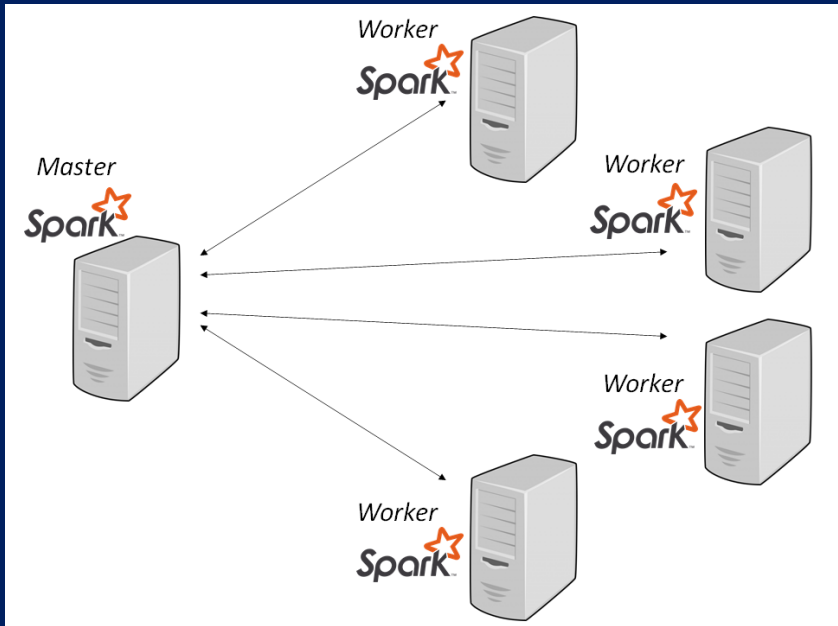


# Databricks

## Clusters



*Bryan Cafferky*  
*Big Data and AI Consultant*

<https://github.com/bcafferky/shared>

# Where Are We Heading?

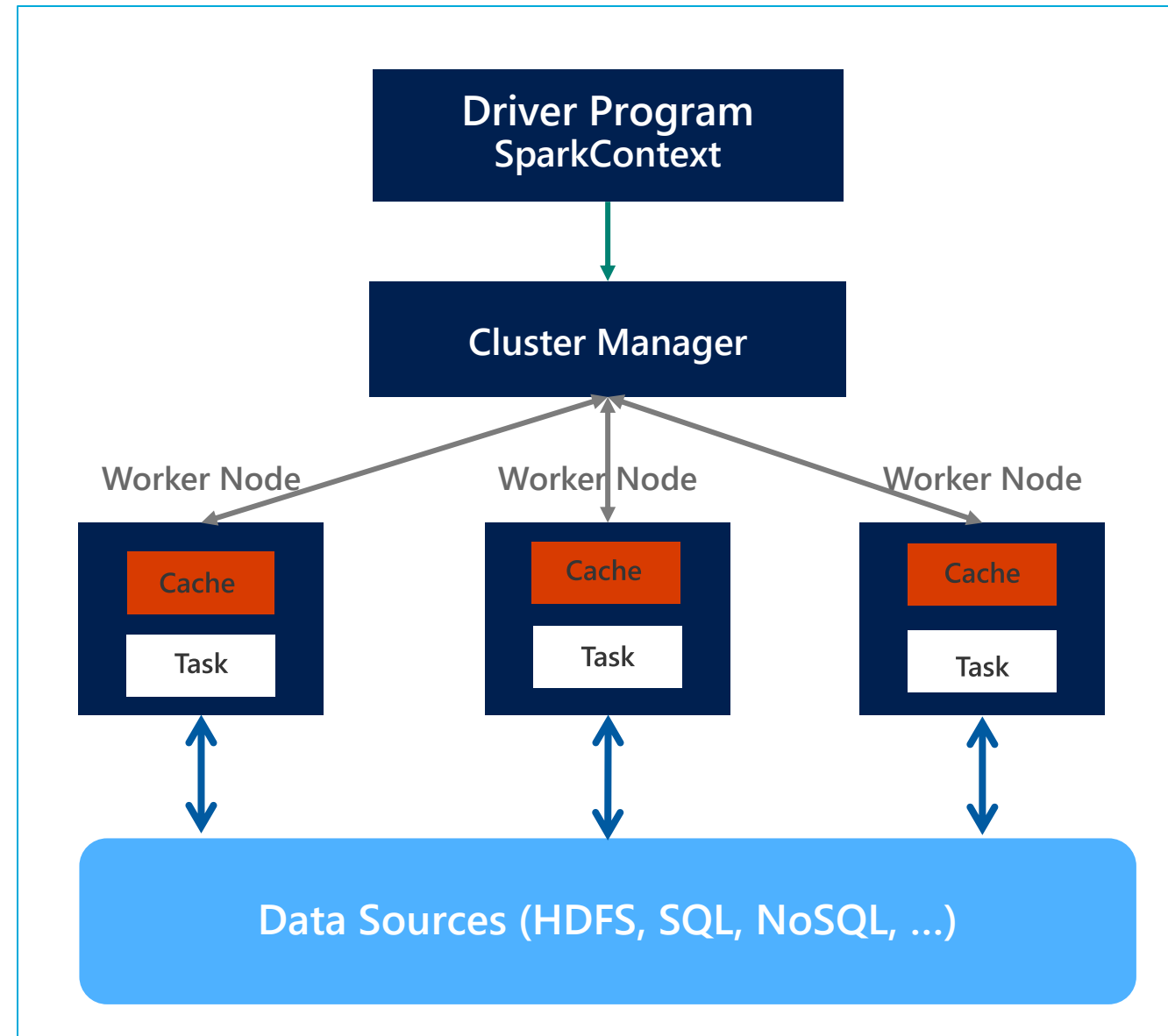
- What is a Databricks Cluster?
- Creating a Databricks Cluster
- Cluster Configuration Choices
- AdventureWorks Data Model
- Data Formats and Sources

# What is a Databricks Cluster?

- An Apache Spark cluster with the Databricks driver and enhanced features.
- A set of computers to do Data Engineering workloads.
- All work is done on a cluster.
- The work is coordinated by a node called the Driver.
- The data is processed by Worker Nodes.
- A Single Node can have multiple parallel processes called Executors.

# GENERAL SPARK CLUSTER ARCHITECTURE

- 'Driver' runs the user's 'main' function and executes the various parallel operations on the worker nodes.
- The results of the operations are collected by the driver
- The worker nodes read and write data from/to Data Sources including HDFS.
- Worker node also cache transformed data in memory as RDDs (Resilient Data Sets).
- Worker nodes and the Driver Node execute as VMs in public clouds (AWS, Google and Azure).



# Databricks Workspace

Microsoft Azure

PORTAL brcaffer@microsoft.com

Demo\_CognitiveServices\_ImageTextExtraction (1) (Scala)

Workspace/Demo/Demo\_CognitiveServices\_ImageTextEx (1)

Blob Account Key : [redacted] Blob Account Name : bcafferkystoragexxx

Blob Container Name : bcafferkyc231

Cmd 1

### Set Parameters to Notebook Variables

```
val blobcontainer = dbutils.widgets.get("containername")
val blobaccount = dbutils.widgets.get("accountname")
val blobaccountkey = dbutils.widgets.get("accountkey")
val bloburl = "wasbs://" + blobcontainer + "@" + blobaccount + ".blob.core.windows.net"
val blobconf = "fs.azure.account.key." + blobaccount + ".blob.core.windows.net"
```

Cmd 2

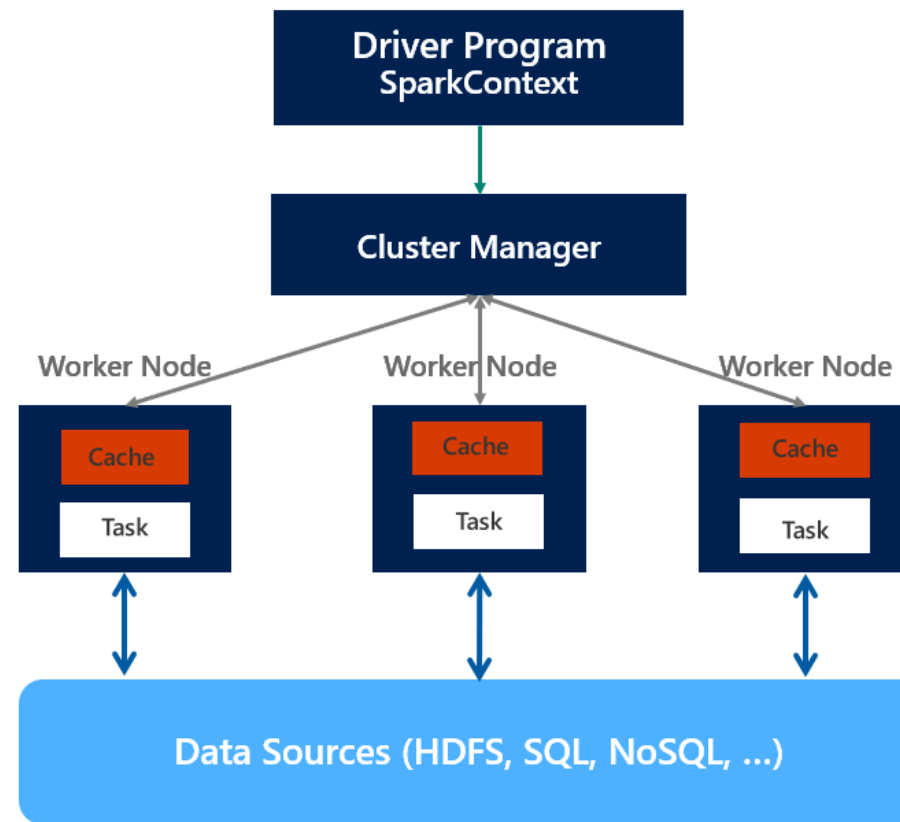
### Mount the Blob Container to DBFS to /mnt/blob

```
dbutils.fs.mount(
  source = bloburl,
  mountPoint = "/mnt/blobimages",
  extraConfigs = Map(blobconf -> blobaccountkey))
```

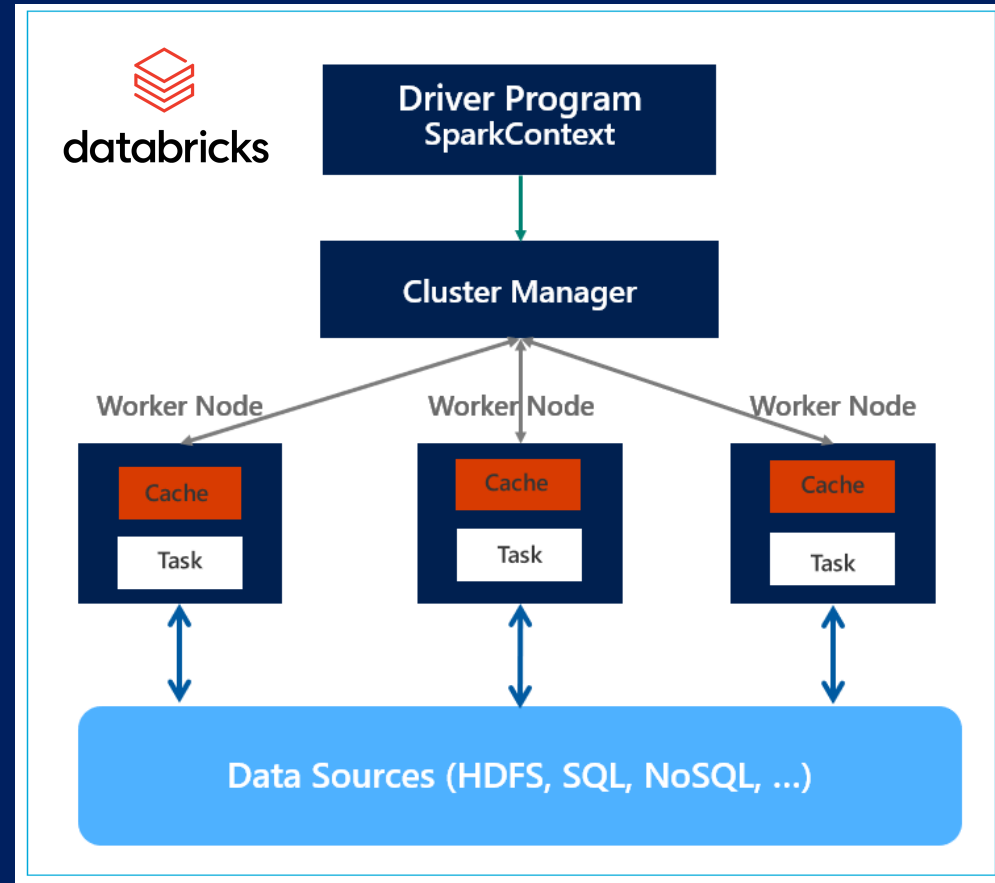
Cmd 3

### View files in the DBFS mount point (which is mounted to Azure Storage)

```
%fs ls /mnt/blobimages
```



# Creating a Databricks Cluster



# Creating a Databricks Cluster – Main Screen

## Create Cluster

### New Cluster

Cancel

Create Cluster

DBU / hour: 2.25 - 6.75 ?

2-8 Workers: 28-112 GB Memory, 8  
1 Driver: 14 GB Memory, 4 Cores

Cluster name

mycluster

Cluster mode ?

Standard

Databricks runtime version ?

[Learn more](#)

Runtime: 9.1 LTS (Scala 2.12, Spark 3.1.2)

**50% promotional discount applied to Photon during preview** ?

X

Autopilot options

☒ Enable autoscaling ?

☒ Terminate after  minutes of inactivity ?

Worker type ?

Standard\_DS3\_v2

14 GB Memory, 4 Cores

Min workers

2

Max workers

8

⚠

☐ Spot instances ?

**New** Configure separate pools for workers and drivers for flexibility. [Learn more](#)

Driver type

Same as worker

14 GB Memory, 4 Cores

DBU / hour: 2.25 - 6.75 ?

Standard\_DS3\_v2

▶ Advanced options

Databricks Cluster Type

Databricks Runtime

Let Databricks Scale Cluster as Needed

Pause the Cluster After X Minutes of No Use


Minimum and Maximum # of Worker Nodes


Keeps VMs Ready to Use for Cluster

Driver Node VM Type

Customer Cluster Configuration

# Creating a Databricks Cluster – Cluster Mode

Cluster mode 

Standard 

- High Concurrency
- Standard
- Single Node

Option	Meaning
High Concurrency	Optimized for sharing the cluster with many concurrent users. Scala is not supported in this cluster mode.
Standard	Standard Spark Cluster mode of a driver and workers
Single Node	Just a Driver Node is created to minimize costs. Good for learning.



# Creating a Databricks Cluster – Autopilot Options

Autopilot options

☒ Enable autoscaling 

☒ Terminate after  minutes of inactivity 

Option	Meaning
Enable autoscaling	Let Databricks automatically increase or decrease the number of worker nodes based on the workload demanded at any given time.
Terminate after	To save money, pause the cluster after the specificized number of minutes have passed with no usage.

# Creating a Databricks Cluster – Worker Type

Worker type ?

Standard\_DS3\_v2 14 GB Memory, 4 Cores ▼

**General Purpose**

- Standard\_DS3\_v2 14 GB Memory, 4 Cores
- Standard\_DS4\_v2 28 GB Memory, 8 Cores
- Standard\_DS5\_v2 56 GB Memory, 16 Cores
- 30 more

**General Purpose (HDD)**

- Standard\_D3\_v2 14 GB Memory, 4 Cores
- Standard\_D4\_v2 28 GB Memory, 8 Cores
- Standard\_D5\_v2 56 GB Memory, 16 Cores
- 4 more



**Memory Optimized (Remote HDD)**

- Standard\_D12\_v2 28 GB Memory, 4 Cores

Option	Meaning
Worker type	Select the Virtual Machine Type to use for the worker nodes.

# Creating a Databricks Cluster – Cluster Mode

Min workers    Max workers

        ☐ Spot instances 

Option	Meaning
Min workers	The smallest number of worker nodes that must be maintained for the cluster.
Max workers	The largest number of worker nodes that may be created for the cluster.
Spot instances	<p>Spot instances provide a way to save money by letting Databricks take unused VMs available in the cloud platform as discounted costs.</p> <p><a href="https://techcommunity.microsoft.com/t5/analytics-on-azure-blog/azure-databricks-and-azure-spot-vms-save-cost-by-leveraging/ba-p/2374187">https://techcommunity.microsoft.com/t5/analytics-on-azure-blog/azure-databricks-and-azure-spot-vms-save-cost-by-leveraging/ba-p/2374187</a></p>

# Creating a Databricks Cluster – Driver Type

Driver type

Same as worker 14 GB Memory, 4 Cores ▼

**General Purpose**

Same as worker 14 GB Memory, 4 Cores

Standard\_DS3\_v2 14 GB Memory, 4 Cores

Standard\_DS4\_v2 28 GB Memory, 8 Cores

31 more

**General Purpose (HDD)**

Standard\_D3\_v2 14 GB Memory, 4 Cores

Standard\_D4\_v2 28 GB Memory, 8 Cores

Standard\_D5\_v2 56 GB Memory, 16 Cores

4 more

**Memory Optimized (Remote HDD)**

Standard\_D12\_v2 28 GB Memory, 4 Cores

Option	Meaning
Driver type	Select the Virtual Machine Type to use for the cluster driver. Sometimes you want this to be different like wanted to do more work on the driver.

# Creating a Databricks Cluster – Advanced options

## ▼ Advanced options

### Azure Data Lake Storage credential passthrough ?

☐ Enable credential passthrough for user-level data access

Spark

Tags

Logging

Init Scripts

### Spark config ?

Enter your Spark configuration options here. Provide only one key-value pair per line.

Example:

spark.speculation true

spark.kryo.registrator my.package.MyRegistrator

### Environment variables ?

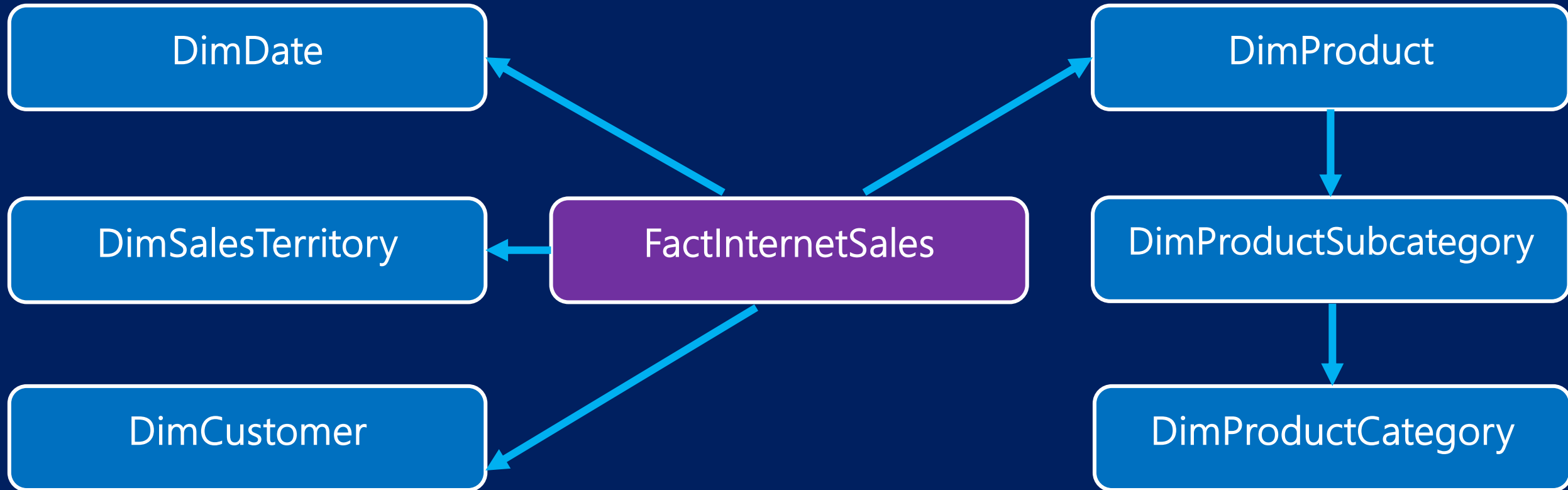
PYSPARK\_PYTHON=/databricks/python3/bin/python3

Option	Meaning
Advanced options	Click on Advanced options to open the Advanced options form. Here you can customize the cluster configuration and perform additional set up work like installing open source libraries.

# About the data



# AdventureWorks Data Model

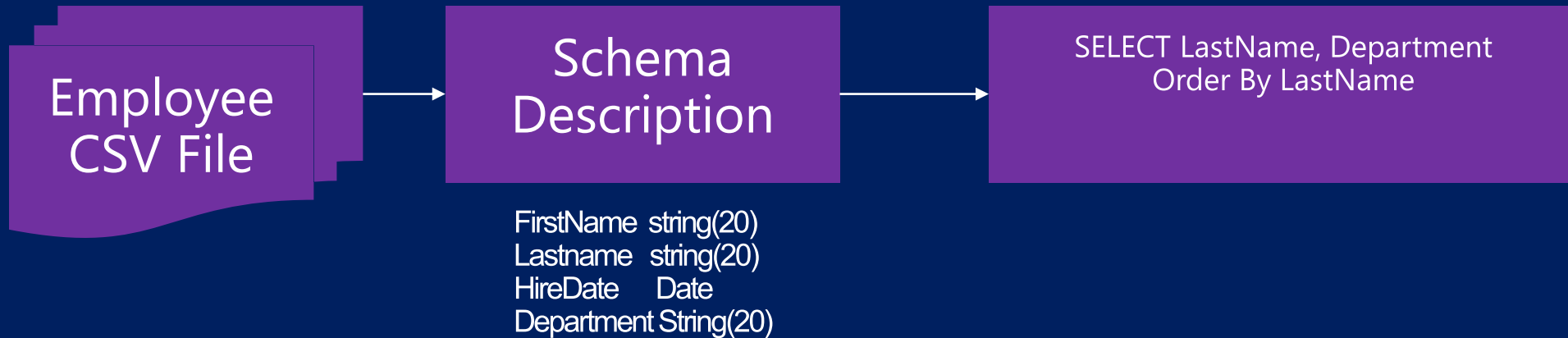


➤ Each table has a common key to join them.



# Schema On Read

- Data Is Not in an RDMS
- External File Is Described Structurally



Built-In Databricks Datasets:

<https://docs.databricks.com/data/databricks-datasets.html>





# Databricks Supported Data Types and Sources

## Text Files

- CSV
- JSON
- XML
- TXT

## Public Cloud

- Databases
- NoSQL
- Cloud Storage
- Data Services

## Big Data Formats

- Parquet
- Delta Lake
- ORC
- Avro

# Where We've Been

- What is a Databricks Cluster?
- Creating a Databricks Cluster
- Cluster Configuration Choices
- AdventureWorks Data Model
- Data Formats and Sources