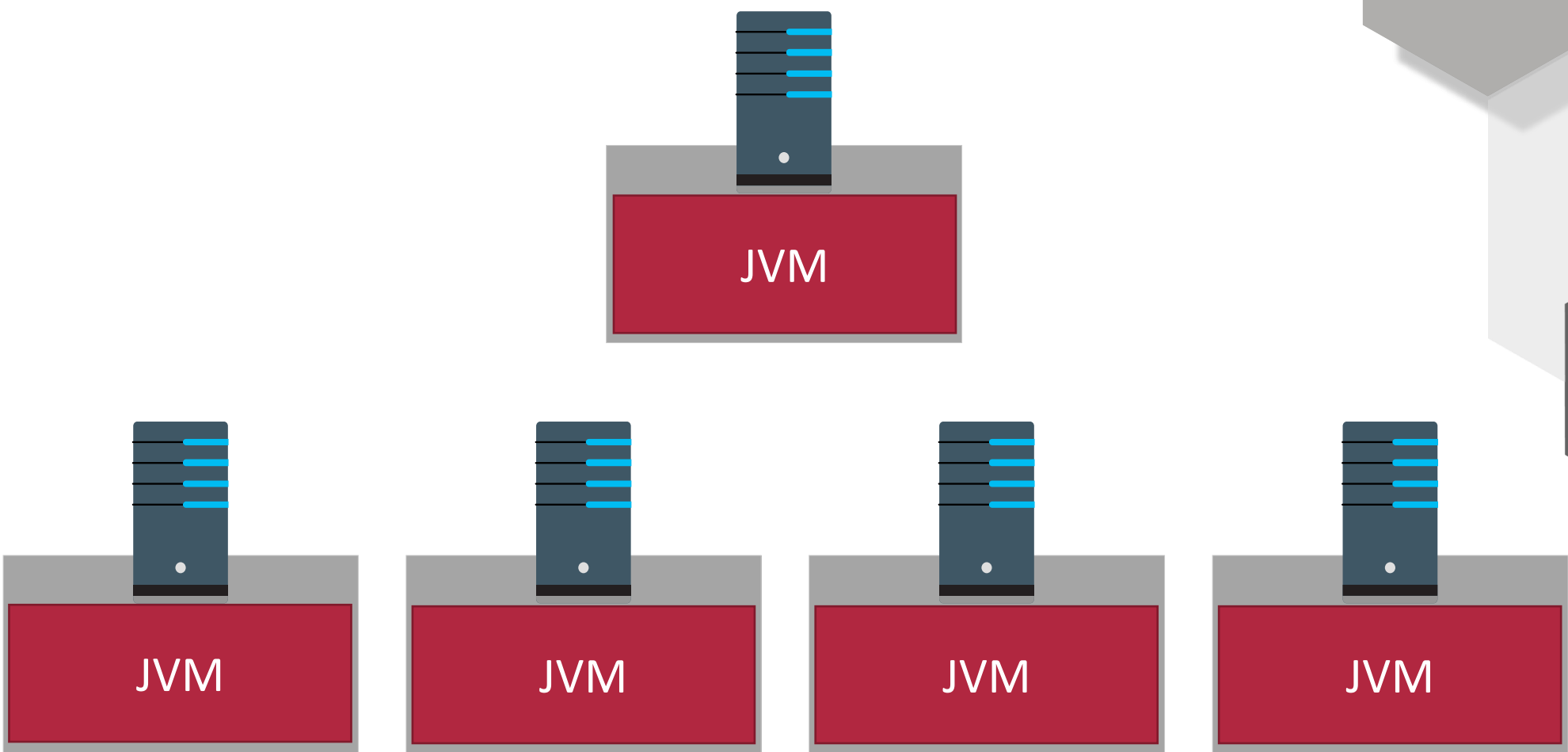
Driver



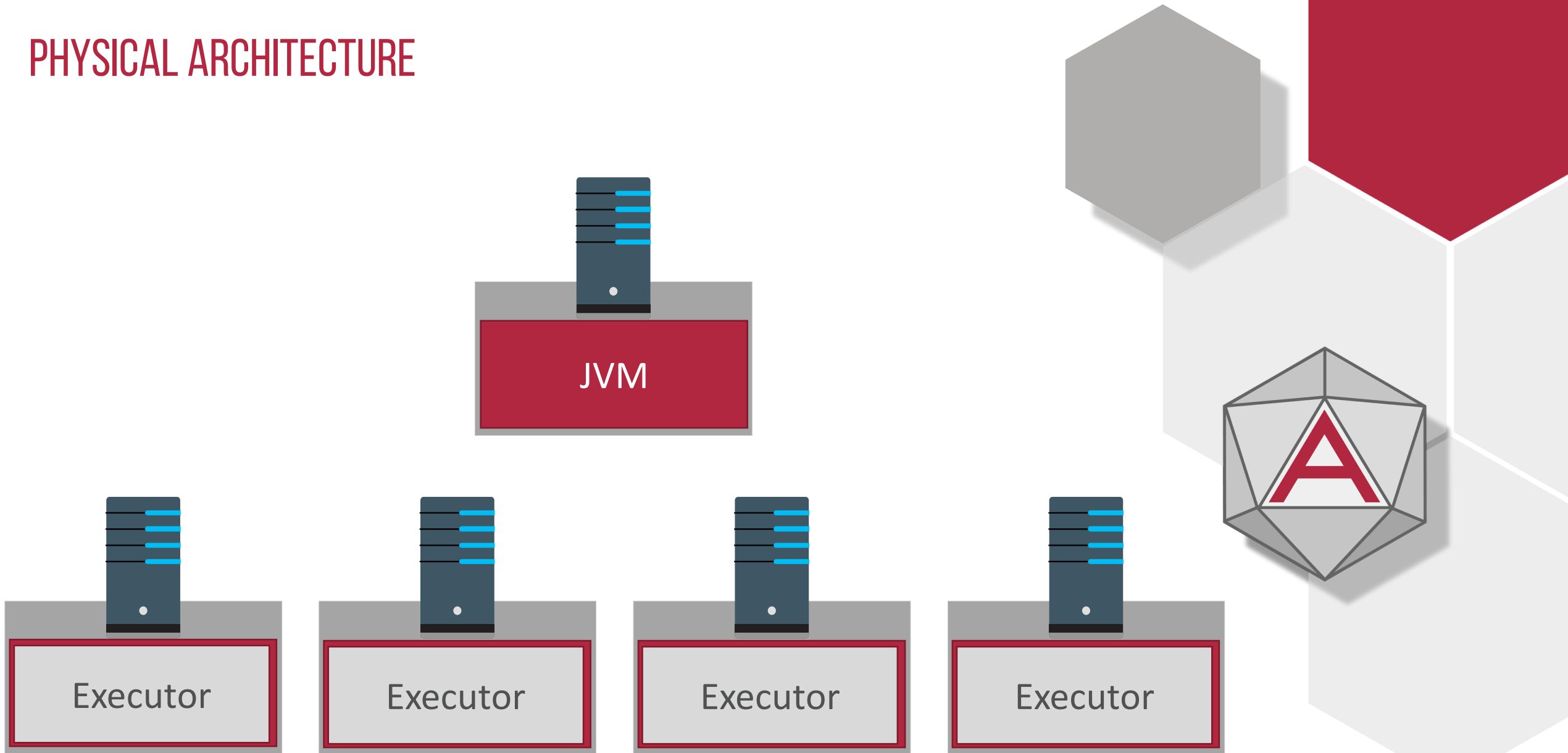
Workers



# PHYSICAL ARCHITECTURE

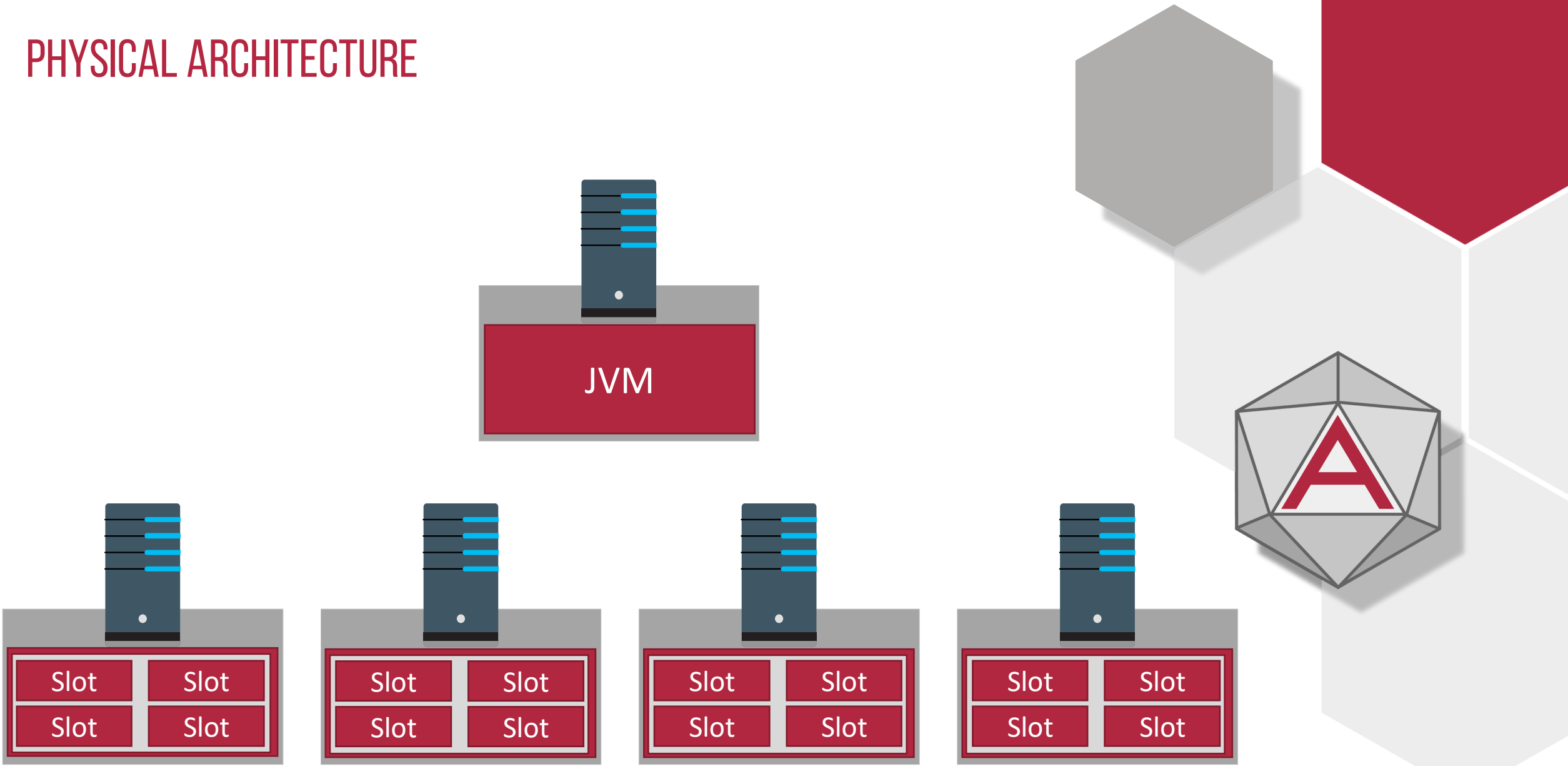


# PHYSICAL ARCHITECTURE



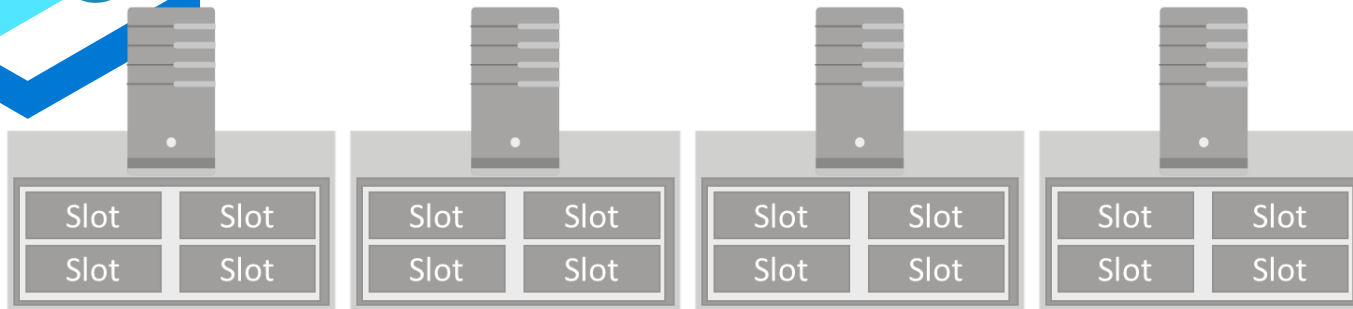
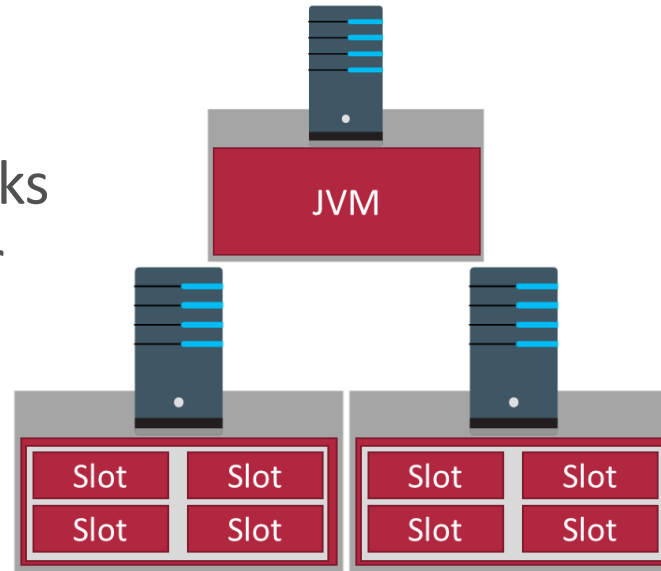


# PHYSICAL ARCHITECTURE





Databricks  
Cluster

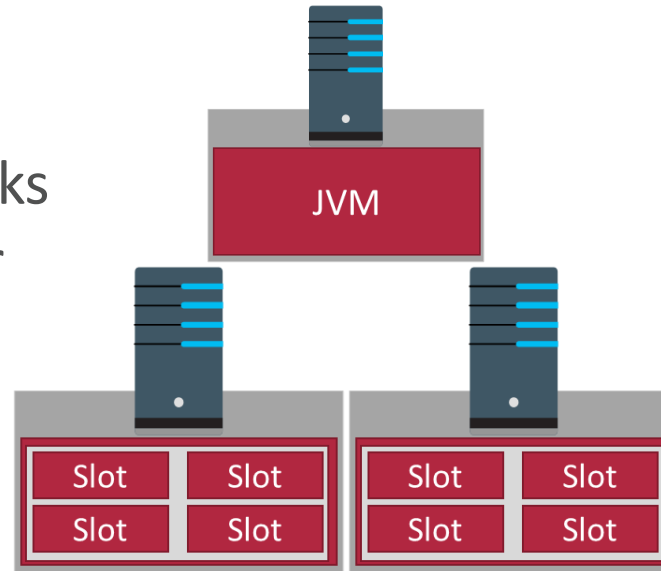


Spark Pool

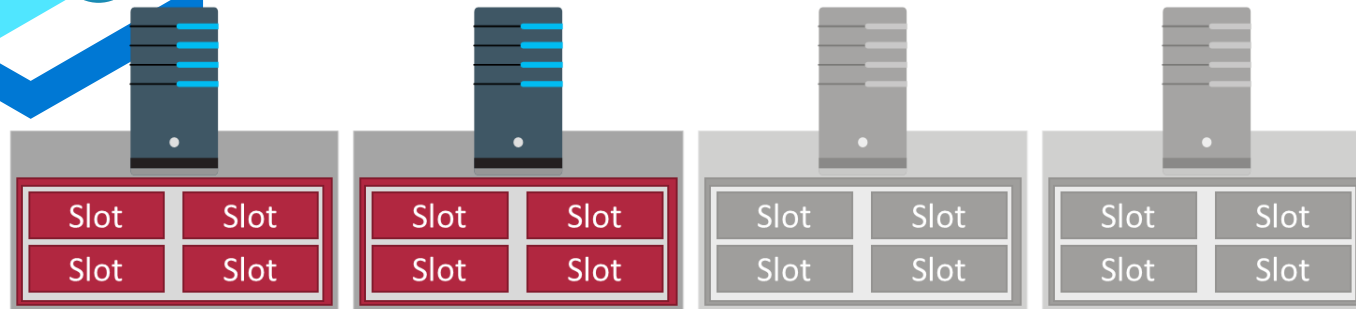




Databricks  
Cluster



Session

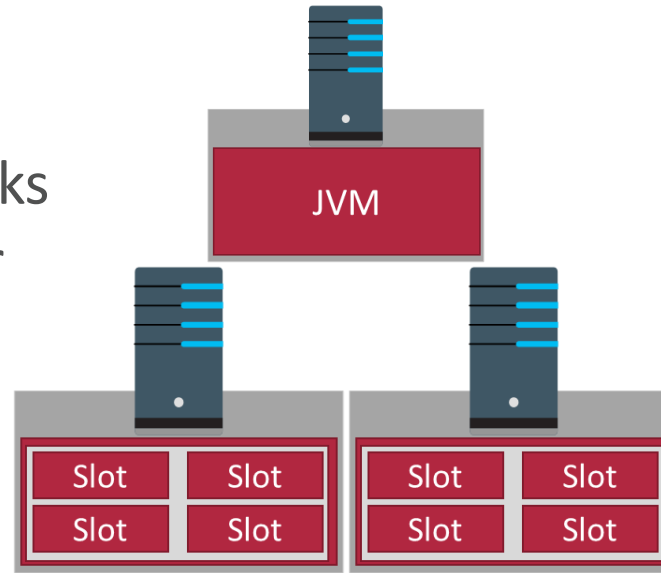


Spark Pool

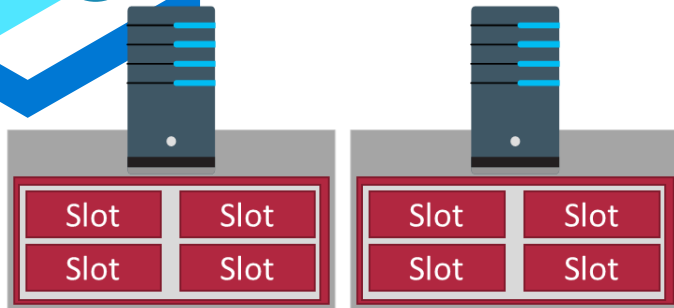




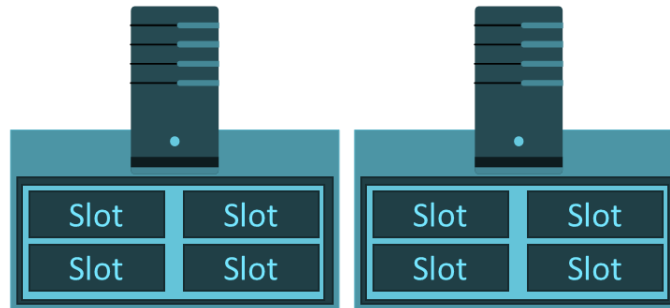
Databricks  
Cluster



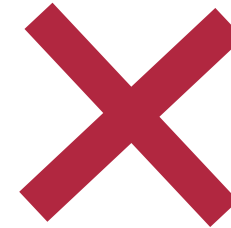
Session 1



Session 2



Session 3



Spark Pool





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

# ROUND TWO - SPECIAL MOVES



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS

# DATABRICKS NOTEBOOKS

Databricks



DBUtils (*Databricks Utilities*)

- Libraries
- File Management
- Secrets
- Widgets

Display / Charting

Version Control

More Natural UI Integration



# SYNAPSE INTEGRATIONS



New SQL script ▾ New data flow New dataset Upload Download + New folder

← → ▾ ↑ root > BASE > Parquet > Adventureworks > SalesLT

NAME

Address

Customer

CustomerAddress

Product

ProductCategory

ProductDescription

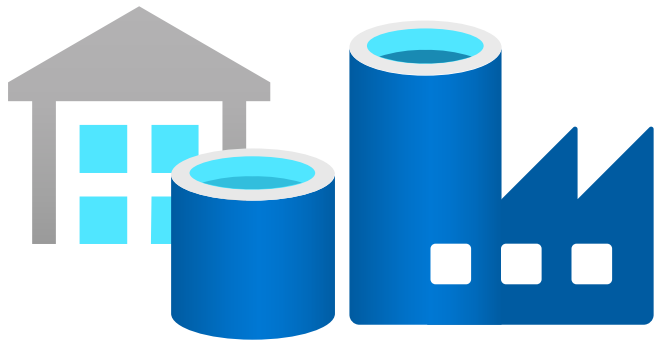
ProductModel

ProductModelProductDescription

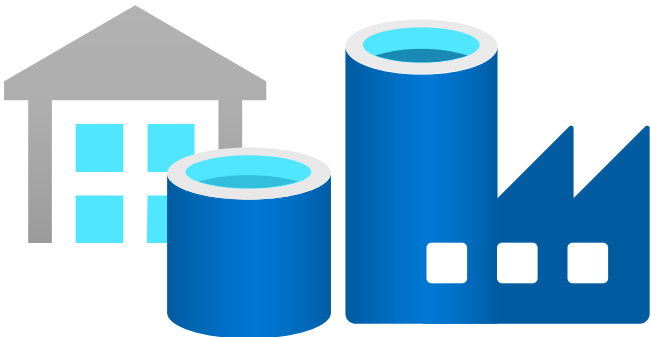
SalesOrderDetail



# SYNAPSE INTEGRATIONS



# SYNAPSE INTEGRATIONS



▲ Cosmos DB	1
▲ 🔗 OnlineSalesHTAP (OnlineSales)	
▶ 📄 Products	
▲ 📄 Sales	



## Hive & Metadata

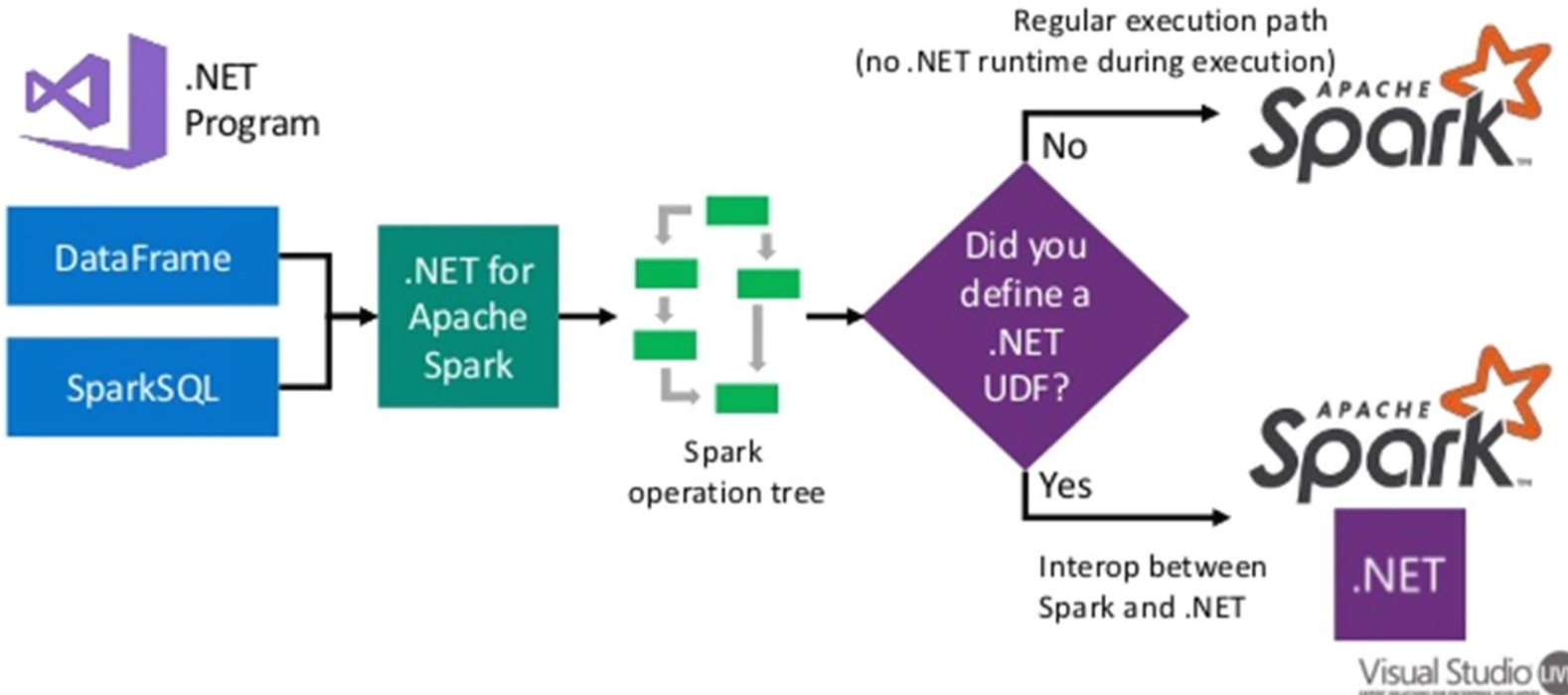


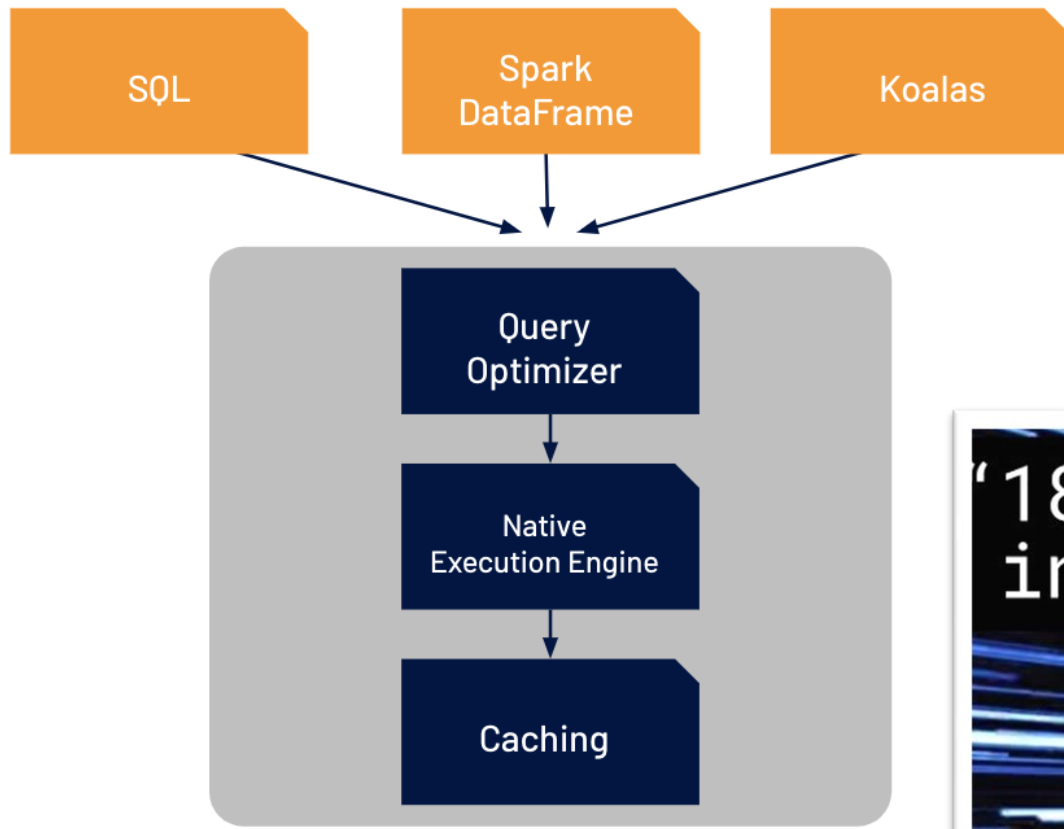
SQL On-Demand





## What is happening when you write .NET Spark code?





**"18x increased performance in star schema workloads"**

# *DELTA ENGINE*

#SparkAISummit

The promotional banner features a dark background with a blue starburst or light trail effect radiating from the center. The Delta logo is positioned in the top right corner. The text is presented in a clean, sans-serif font, with the product name in a large, italicized, yellow font.





# HYPERSPACE

An extensible indexing subsystem for Apache Spark™

Scala

Python

.NET



## Our first hyperspace: the covering index

Table A

a	b	c
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

Without  
Indexes

Table B

a	p	q
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

```
SELECT b, c  
FROM Table A, B  
JOIN ON A.a = B.a
```

With Covering  
Indexes

Step 1: Shuffle  
(data is not sorted)

a	b	c
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

Table A

Table B

Step 2: Sort both sides

a	b	c
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

Table A

Table B

Step 3: Merge

a	p	q
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

Result

Step 1: Optimizer picks index  
(pre-shuffled, pre-sorted)

a	b	c
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

Idx A

Idx B

Step 2: Merge

a	p	q
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9
10	10	10

Result



**Shuffle eliminated**  
Since shuffle is the most  
expensive step, this query  
might run faster at scale



mlflow



**DELTA LAKE**



**DELTA LAKE**







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

# ROUND THREE - PRICE



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS

## Apache Spark Pools in Azure Synapse

Memory Optimised

NODE SIZE:

Small (4 vCPU(s)/32 GB)

2

×

730

Hours

×

£0.13

=

£774.76

Instances

Per vCore-hour

INSTANCE:

L4s: 4 Core(s), 32 GB RAM, 1 Databricks Unit(s), £0.270/hour

3

×

730

Hours

Virtual machines

=

£590.87

Average per month  
(£0.00 charged upfront)

DBU (Databricks Unit) ⓘ

3.00

×

£0.410

×

730

Hours

DBU

Per DBU per hour

=

£897.72

Upfront cost

£0.00

Monthly cost

£1,488.59

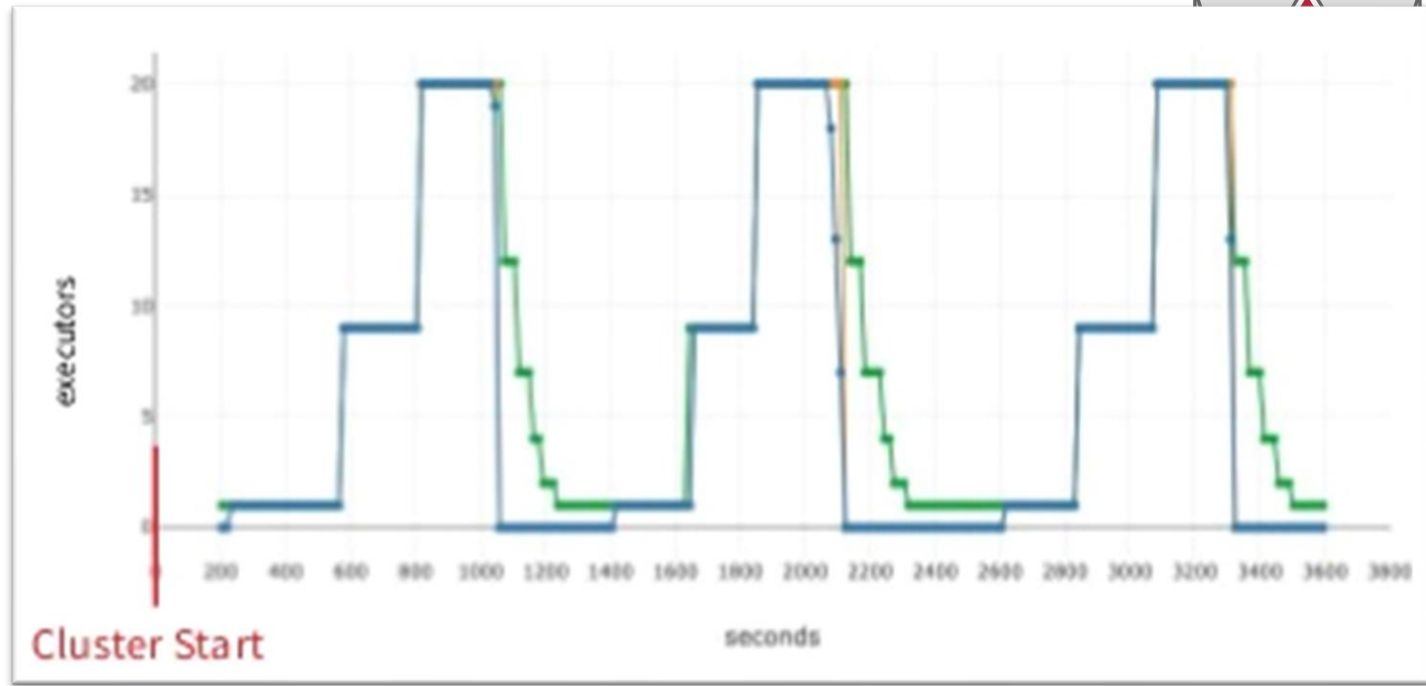


DBX costs more,  
but is more efficient  
for parallel use &  
better optimised





Databricks Premium has  
advanced **Autoscale** feature



How Synapse will handle session-based autoscaling  
or high-concurrency workloads is still unknown



AND THE WINNER IS...



## FOR PURE SPARK PERFORMANCE



For large volume, high-performance or high-concurrency spark usage, Databricks is a more mature, capable offering



But that's like comparing SQL Server **Enterprise** to SQL Server **Standard**.

Of course Enterprise does more, performs better, has additional features - it's the Premium offering



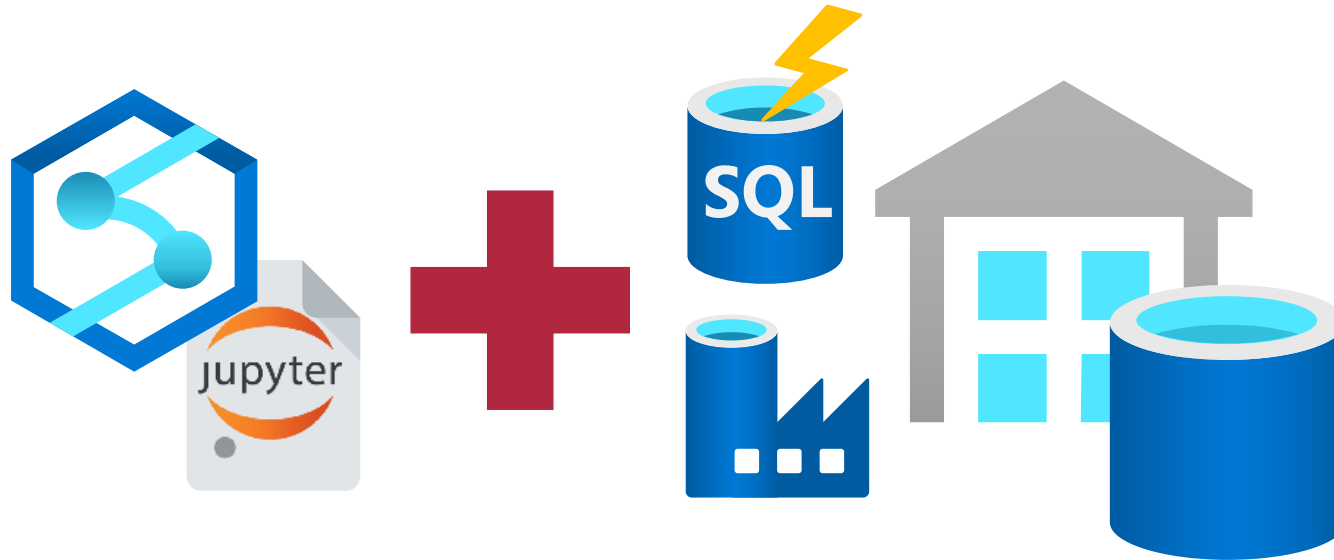
Premium Azure Spark

Standard Azure Spark



## SCENARIOS:

You're building a platform that is largely SQL-based, but may have one or two Spark edge-cases



**Synapse Spark Pools** will be far easier and less complex to set up



## SCENARIOS:

You're a Spark shop with a complex spark application and looking for a place to house it in the cloud



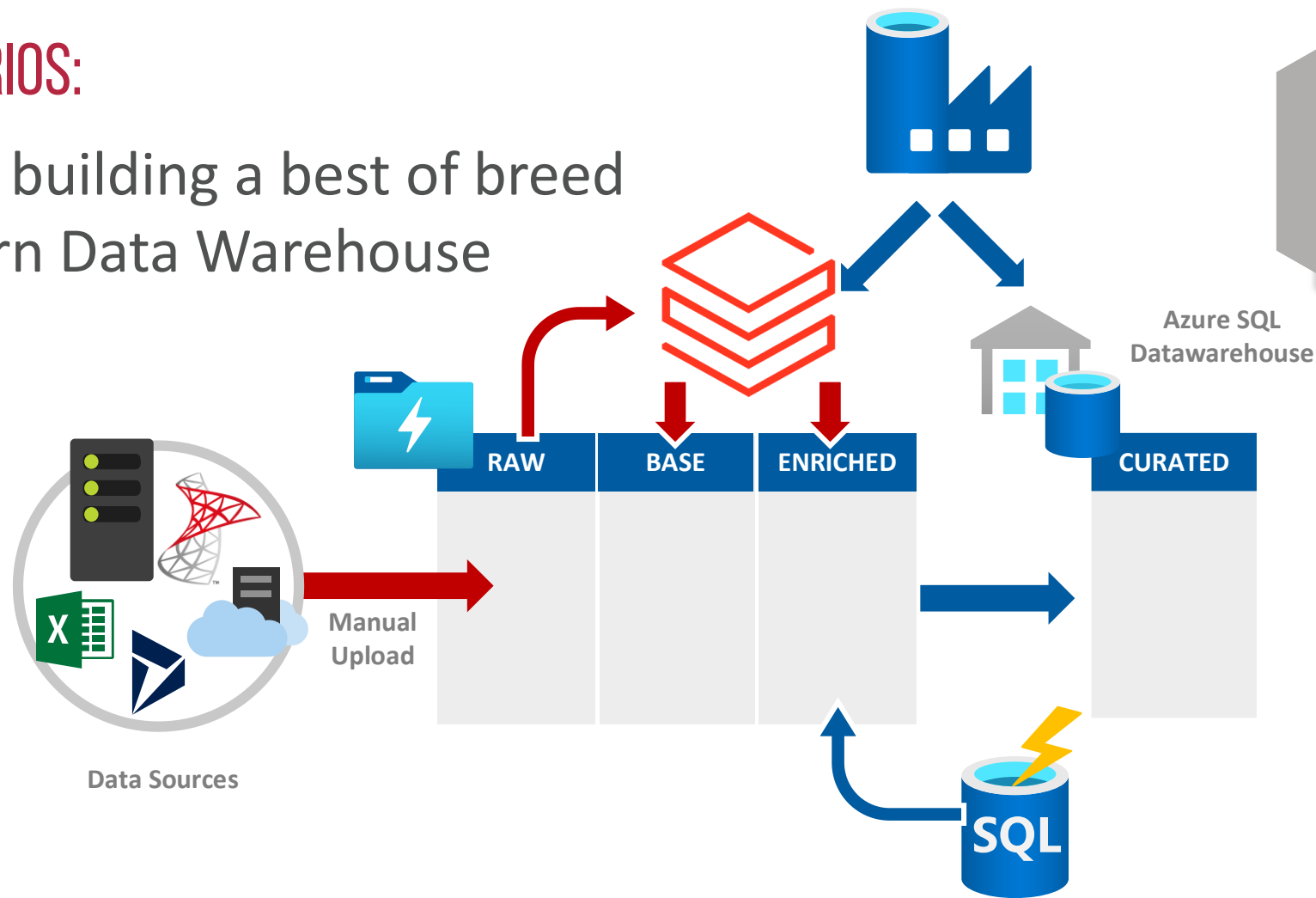
**Databricks** will have more features and provide better performance optimisations





## SCENARIOS:

You're building a best of breed  
Modern Data Warehouse



**Databricks** and **Synapse Analytics** work **Better Together**