

Python Pipeline Primer

ETL in Azure

Simon Whiteley | Adatis

10/10/2018



Gold Data Analytics
Gold Data Platform
Gold Cloud Platform



Agenda

Why Big Data
processing is
needed in the
cloud

What is
Databricks?

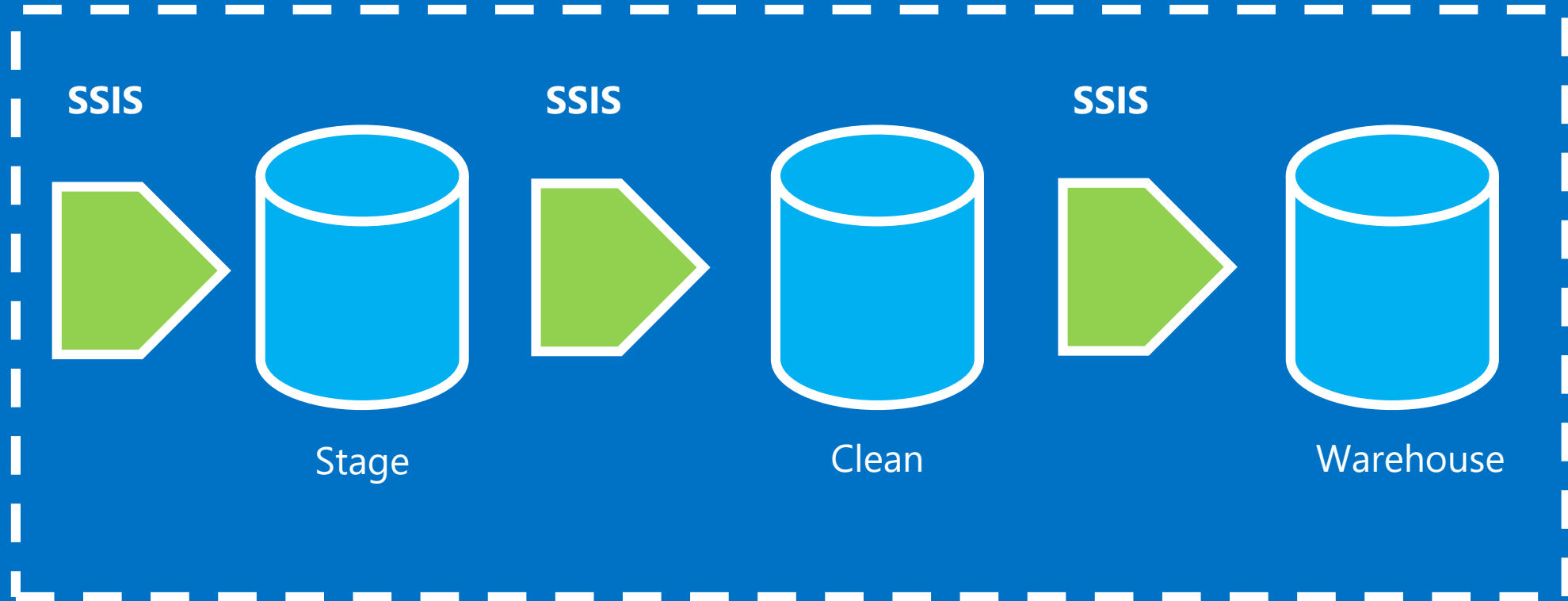
Patterns &
Implementation

Orchestration

Data Factory
Dataflows

WBDPINITC

One-Box SQL BI Architecture



On-Prem SQL Server



"Big Data" Solutions



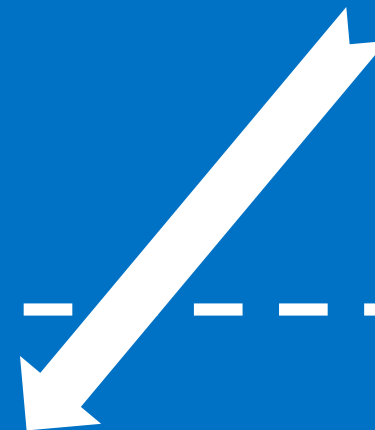
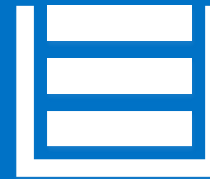
Technical Barriers



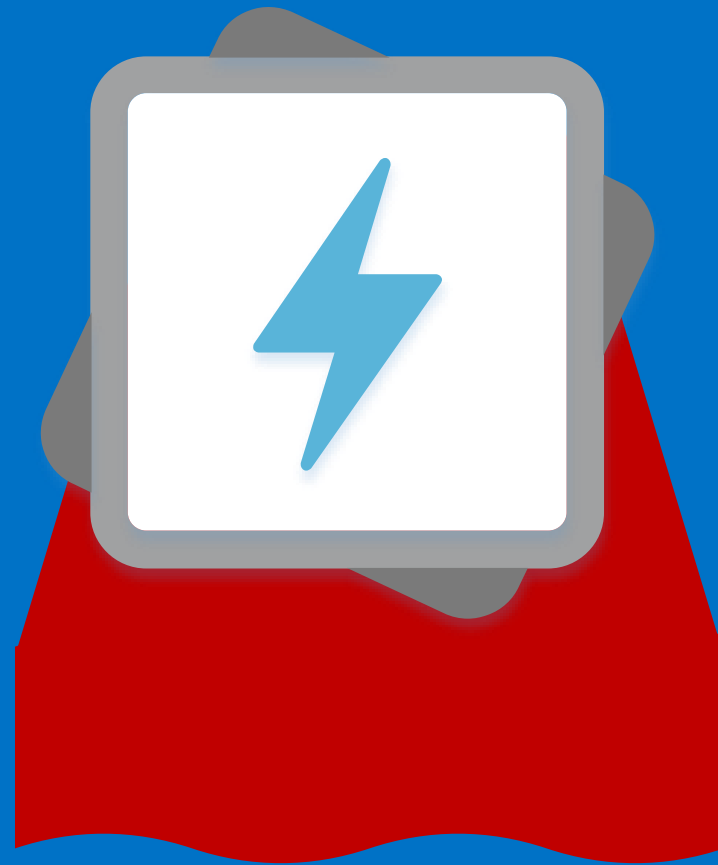
On-Prem SQL Server



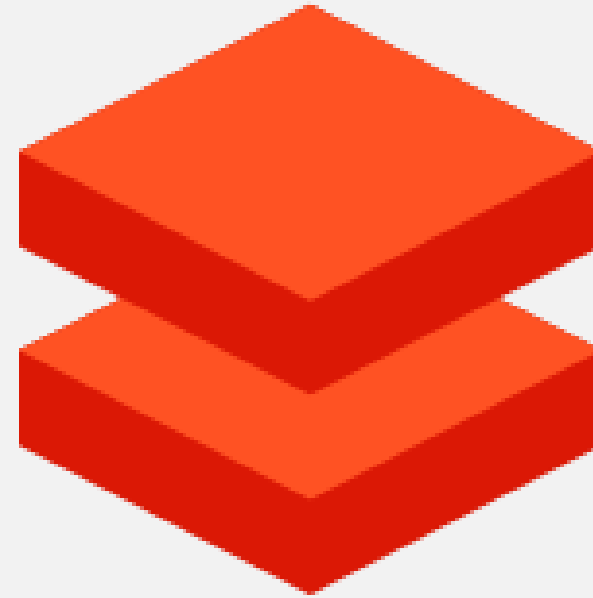
"Big Data" Solutions



Modern Analytics Platform



DATA LAKE ANALYTICS



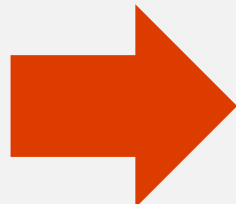
Azure
Databricks

Databricks?



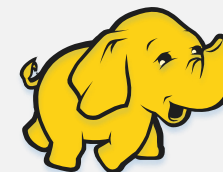
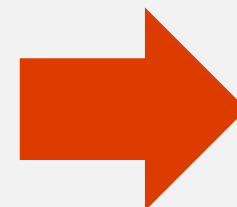
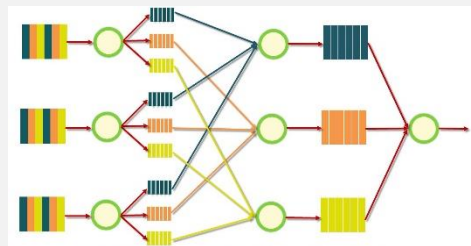
Google File System Papers
Released

2003



Google MapReduce Papers

2004



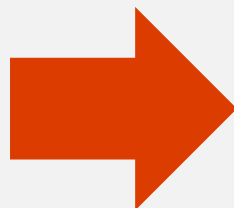
2006

Apache Hadoop
project created



Matei Zaharia starts Spark
project

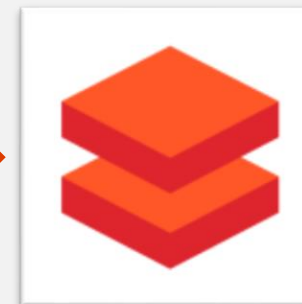
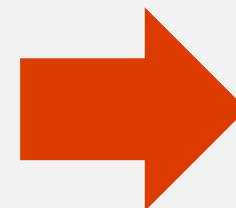
2012



THE
APACHE
SOFTWARE FOUNDATION

Project donated to Apache
Foundation

2013



Databricks founded by
Matei

2013



2016

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

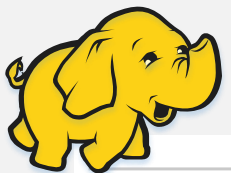


Microsoft
Azure

2018

So What?

- Most up to date Spark optimisations
- Doesn't need specialist hardware
- Quicker than traditional MapReduce
- Cluster Management, Notebooks, Jobs...



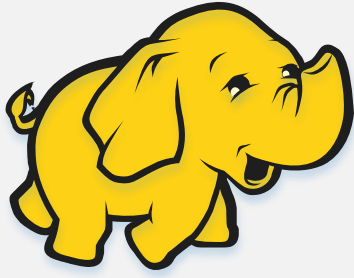
HDInsight

| INSTANCE | CPU | RAM | OS | HDINSIGHT PRICE | TOTAL PRICE++ |
|----------|-----|-------|-------------|-----------------|---------------|
| D3 v2 | 4 | 14 GB | £0.171/hour | £0.05/hour | £0.22/hour |



Databricks

| INSTANCE | vCPU | RAM | DBU COUNT | LINUX VM PRICE | DBU PRICE | PAY AS YOU GO TOTAL PRICE |
|----------|------|-----------|-----------|----------------|-------------|---------------------------|
| D3 v2 | 4 | 14.00 GiB | 0.75 | £0.203/hour | £0.308/hour | £0.511/hour |



Open Source

20 min provisioning

Integrates Well

Secure

Hadoop, Spark, Kafka,
Hbase, HIVE, Storm...

Slow Release Cycle



Open Source

5 min provisioning

Integrates Well

Secure

Spark (Python/Scala/R)

Fast Release Cycle



Proprietary

1 min provisioning

Integrates Poorly

Secure

U-SQL

Slow Release Cycle



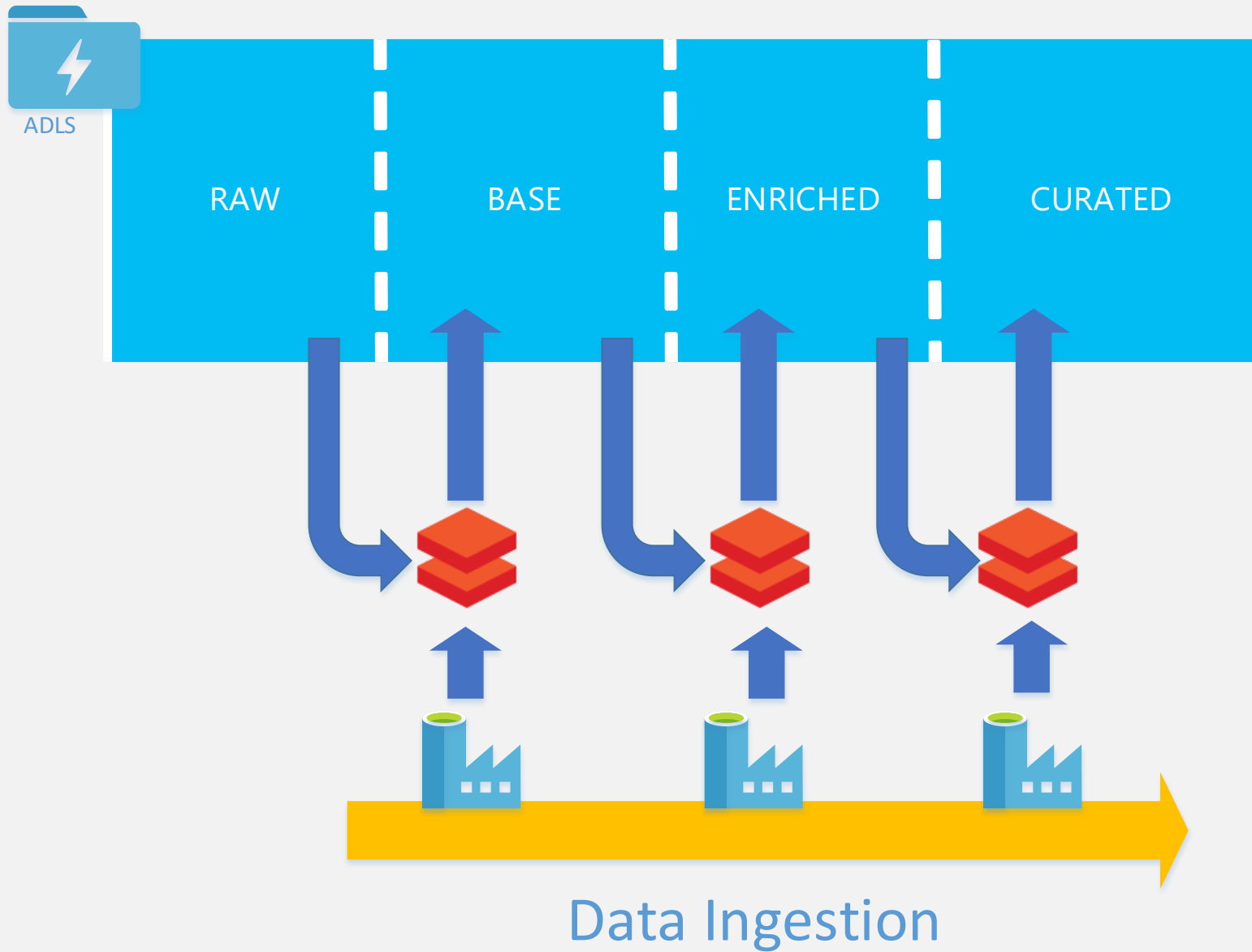
Data
Scientists

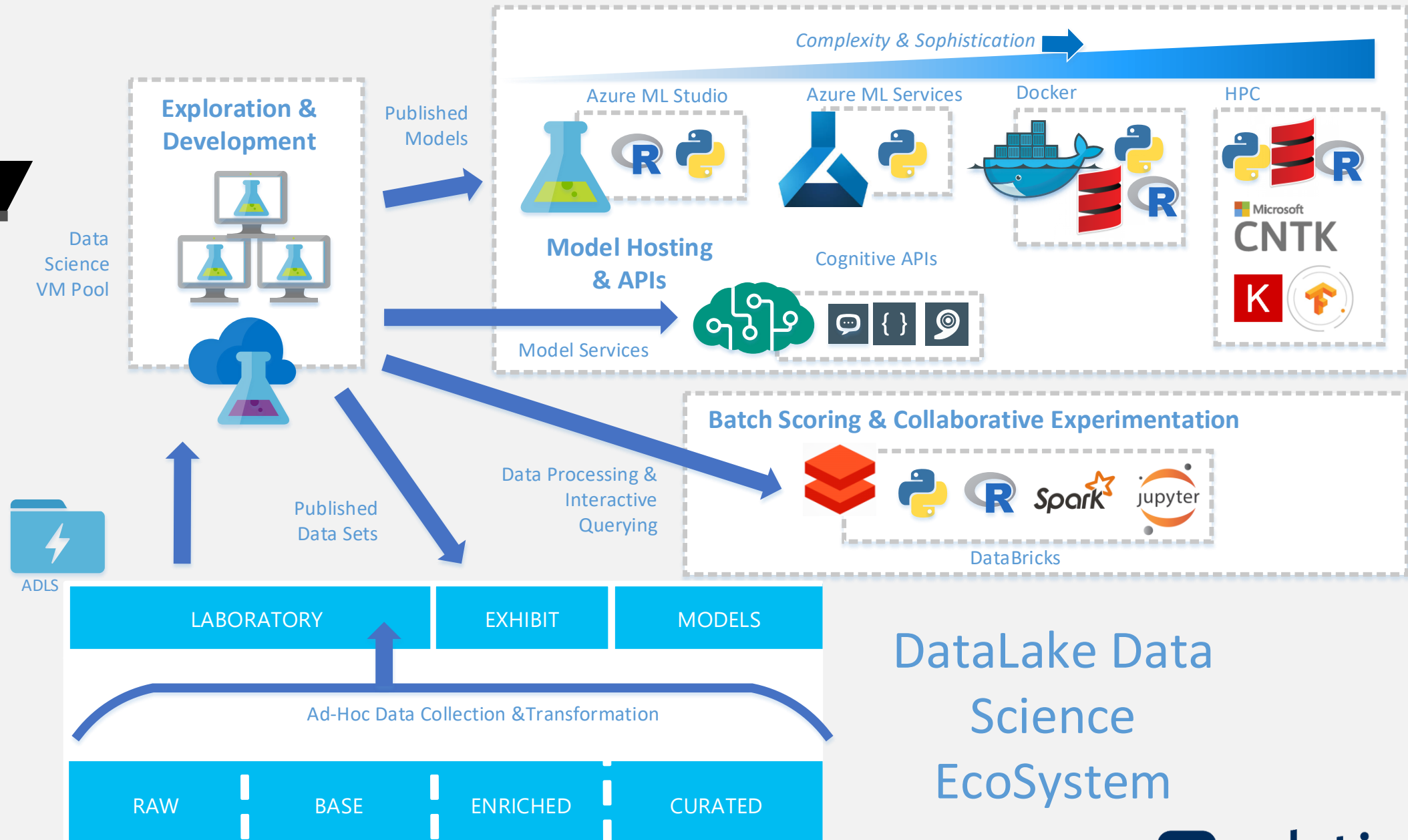


Data
Analysts

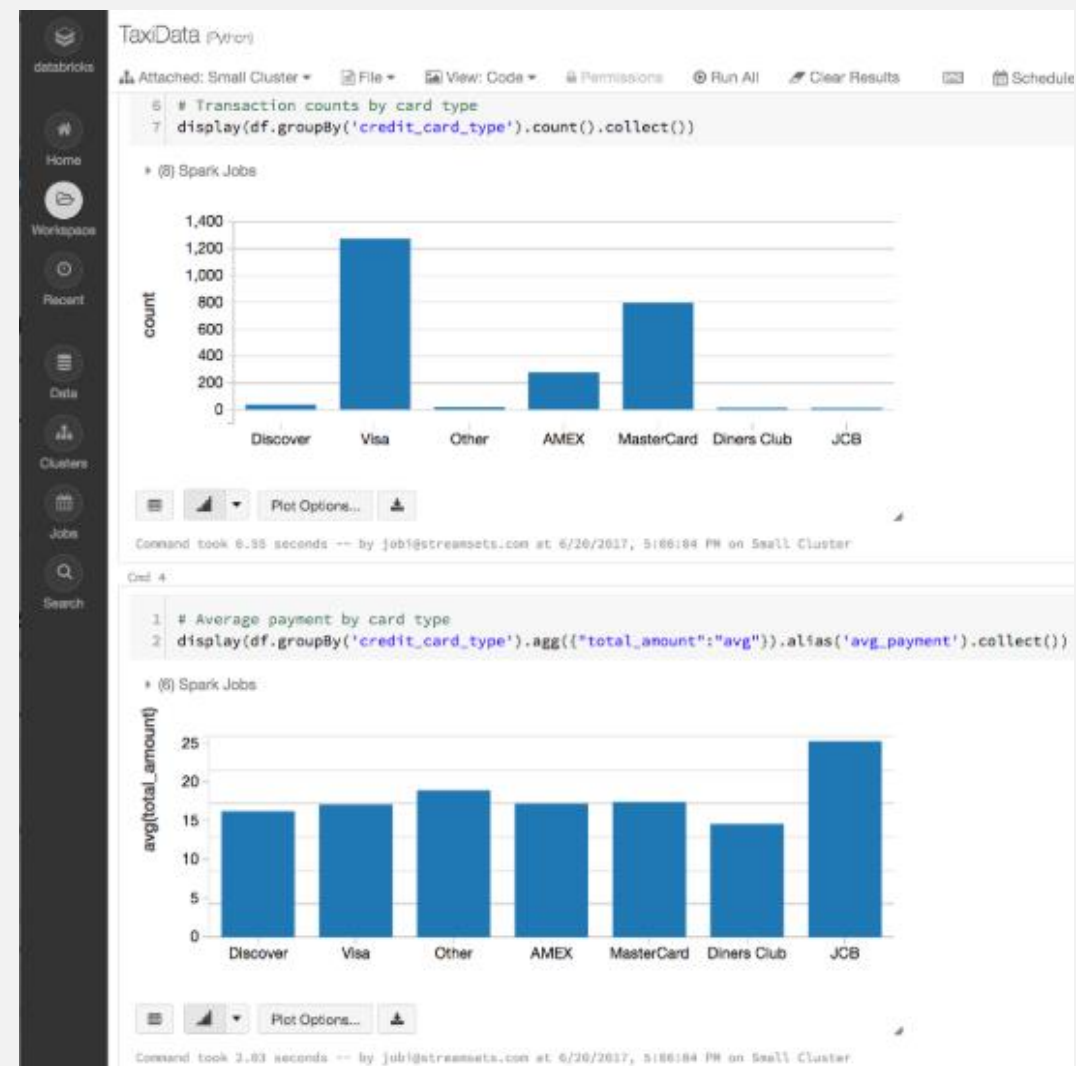
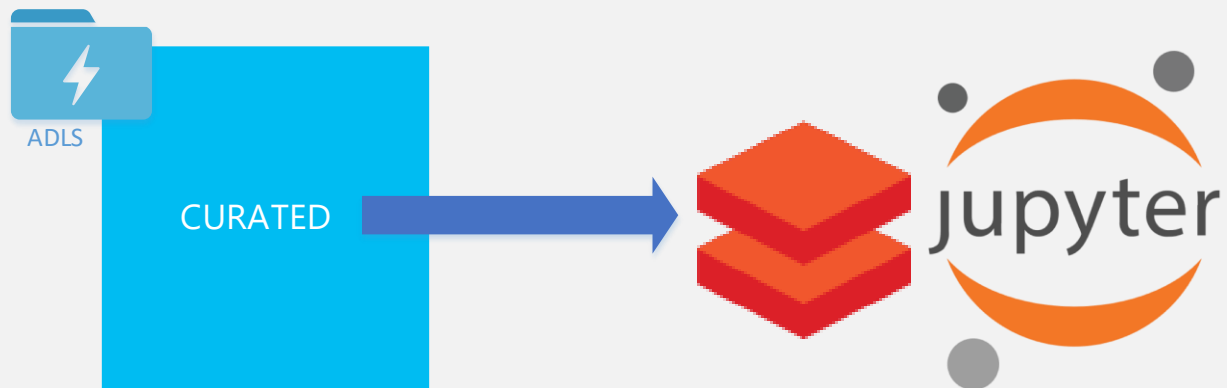
Data
Engineers







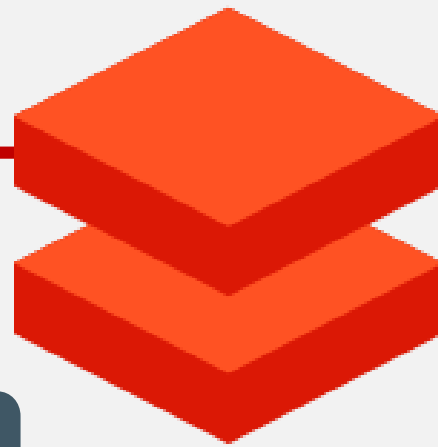
DataLake Data
Science
EcoSystem



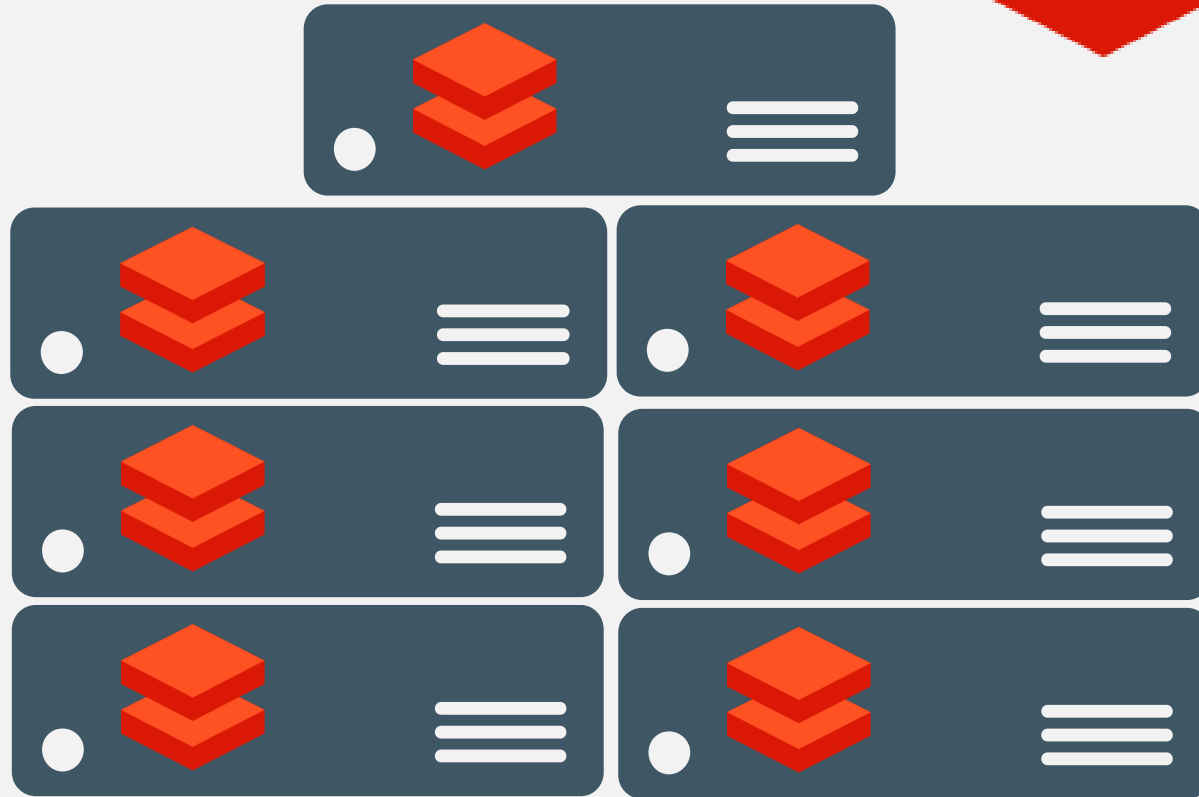
Databricks Basics

Patterns & Implementation

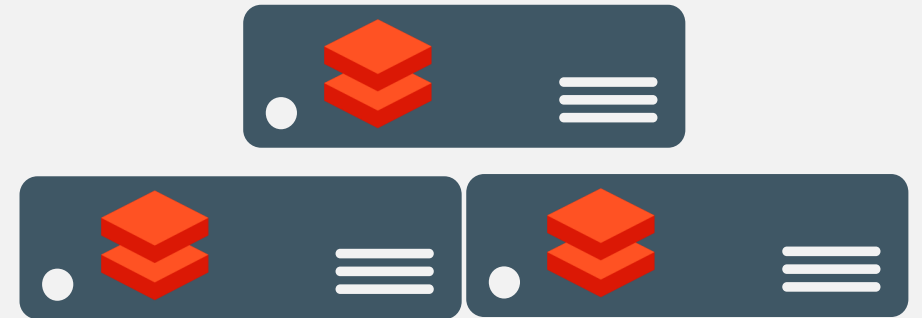
Workload Isolation



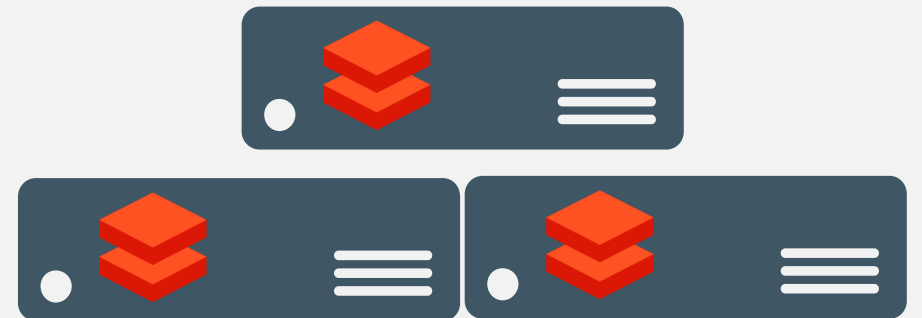
Processing Cluster



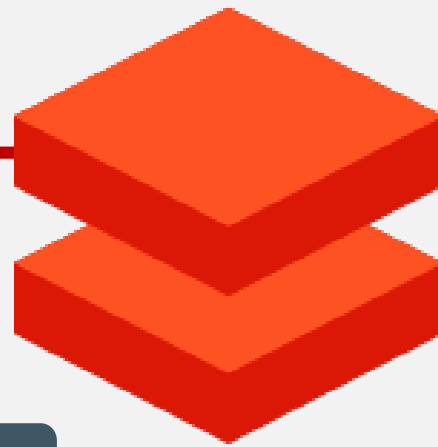
Streaming Cluster



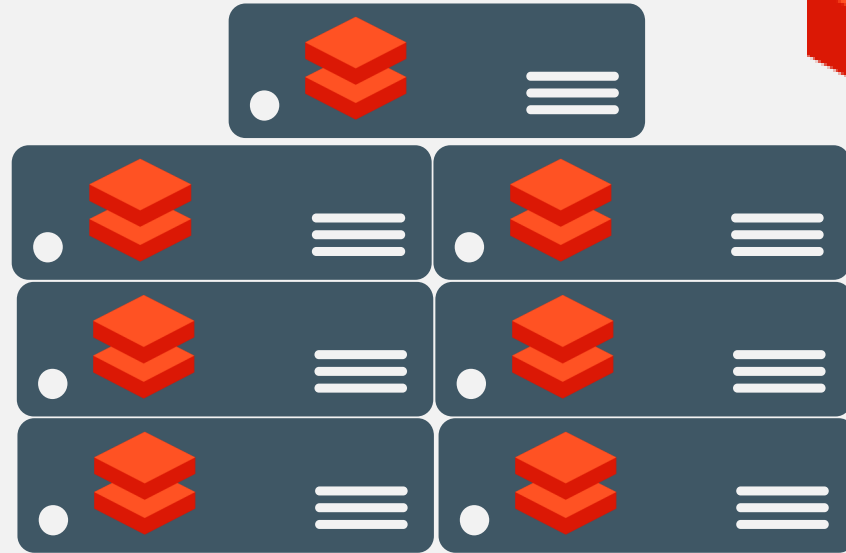
Interactive Cluster



Workload Isolation



Fact Cluster



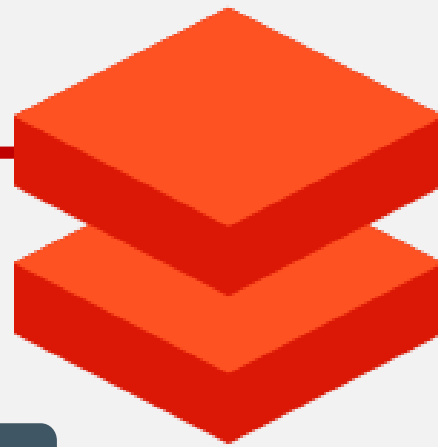
Tensorflow Cluster



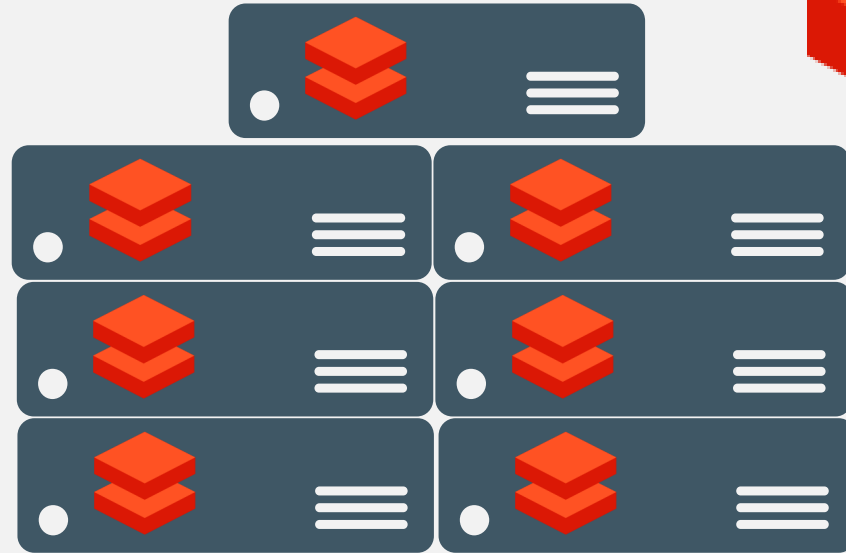
Dim Cluster



Workload Isolation



Fact Cluster



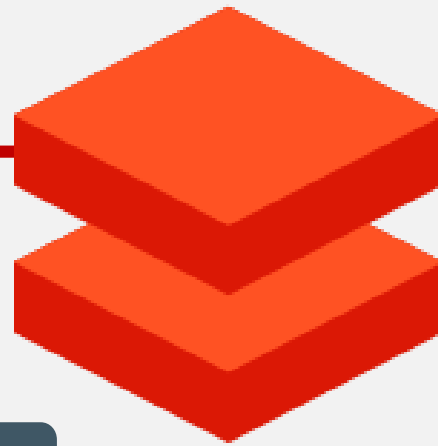
Tensorflow Cluster



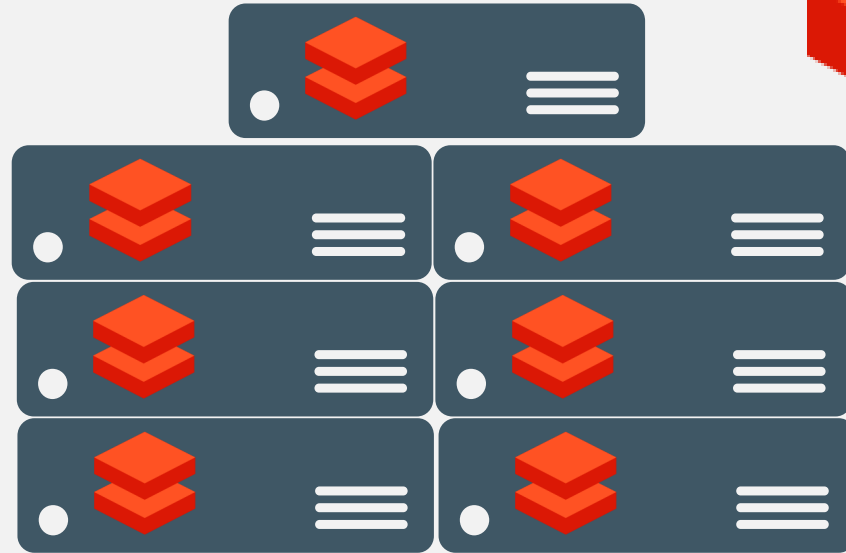
Dim Cluster



Workload Isolation



Fact Cluster

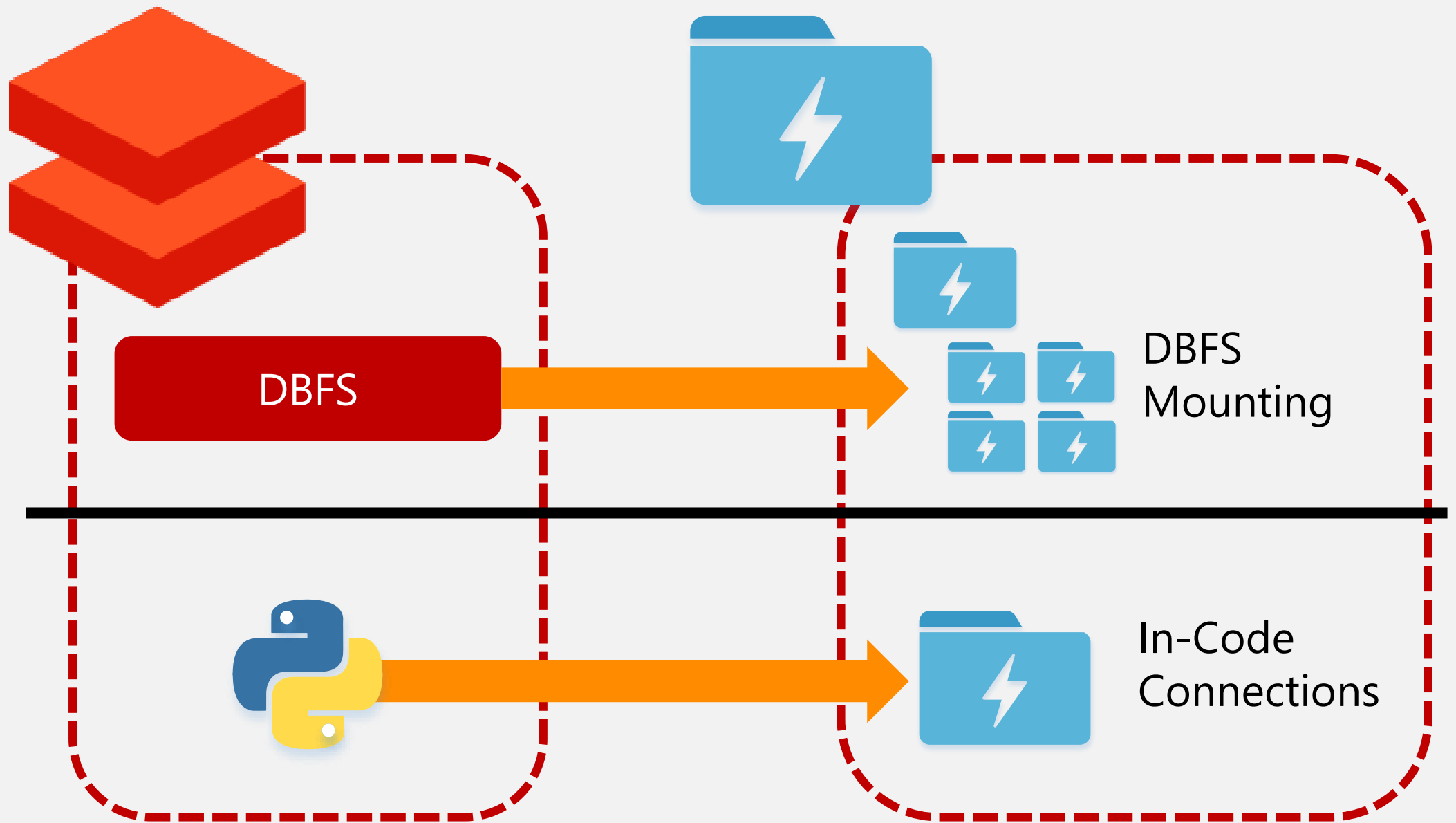


Tensorflow Cluster

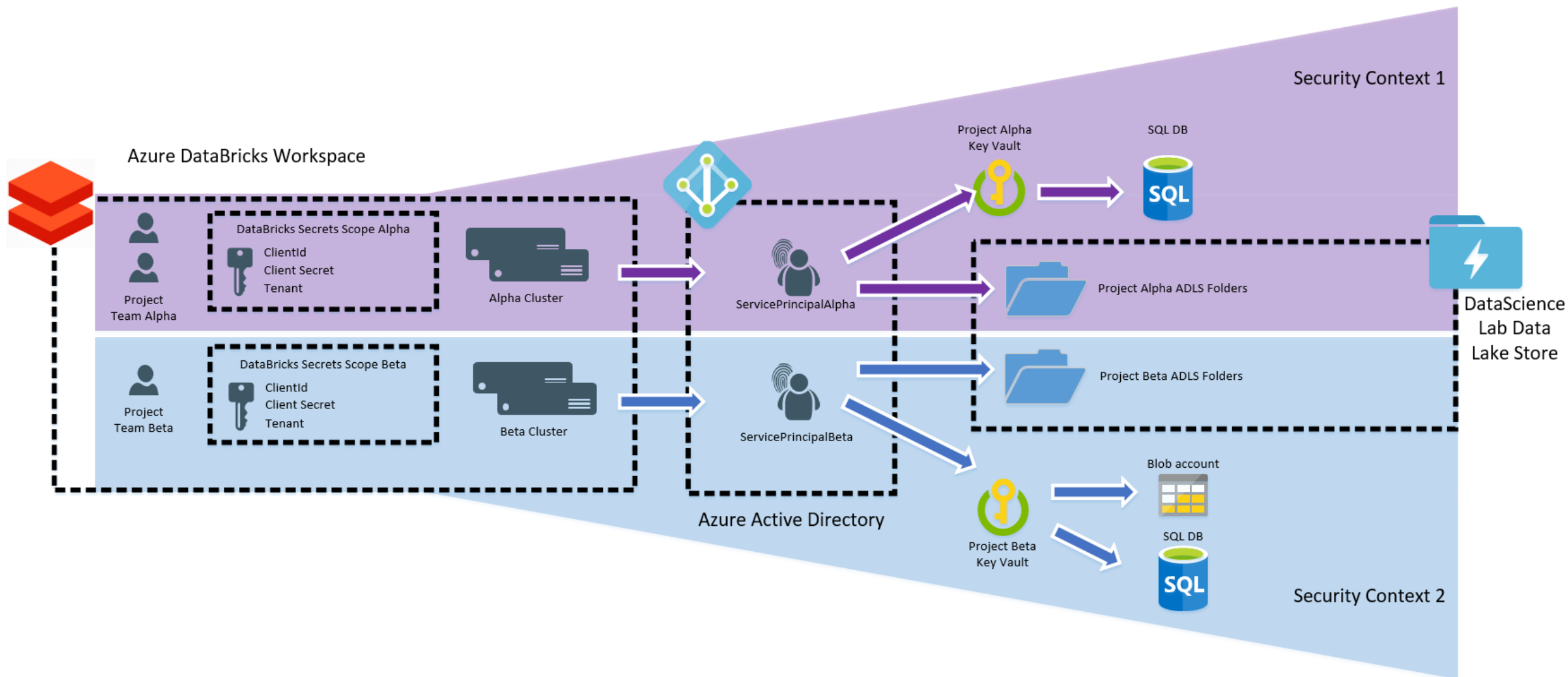


Dim Cluster





Secrets



Standard Clusters Vs Serverless pools

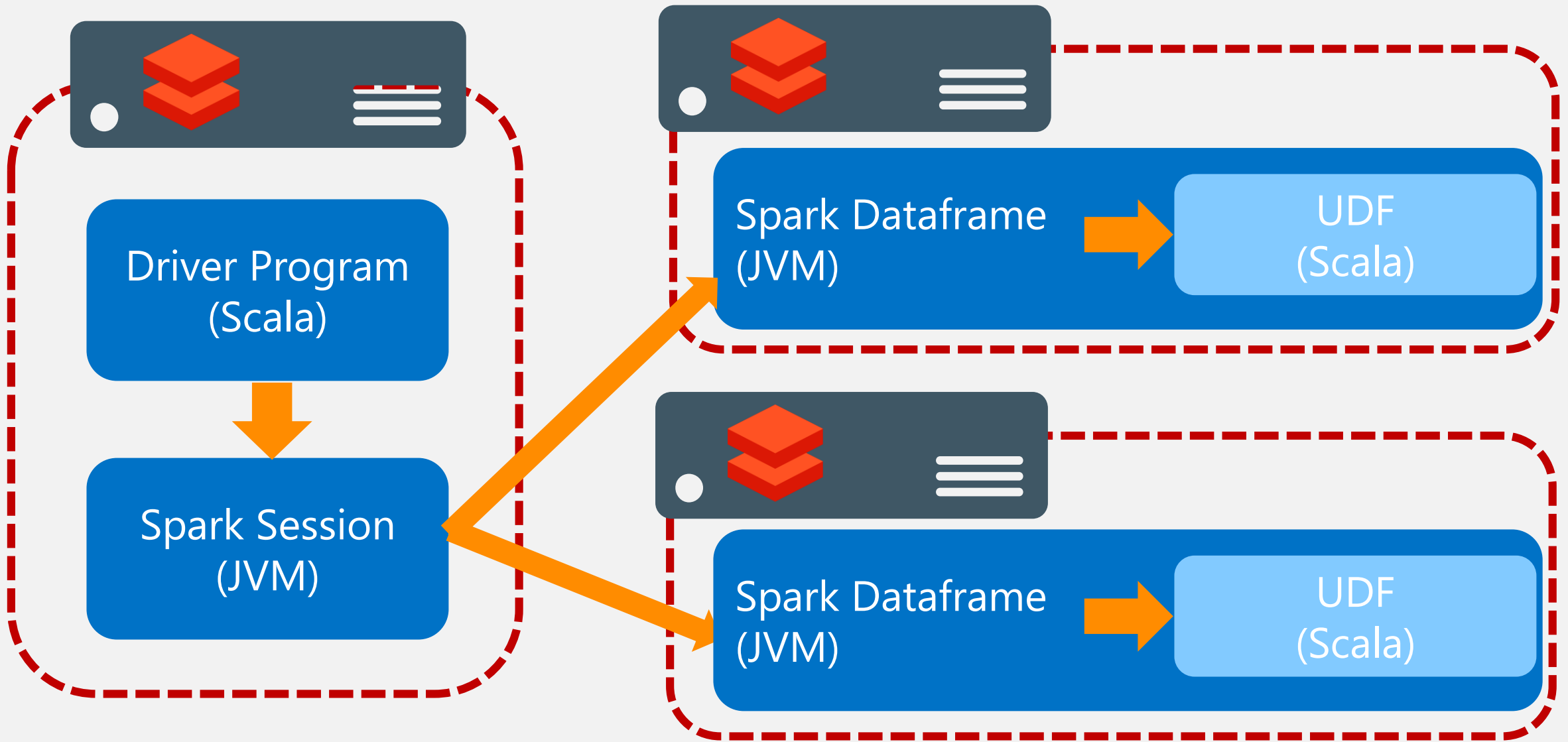
Serverless pools / High Concurrency

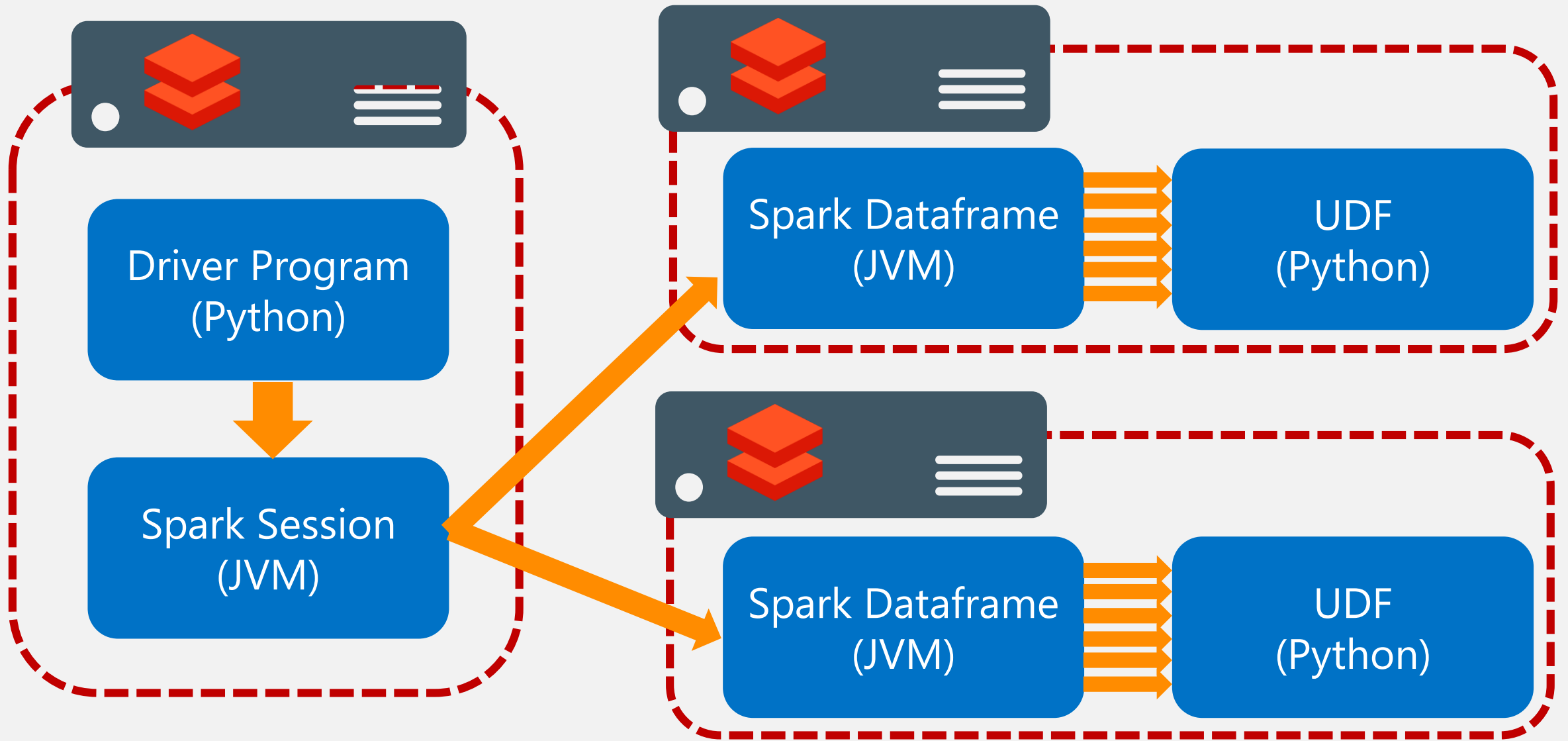
- Use SQL, Python, or R
- Want Azure Databricks to manage worker selection

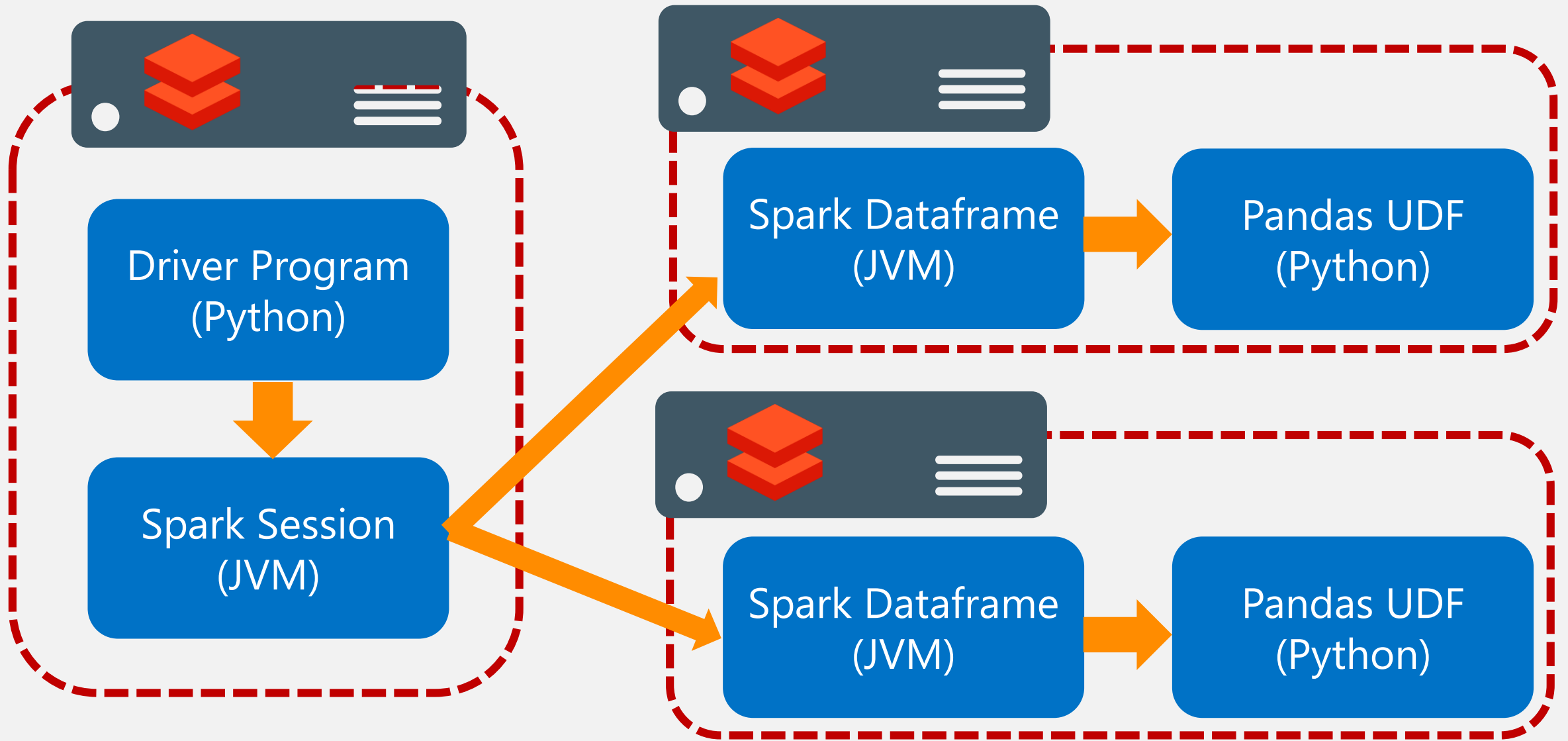
Standard cluster

- Use Scala
- Require a specific Spark version or want to configure Spark
- Want to control some advanced parameters

UDFs

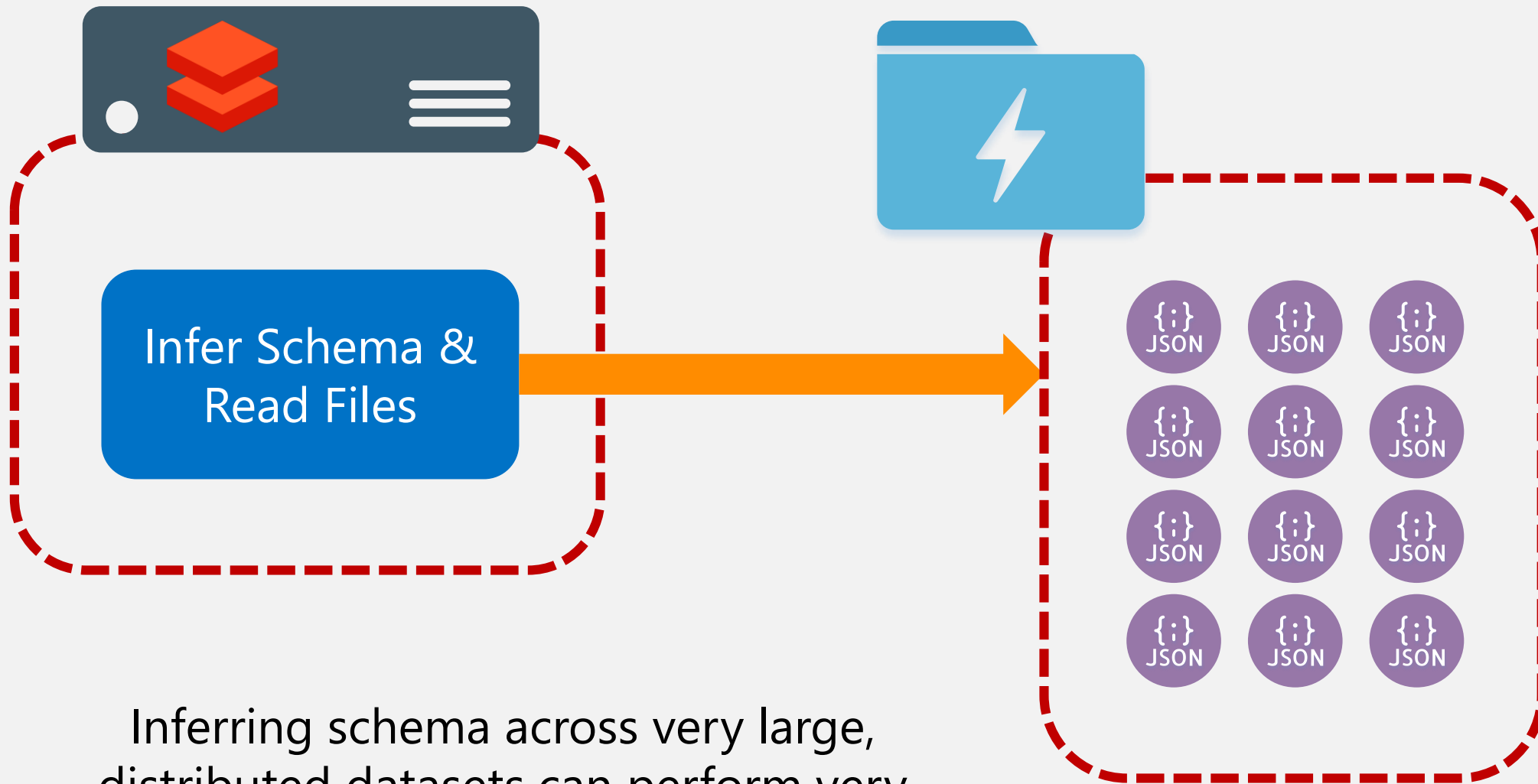




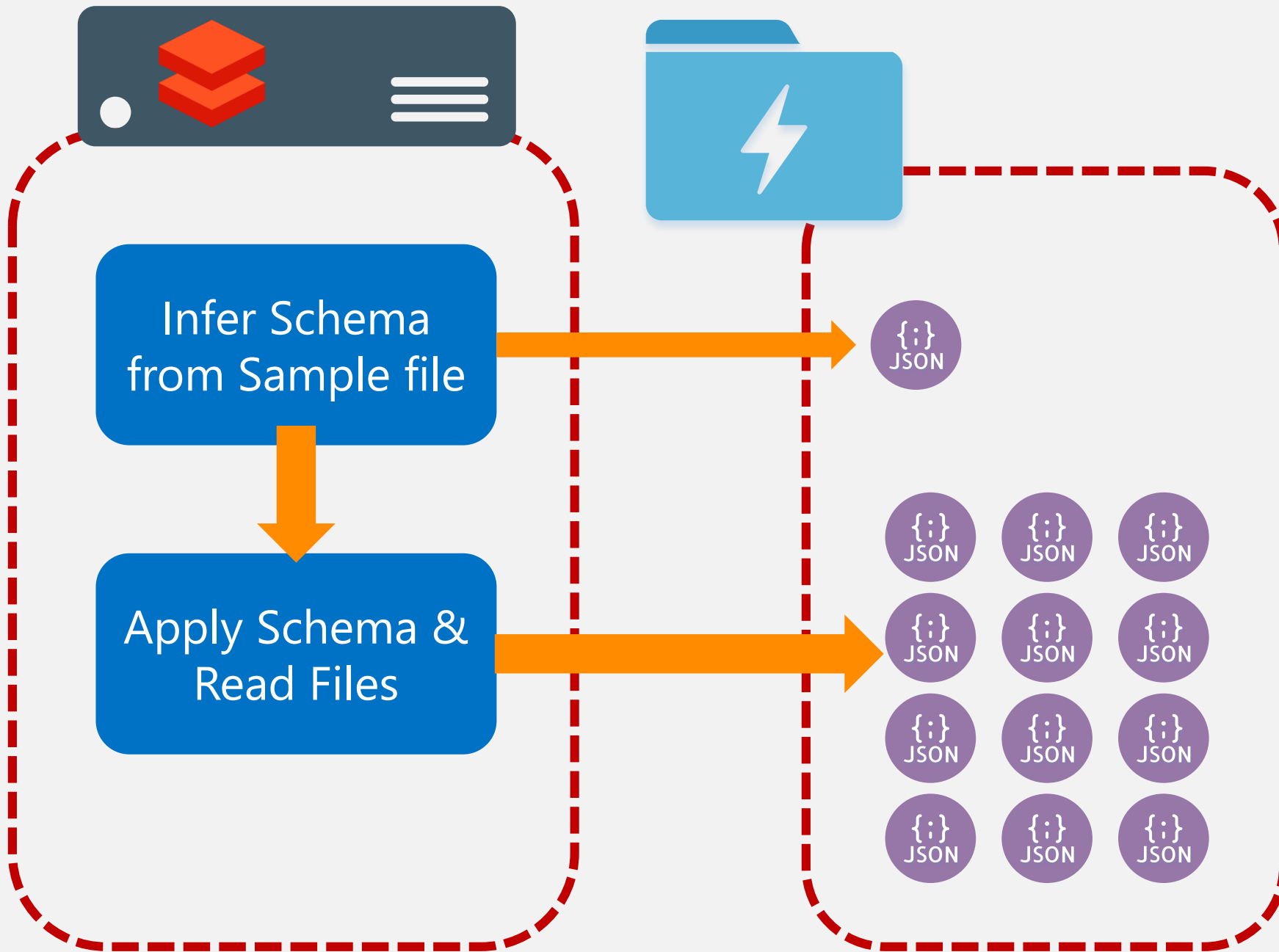


Performance comparison of different UDF methods in Databricks

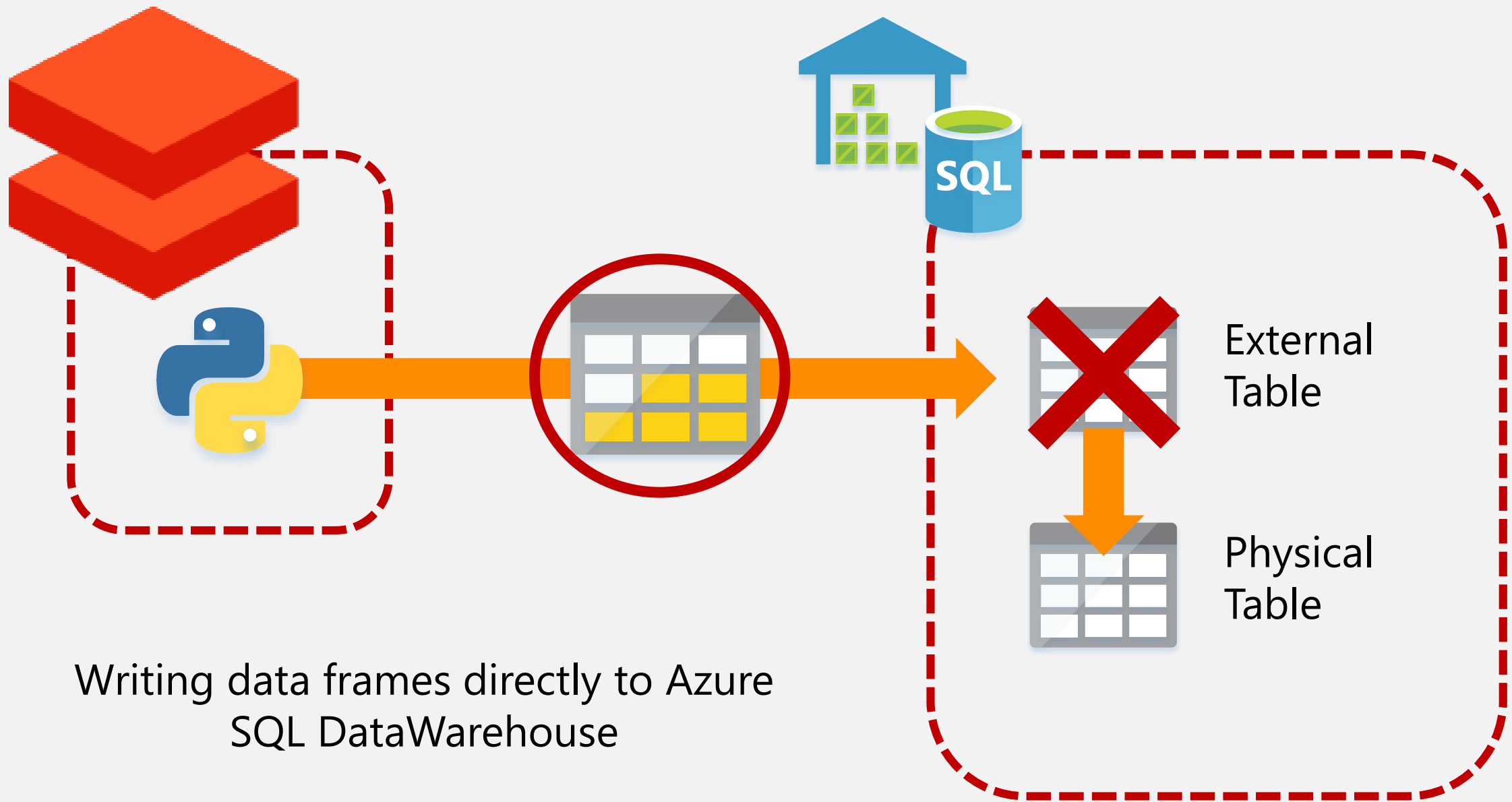
<https://bit.ly/2CAXkVI>



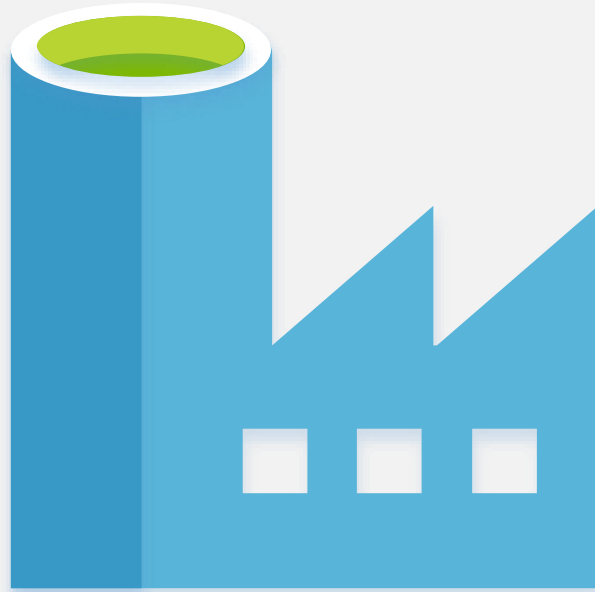
Inferring schema across very large,
distributed datasets can perform very
badly!



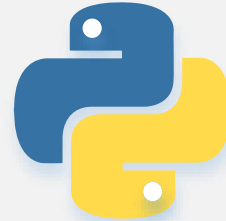
Acquire schema metadata (or infer from sample file) before reading large datasets!



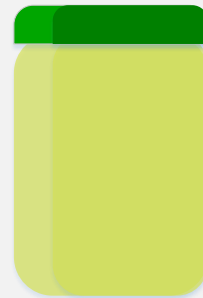
Orchestration



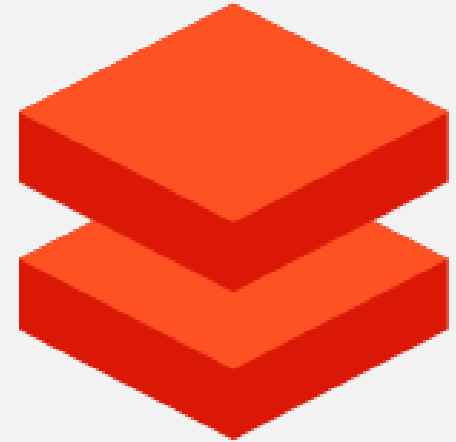
Notebook



Python
Script



Jar File

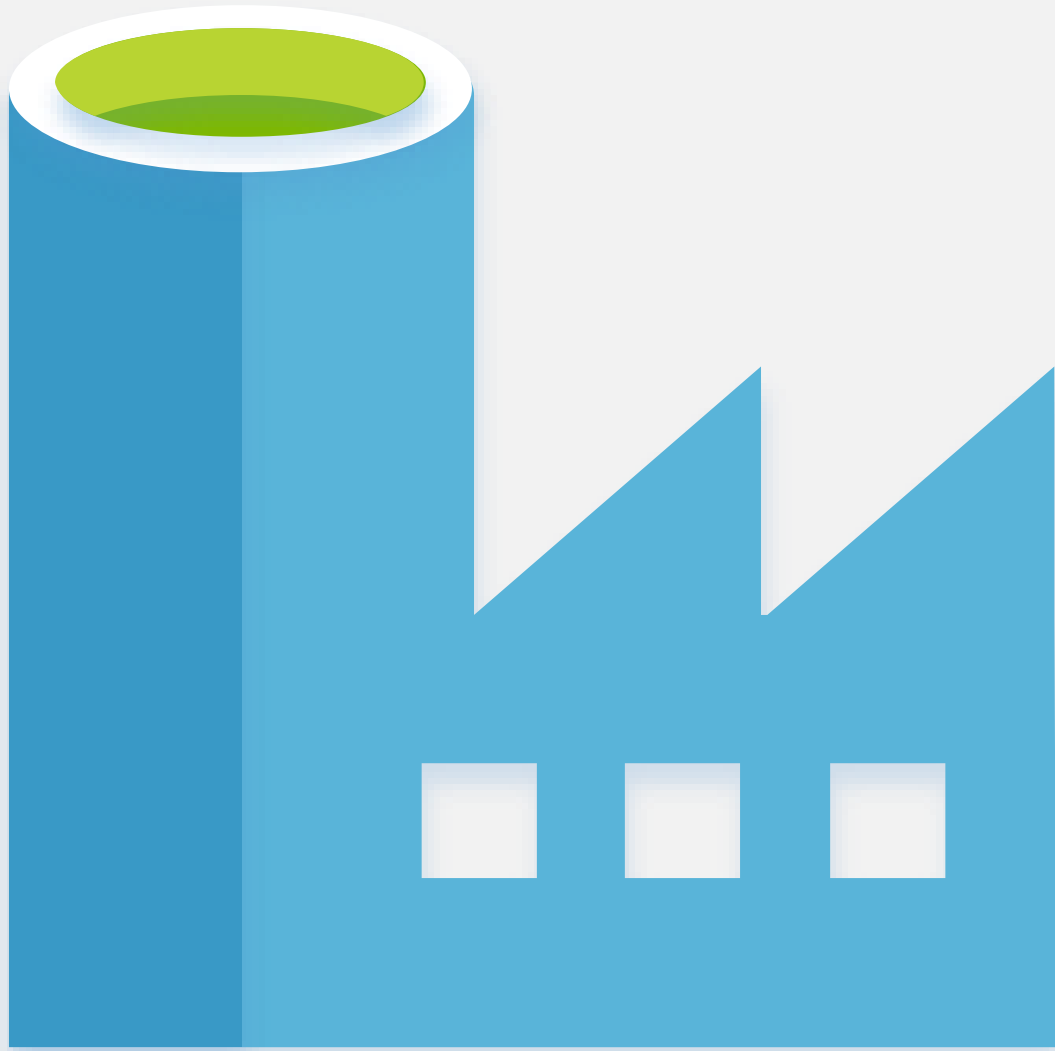


Data Factory

Demo

But what if I don't
want to write any
code?

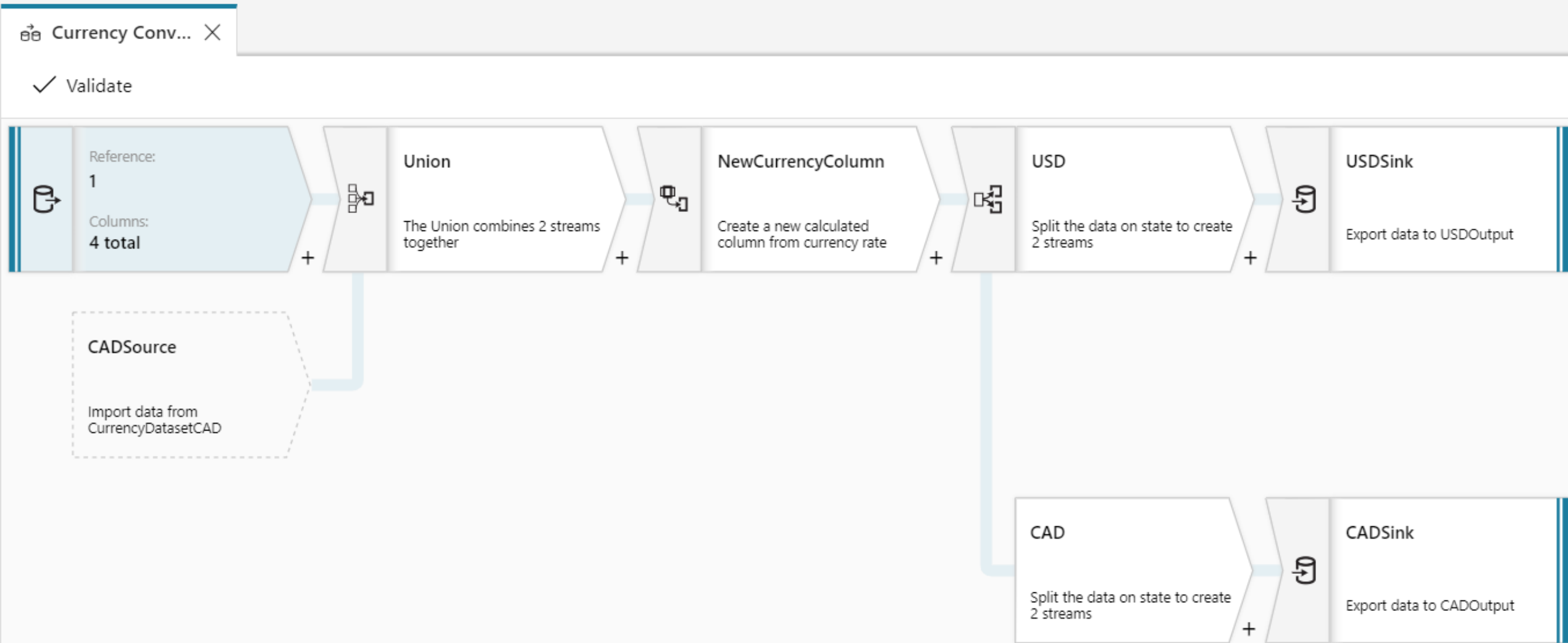


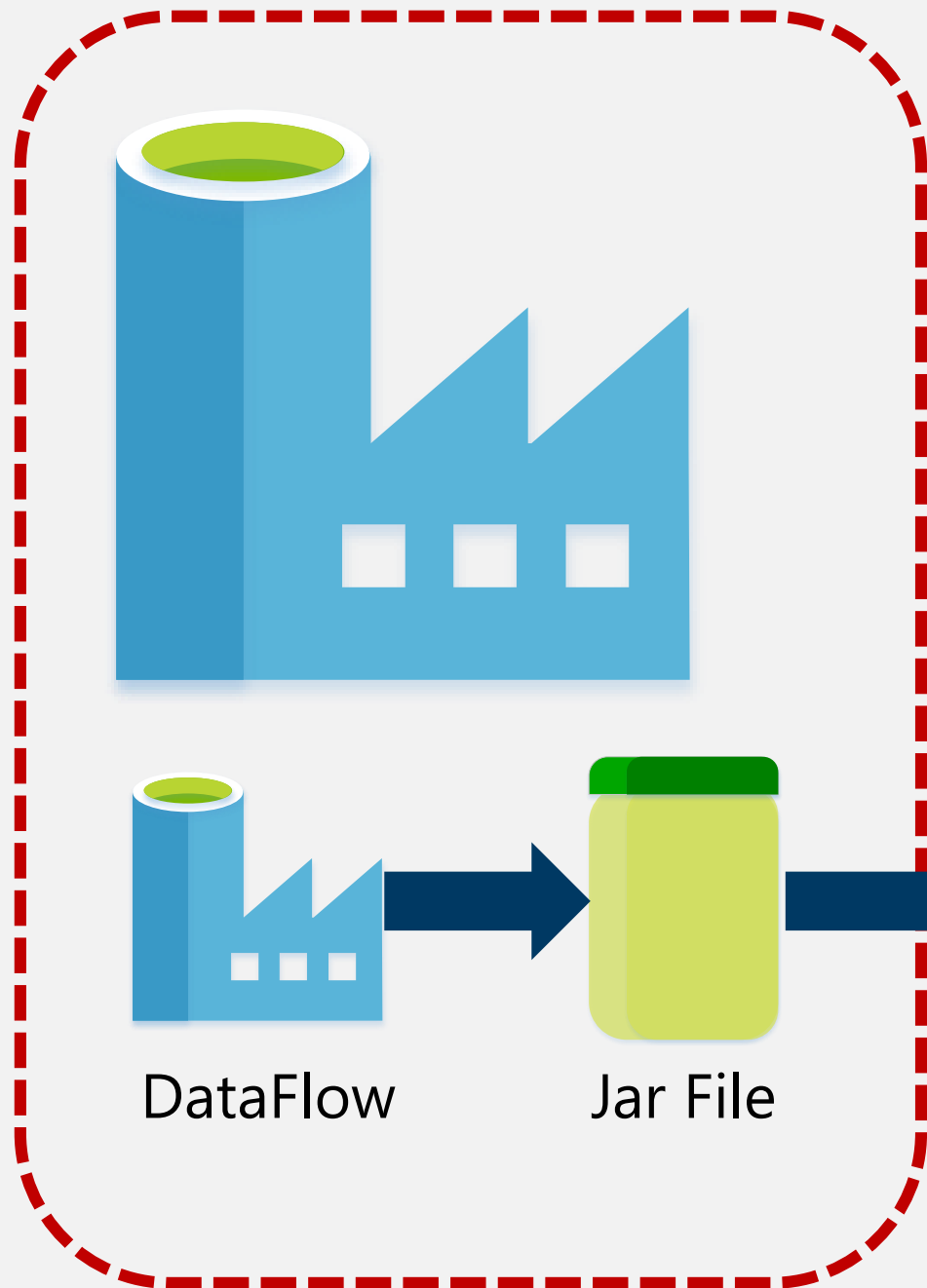


Azure Data Factory

Data Flows

New Data Factory DataFlows can write Databricks processing packages for you!!





Dataflows will compile down to a JAR file which will be sent to the Databricks cluster for execution

This means it uses Scala!

Thanks for Listening

Simon Whiteley

 @MrSiWhiteley



<http://blogs.adatis.co.uk>