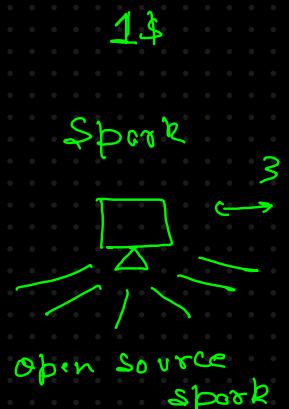
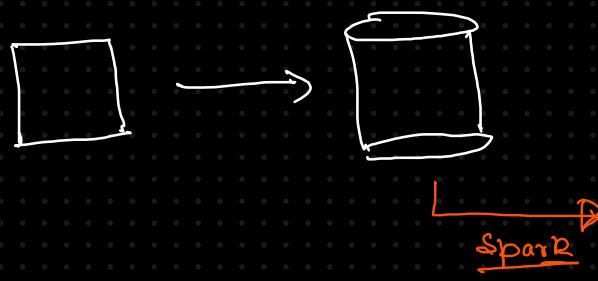


Databricks

Databricks is a cloud-based, managed data analytics platform built on Apache Spark

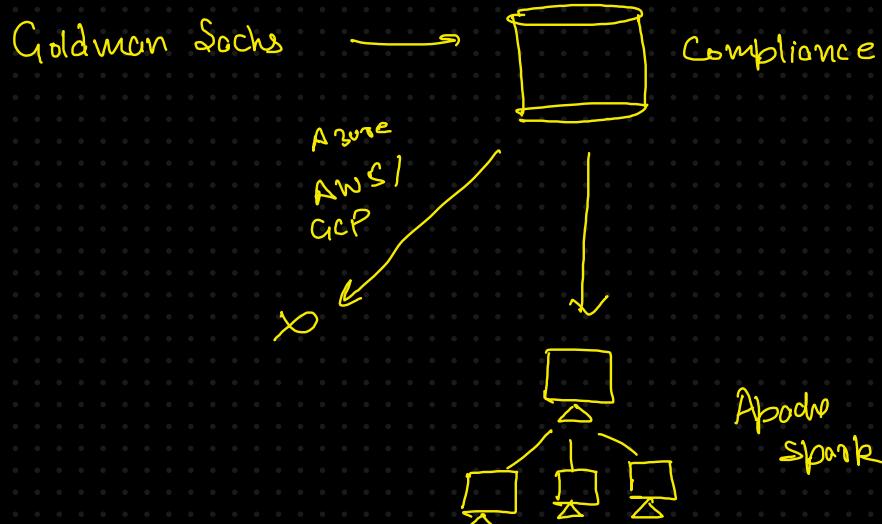
interactive workspace ← data eng.
data analyst



Databricks is created by inventors of Spark.

Spark → Apache Foundation

↓
Databricks (Part)



Open Source Spark

Complex Infrastructure setup

Manual software install & updates

Lack of user interface

Difficult security management

Version compatibility issue

Databricks

Fully managed cluster

Auto configured environment

Web-based notebooks

Built-in security & governance

Optimized spark runtime

Datalake architecture

Hardware + Software

AWS {
GCP }
Azure } m/c
datawrkhs

Workspace

Notebooks

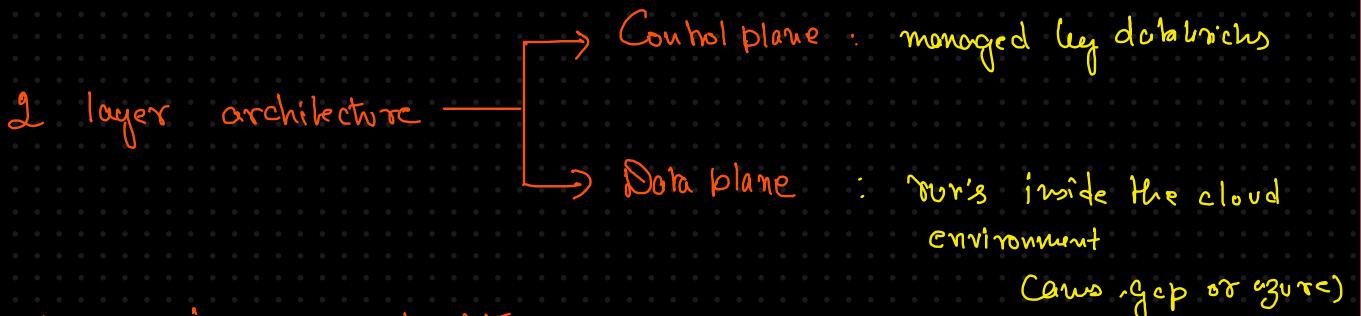
Feature	Description
Workspace	Organizes notebooks, libraries, and files.
Clusters	Manages Spark clusters for processing.
Jobs	Runs scheduled workloads and ETL tasks.
Data	Stores datasets and connects external storage.
Notebooks	Interactive coding interface (supports Python, Scala, SQL, R).

Cluster Type	Purpose
All-Purpose Cluster	Interactive work (for development, shared by multiple users).
Job Cluster	Runs specific jobs and shuts down after completion (cost-effective).
Cluster Pool	Manages pre-warmed clusters for faster job execution.

Understanding the databricks architecture

1. How computations are executed
2. How data is stored and accessed
3. How databricks differ from open source spark

High level architecture:



* However in Community edition, a free tier version is provided by databricks & it's not integrated with any service

Control Plane

UI

cluster manager

Job scheduler

Notebook execution

Workspace management

Data Plane

Driver node
worker node

DBFS

Compute resources

Feature	Community Edition	Cloud-Based Databricks
Cluster Type	Single-node cluster (only a Driver)	Multi-node cluster with Worker nodes
Data Processing	Limited to small datasets	Scales for big data workloads
Storage	Uses DBFS (local filesystem)	Can integrate with AWS S3, Azure Blob, GCS
Cloud Connectivity	No cloud support	Supports AWS, GCP, Azure
Best For	Learning, small-scale testing	Production workloads, enterprise applications

Data Bricks File System (DBFS)

fat32 ↔ ntfs
ext
app
hdfs

DBFS is a distributed file system in databricks that allows users to store, manage and interact with files

has to be enabled
& then refresh.

Why use DBFS?

- act as an abstraction layer over the cloud storage (GCS, AWS, Azure)
- supports structured and unstructured data (csv, parquet, images)
- integrates very nicely with spark, DB Notebook & delta lake
- provides a simple file system interface

