

Big Data Clusters for the Absolute Beginner

Abstract: Are you a DBA or data professional working with outdated technology or feel your skills are quickly becoming irrelevant? Stay ahead in today's ever changing tech space and take your career to the next level.

Big Data Clusters for the Absolute Beginner

Mohammad Darab, MaeTek, LLC

Email: mo@maetek-llc.com

About Me

- From Washington, DC
- 19+ years in IT, 10+ years working with SQL Server
- Idera ACE 2019
- Online:

MohammadDarab.com

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



Agenda

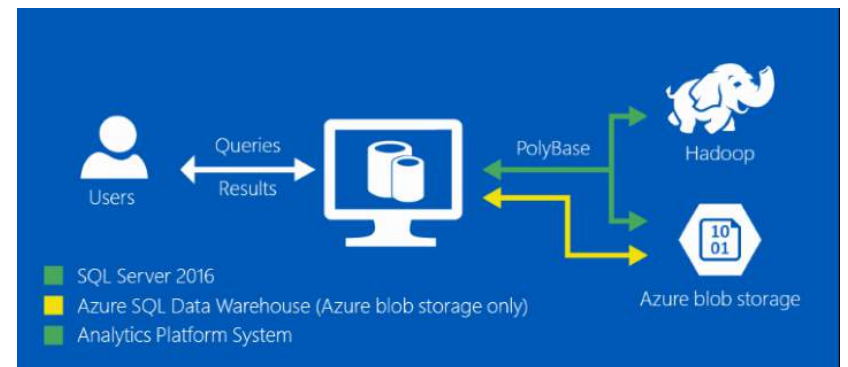
- History
- What is a BDC
- Architecture
- Features
- Mo's Thoughts
- Learning Path

<https://MohammadDarab.com/bdc>

History

2016

- Introduction of **Polybase**
 - Combines non-relational and relational data
 - Connectors for Hortonworks, Cloudera, Azure blob Storage
 - Read data from Hadoop

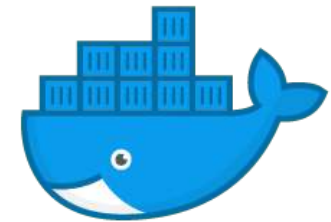
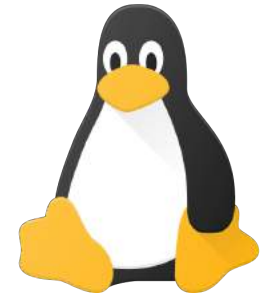


Source: Microsoft

<https://MohammadDarab.com/bdc>

2017

- Added support to run on **Linux**
- Run SQL Server on **Containers**

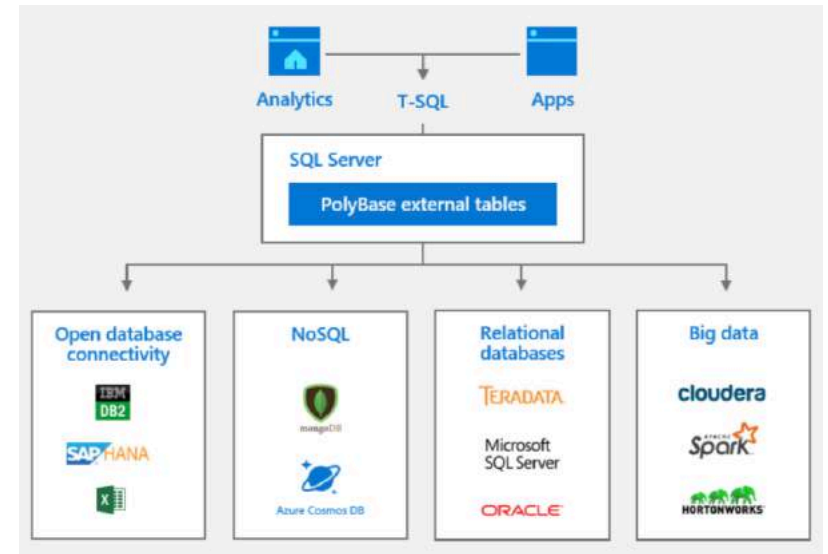


docker

<https://MohammadDarab.com/bdc>

2019

- Enhanced Polybase
- Availability Groups on Kubernetes



Source: Microsoft

<https://MohammadDarab.com/bdc>

Polybase vs Linked Servers

Polybase

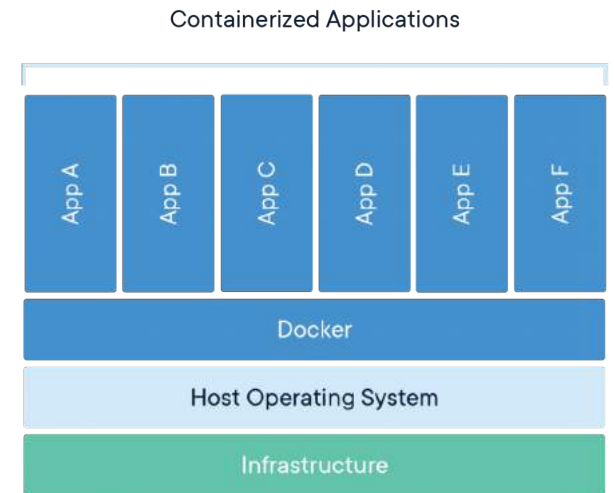
- Database scoped
- Supports read-only operations for all data sources and insert operations for HADOOP and data pool sources
- No separate config needed for AG
- Suitable for analytic queries processing large numbers of rows

Linked Servers

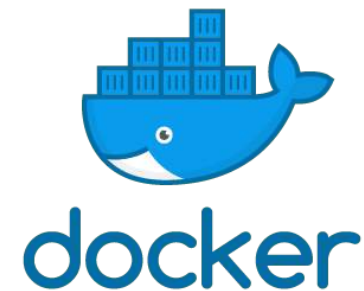
- Instance scoped
- Supports both read and write operations
- Separate config needed for each instance in an AG
- Suitable for OLTP queries returning single or few rows

Containers

A **container** is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another



Source: Docker



Kubernetes

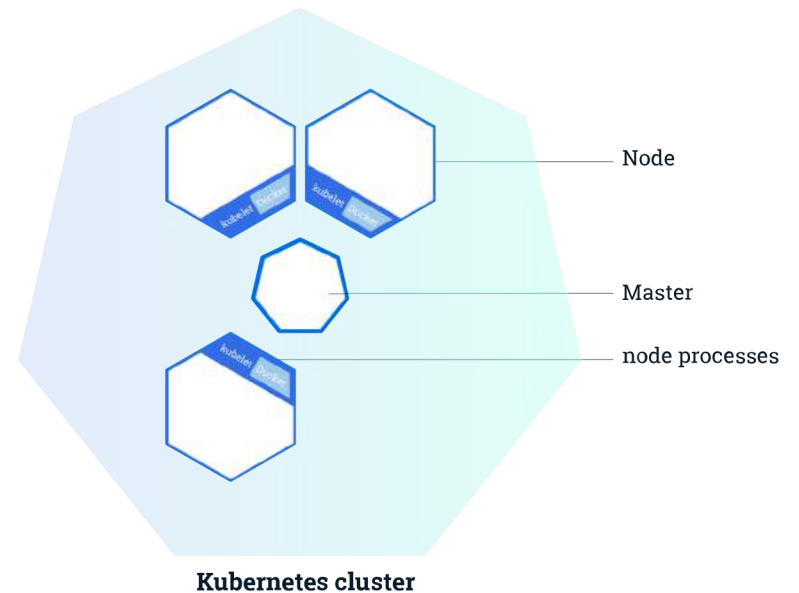
- Is a container "orchestrator"
- Can scale container deployment according to need



<https://MohammadDarab.com/bdc>

Kubernetes Cluster

- Set of machines called nodes
- Master node controls cluster
- The rest are worker nodes
- Master is in charge of work distribution between workers and monitoring cluster health



Source: Kubernetes.io

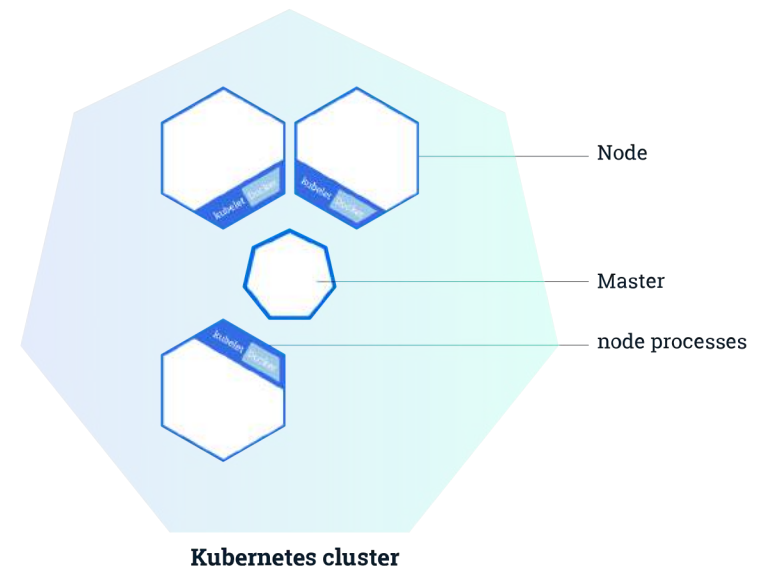
<https://MohammadDarab.com/bdc>

Kubernetes

Node: Runs containerized applications. It can be a VM or physical machine.

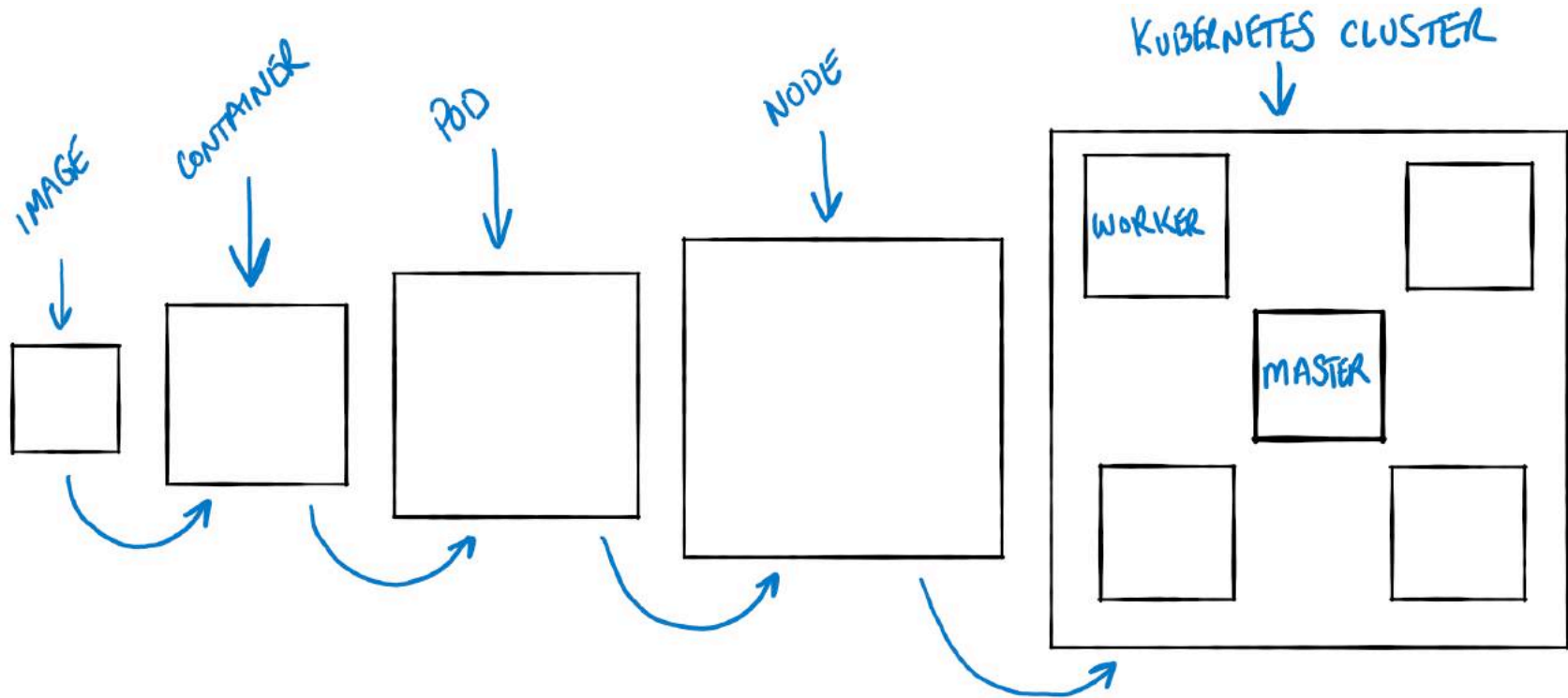
Pod: Logical group of one or more containers and associated resources.

A pod runs on a node, and a node can run one or more pods.



Source: Kubernetes.io

The Flow



Why?

Data Never Sleeps

90% of all data today was generated in the last two years.

<https://MohammadDarab.com/bdc>

Walmart

- Collects **2.5 Petabytes** of data from 1 mil customers every hour
- Billions of Facebook messages, tweets, YouTube videos, etc.

<https://MohammadDarab.com/bdc>

\$\$\$

“We want to know what every product in the world is. We want to know who every person in the world is. And we want to have the ability to connect them together in a transaction.”

-Walmart's CEO 2013

<https://MohammadDarab.com/bdc>

A close-up photograph of a person's hand raised in the air, palm facing forward. The hand is positioned on the right side of the frame. The person is wearing a watch with a dark, textured strap on their left wrist. The background is dark and out of focus, featuring several warm, glowing bokeh lights, suggesting an indoor party or social gathering. The overall mood is celebratory and energetic.

Pick Me!

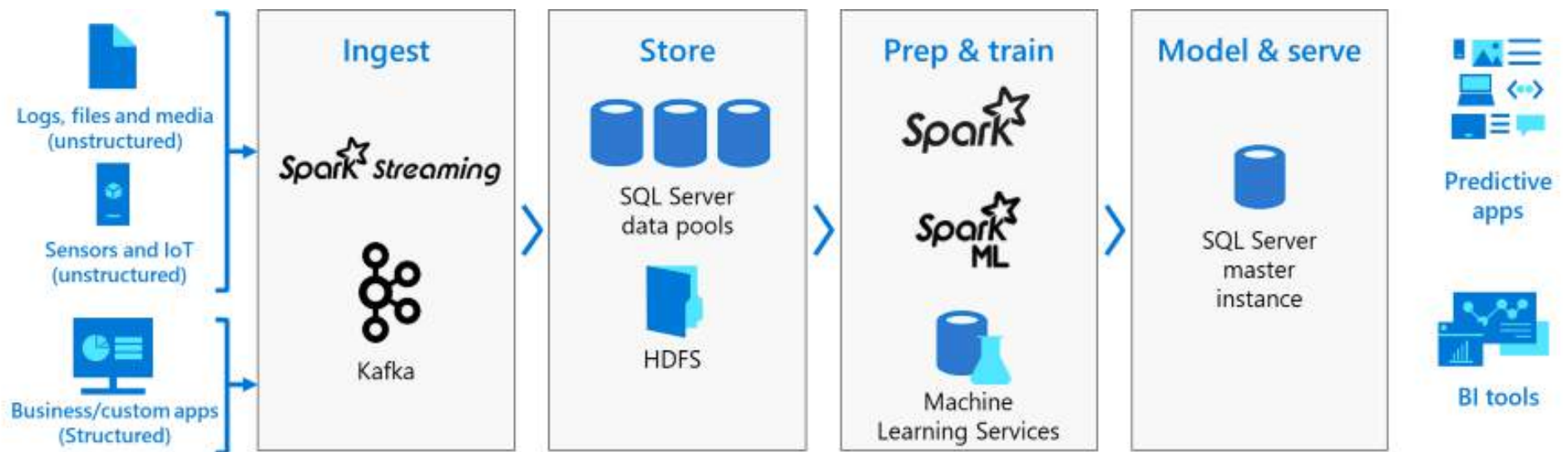
What is a BDC?

MS Definition

“BDCs allow you to deploy scalable clusters of SQL Server, Spark, & HDFS containers running on Kubernetes. These components are running side by side to enable you to read, write, & process big data from T-SQL or Spark, allowing you to easily combine & analyze your high-value relational data with high-volume big data.” ([link](#))

<https://MohammadDarab.com/bdc>

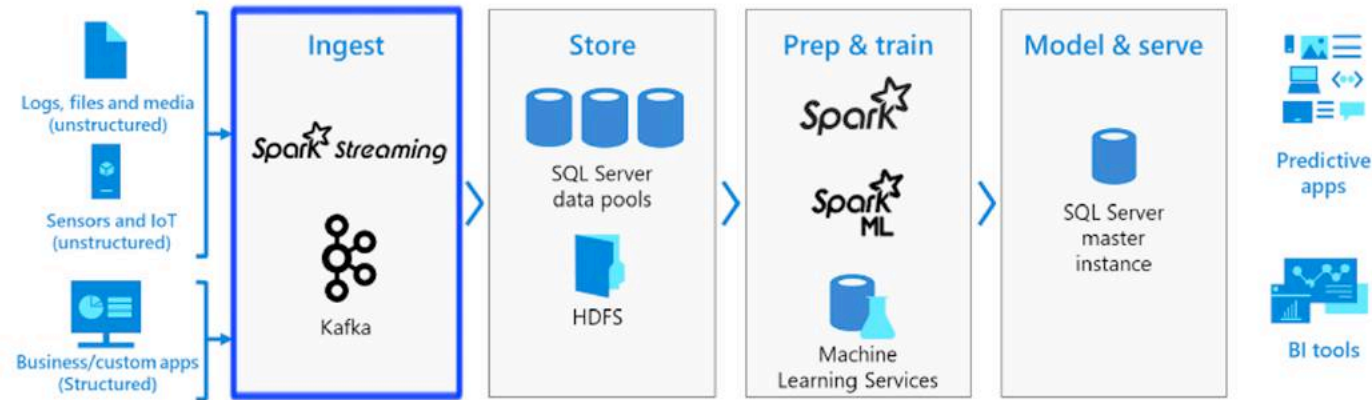
AI & ML Platform



Source: Microsoft

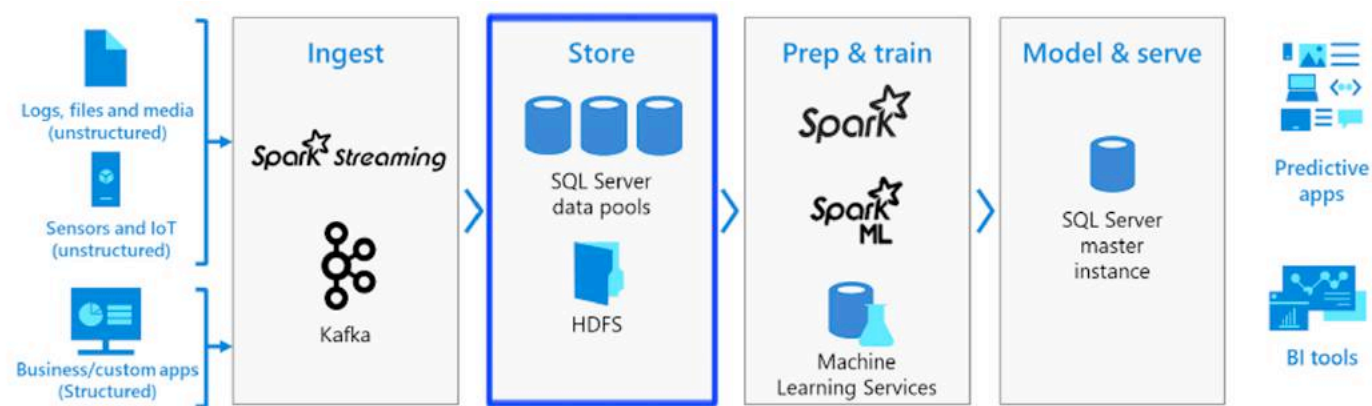
<https://MohammadDarab.com/bdc>

Ingest Data



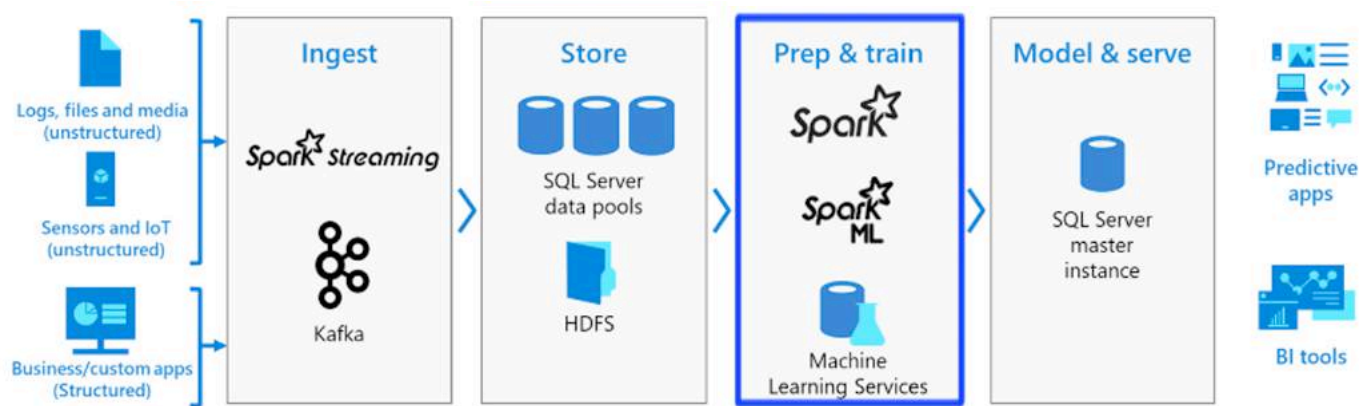
- Data can be ingested by either Spark Streaming, by inserting directly into HDFS, or,
- Directly into SQL Server by T-SQL queries

Store Data



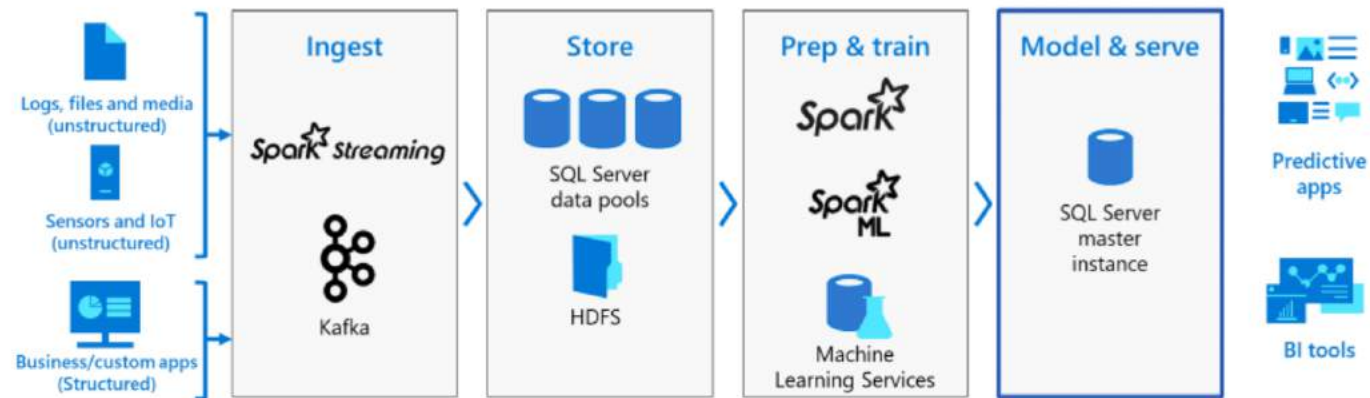
- Data can be stored in files in HDFS, or partitioned and stored in data pools, or,
- Stored in the SQL Server master instance, in tables, graphs or JSON/XML

Prep & Train



- Data can be prepared by Spark jobs or T-SQL queries in files in HDFS, or partitioned and stored in data pools, or,
- Java, Python, R, and Scala to feed into machine learning model training routines

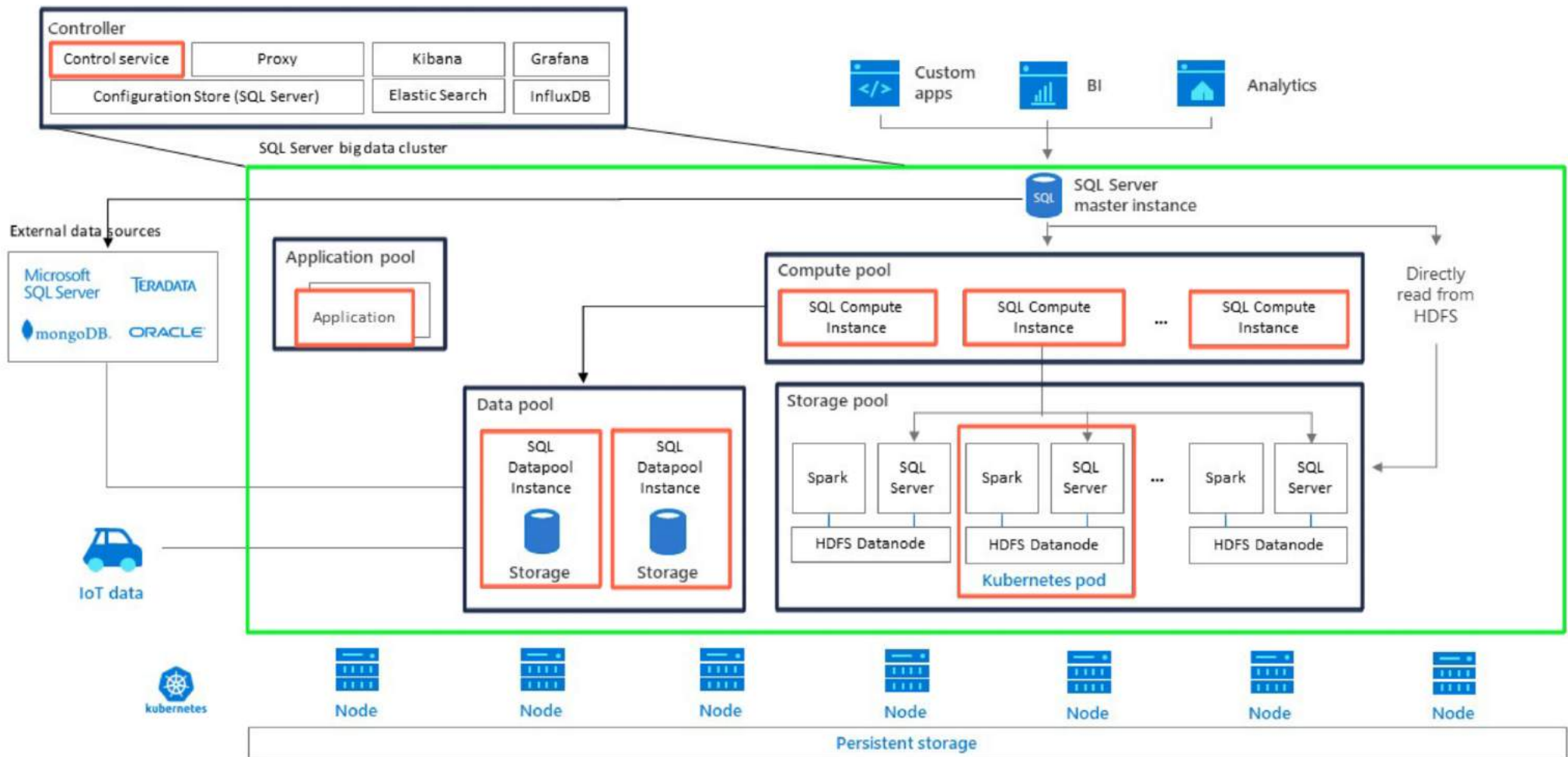
Model & serve



- PREDICT function in stored procedure in the SQL master instance, or,
- Batch scoring over the data in HDFS using Spark

Architecture

Big Data Cluster Architecture



Source: Microsoft

Controller

Controller hosts the core logic for deploying and managing a big data cluster. It takes care of all interactions with Kubernetes, SQL Server instances that are part of the cluster and other components like HDFS and Spark.



Controller (cont.)

Controller provides the following functionality:

- Manages cluster lifecycle, cluster bootstrap, delete, update, etc.
- Manages master SQL Server instance
- Manages compute, data, and storage pools
- Manages cluster security

Master Instance

Connectivity - Provides an externally accessible endpoint for the cluster (you can connect with ADS or SSMS).

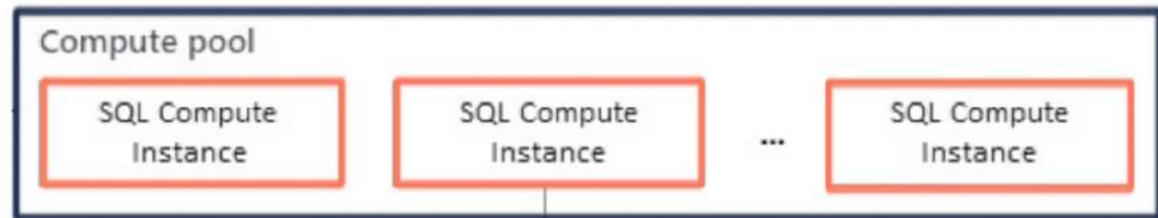
Scale-out query management - Contains the scale-out query engine that is used to distribute queries across SQL Server instances on nodes in the compute pool.

Master Instance (cont.)

ML - machine learning services is an add-on feature to the database engine. Once external script execution is enabled on the master instance, you can execute Java, R and Python scripts using `sp_execute_external_script`.

Compute Pool

Compute Pool provides computational resources to the cluster.

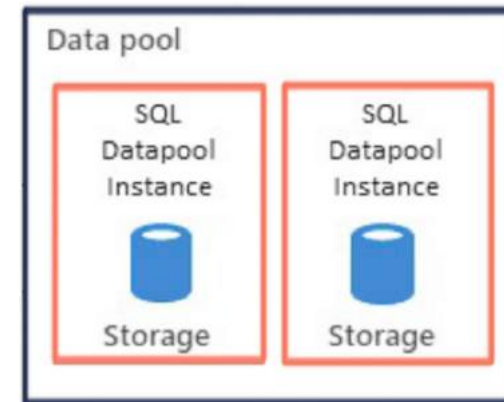


Compute Pool (cont.)

- The pods in the compute pool are divided into SQL Compute instances for distributed queries (specifically Polybase).
- Main role of the compute pool is to perform intermediate joins / aggregations of multiple external tables.

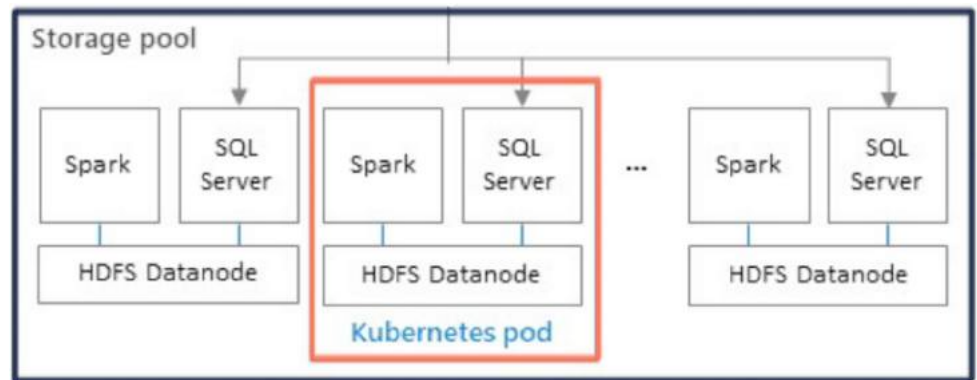
Data Pool

Data Pool is used for data persistence and caching. It is used to ingest data from SQL queries or Spark jobs. SQL Server big data cluster data marts are persisted in the data pool.



Storage Pool

Storage Pool consists of pods comprised of SQL Server on Linux, Spark, and HDFS. All the storage nodes in a SQL Server big data cluster are members of an HDFS cluster.



Storage Pool (cont.)

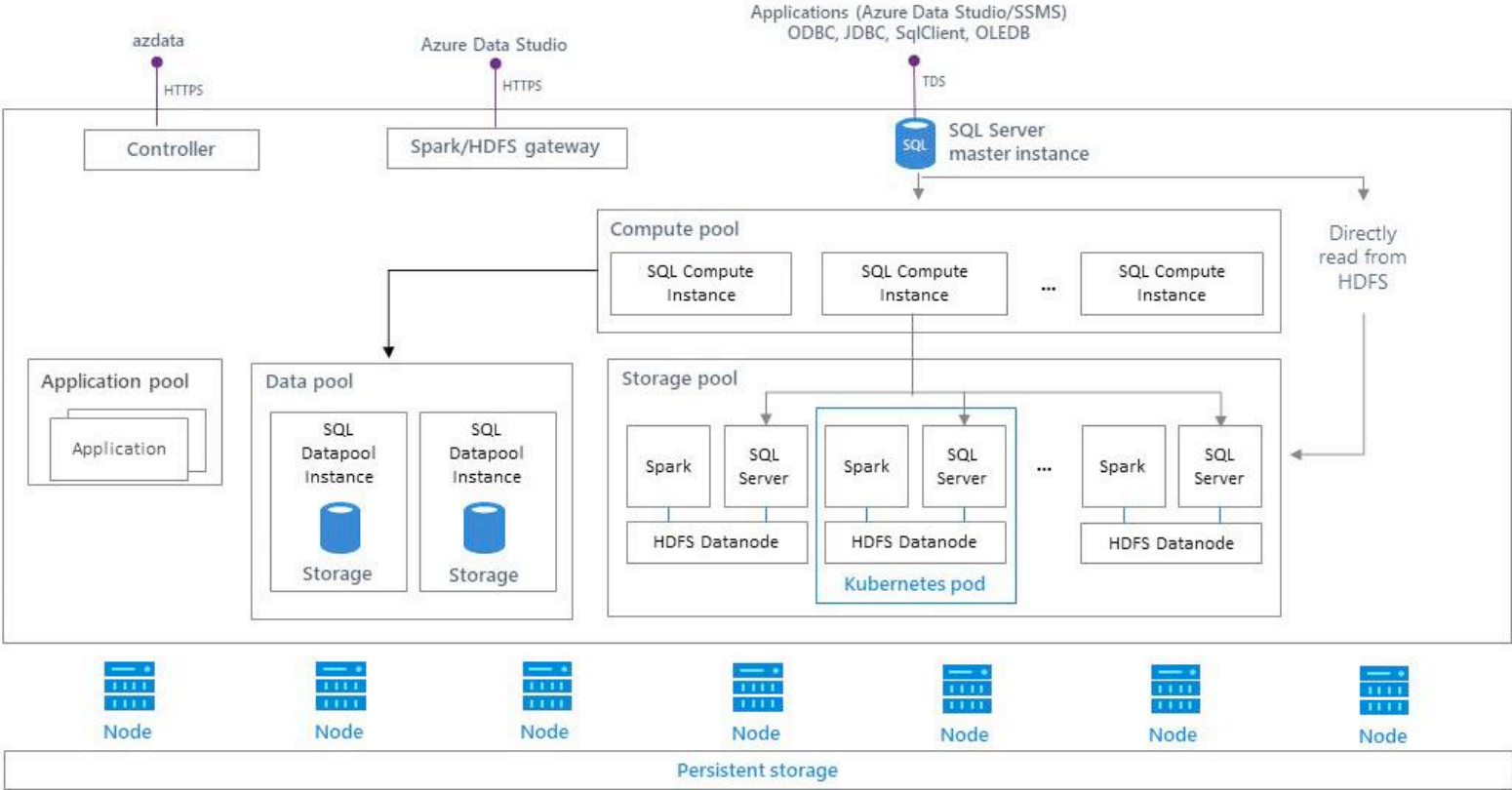
Storage nodes are responsible for:

- Data ingestion through Spark
- Data storage in HDFS (Parquet format). HDFS data is spread across all storage nodes in the BDC for persistency
- Data access through HDFS and SQL Server endpoints

Kubernetes Storage

- **Persistent Volumes (PV)** are how we map external storage onto the Kubernetes cluster
- **Persistent Volume Claims (PVC)** are like tickets that authorize applications (pods) to use a PV

Cluster Endpoints



Source: Microsoft

Points of Entry

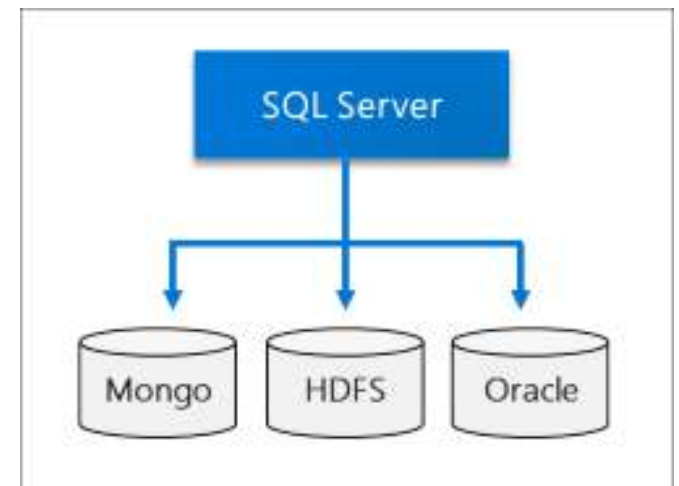
- **Controller endpoint** – Big data cluster management service that exposes REST APIs for managing the cluster
- **HDFS/Spark (Knox) gateway** – HTTPS-based endpoint used for accessing services like webHDFS and Livy.
- **Master Instance** – TDS endpoint for database tools and application connections.

Features

Data Virtualization

Integrates data from different sources, location, and formats without moving the data, to create a single “virtual” data layer.

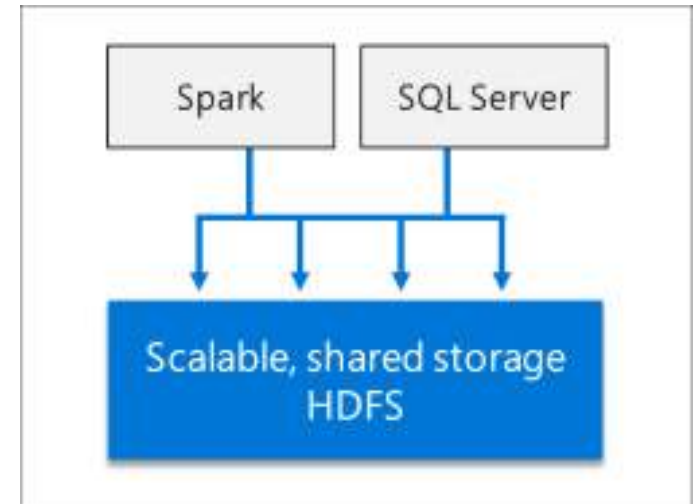
This virtual data layer, also known as “**data hub**”, allows users to query data from many sources via a single interface (i.e. Azure Data Studio)



Source: Microsoft

Data Lake

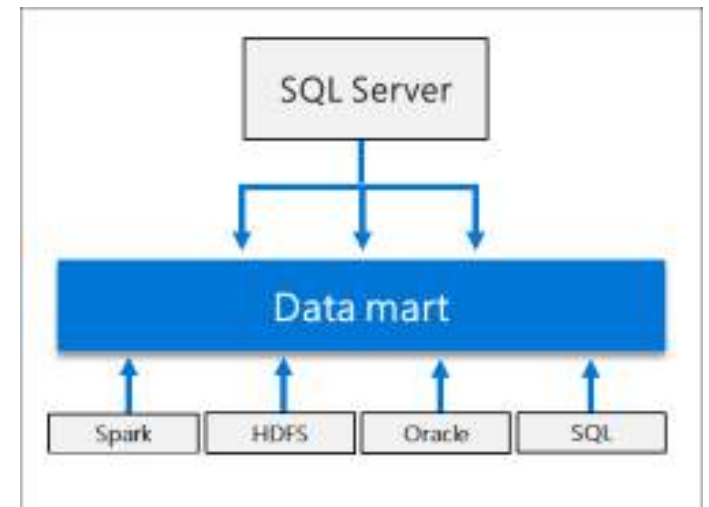
A SQL Server big data cluster scalable HDFS *storage pool* can be used to store big data, ingested from multiple external sources.



Source: Microsoft

Data Mart

Data from a variety of sources can be ingested and distributed across *data pool* nodes as a cache for further analysis.



Source: Microsoft

Mo's Thoughts

Why You Should Be Excited?

- One Stop Shop
- What's Most Unique about BDCs
- Becoming Relevant

<https://MohammadDarab.com/bdc>

Learning Curve



Learning Path

Call To Action

What you will find at <https://MohammadDarab.com/bdc>

- 1.This session
- 2.A curated list of my Big Data Cluster blog posts
- 3.Excel checklist, “Learning Path to BDC”, that itemizes the following topics:
 - 1.Linux
 - 2.Docker Containers
 - 3.Kubernetes

Thank You!

Questions & Answers

Blog: MohammadDarab.com

Twitter: [@mwdarab](https://twitter.com/mwdarab)

Feedback: MohammadDarab.com/feedback