



Get Started with Databricks Data Science & Engineering Workspace

>_



Module Agenda

Get Started with Databricks Data Science and Engineering Workspace

Introduction to the Databricks Lakehouse Platform

Databricks Architecture and Services

Demo – Navigating the Workspace

DE 1.1 – Create and Manage Clusters Interactively

DE 1.2 – Notebook Basics

Git Versioning with Databricks Repos

Demo – Using Databricks Repos

DE 1.3L – Getting Started with the Databricks Lakehouse Platform Lab

Using the Databricks Lakehouse Platform

Learning Objectives

- Describe the components of the Databricks Lakehouse
- Complete basic code development tasks using services of the Databricks Data Science and Engineering Workspace
- Perform common table operations using Delta Lake in the Lakehouse

Using the Databricks Lakehouse Platform

Agenda

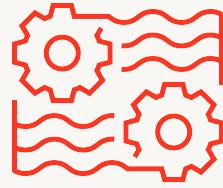
- Introduction to the Databricks Lakehouse Platform
- Introduction to the Databricks Workspace and Services
 - Using clusters, files, notebooks, and repos
- Introduction to Delta Lake
 - Manipulating and optimizing data in Delta tables





The Databricks Lakehouse Platform





Lakehouse

One simple platform to unify all of
your data, analytics, and AI workloads

Customers

7000+

across the globe



Original creators of:



Supporting enterprises in every industry

Healthcare & Life Sciences



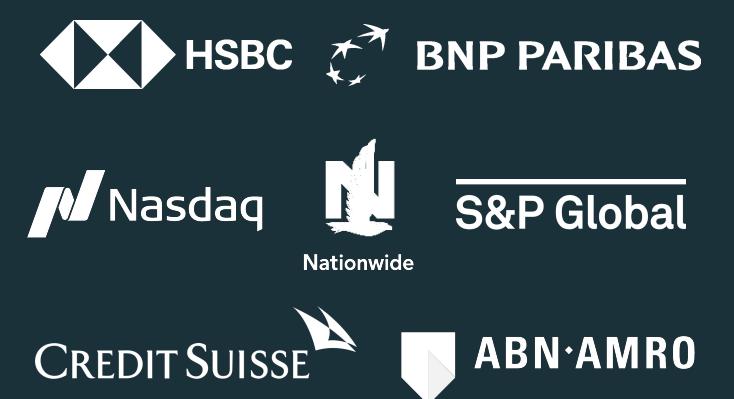
Manufacturing & Automotive



Media & Entertainment



Financial Services



Public Sector



Retail & CPG



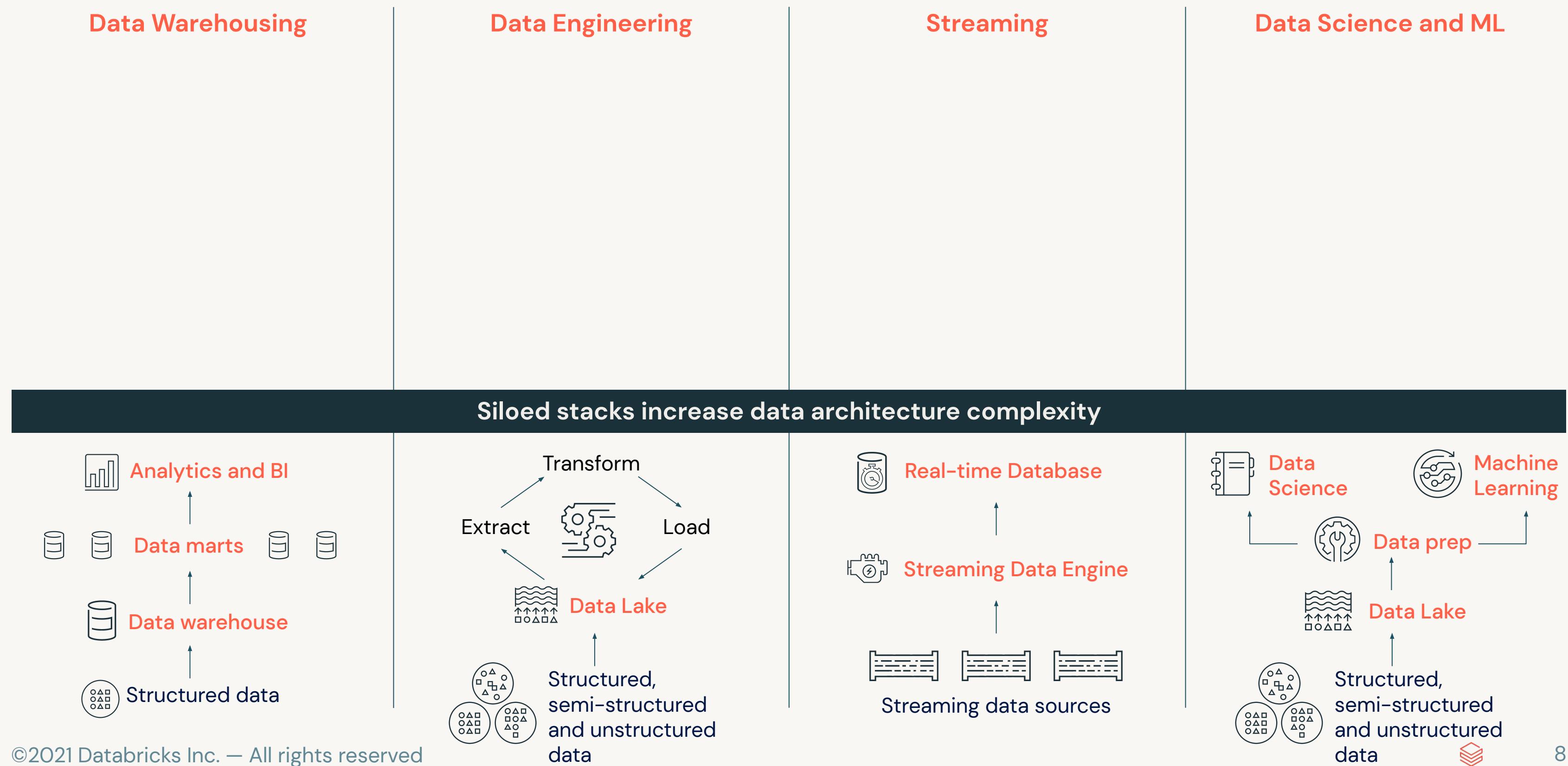
Energy & Utilities



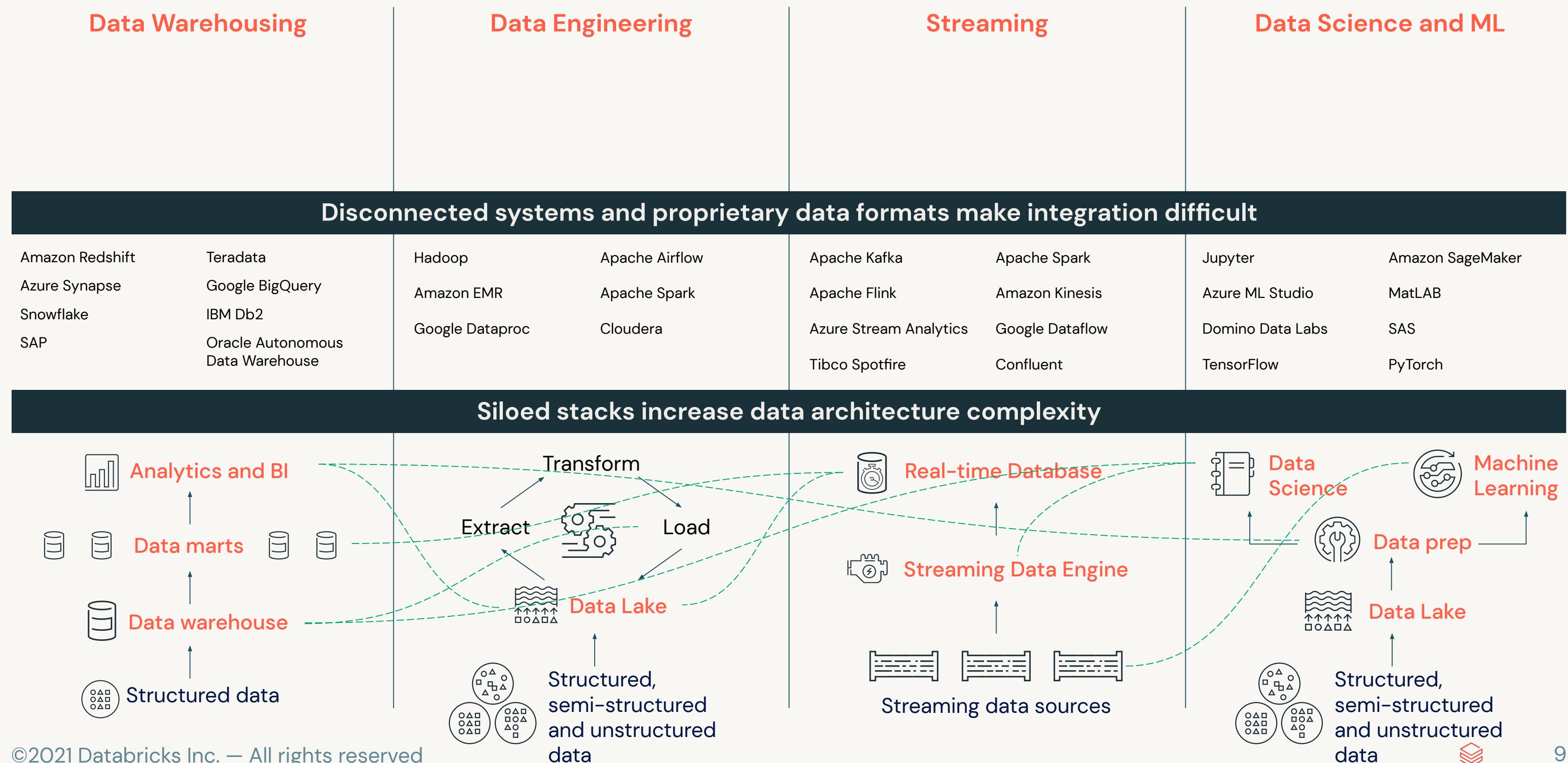
Digital Native



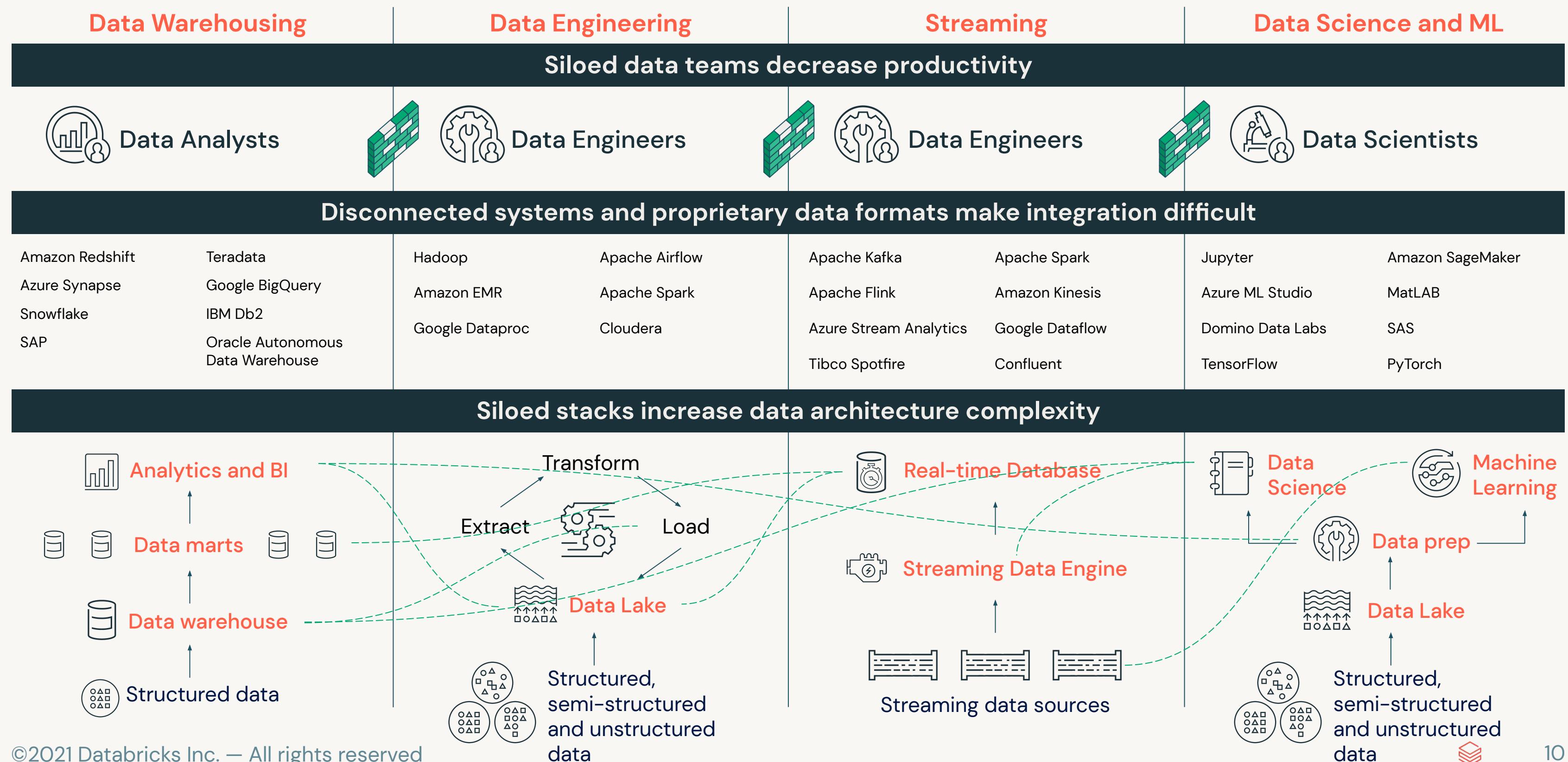
Most enterprises struggle with data

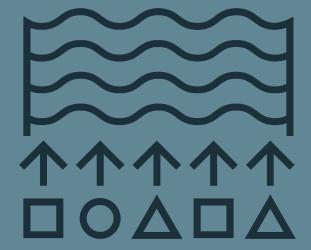


Most enterprises struggle with data



Most enterprises struggle with data

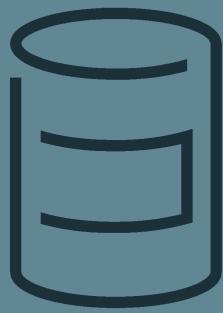




Data
Lake

Lakehouse

One platform to unify all of
your data, analytics, and AI
workloads



Data
Warehouse



Data Lake



Data Warehouse



DELTA LAKE

An open approach to bringing
data management and
governance to data lakes

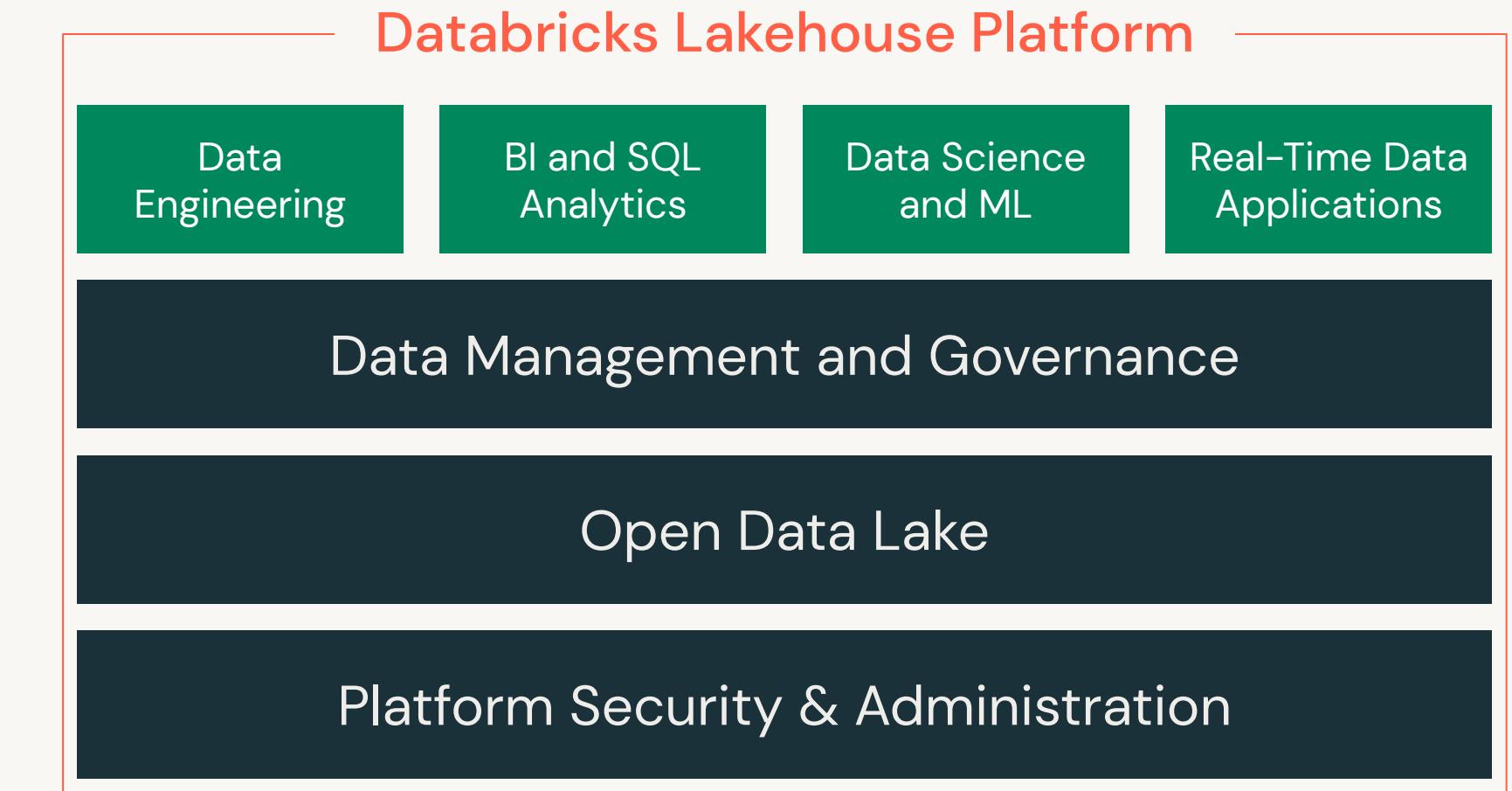
Better reliability with transactions

48x faster data processing with
indexing

Data governance at scale with
fine-grained access control lists

The Databricks Lakehouse Platform

-  Simple
-  Open
-  Collaborative



Unstructured, semi-structured, structured, and streaming data



The Databricks Lakehouse Platform



Simple

Unify your data, analytics, and AI on one common platform for all data use cases

Databricks Lakehouse Platform

Data Engineering

BI and SQL Analytics

Data Science and ML

Real-Time Data Applications

Data Management and Governance

Open Data Lake

Platform Security & Administration



Unstructured, semi-structured, structured, and streaming data



Microsoft Azure



Google Cloud

The Databricks Lakehouse Platform



Open

Unify your data ecosystem with open source standards and formats.

Built on the innovation of some of the most successful open source data projects in the world

30 Million+
Monthly downloads



mlflow™



re'dash

The Databricks Lakehouse Platform

Open

Unify your data ecosystem with open source standards and formats.

450+

Partners across the data landscape

Visual ETL & Data Ingestion



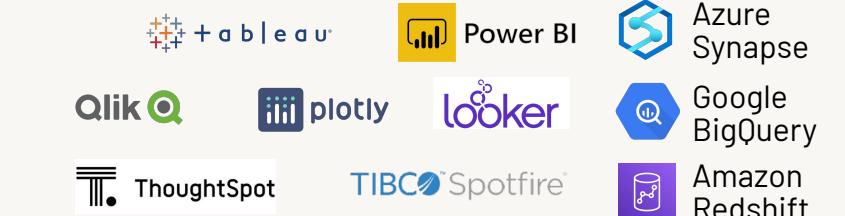
Data Providers



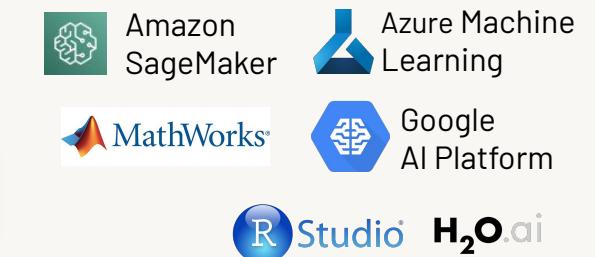
Top Consulting & SI Partners



Business Intelligence



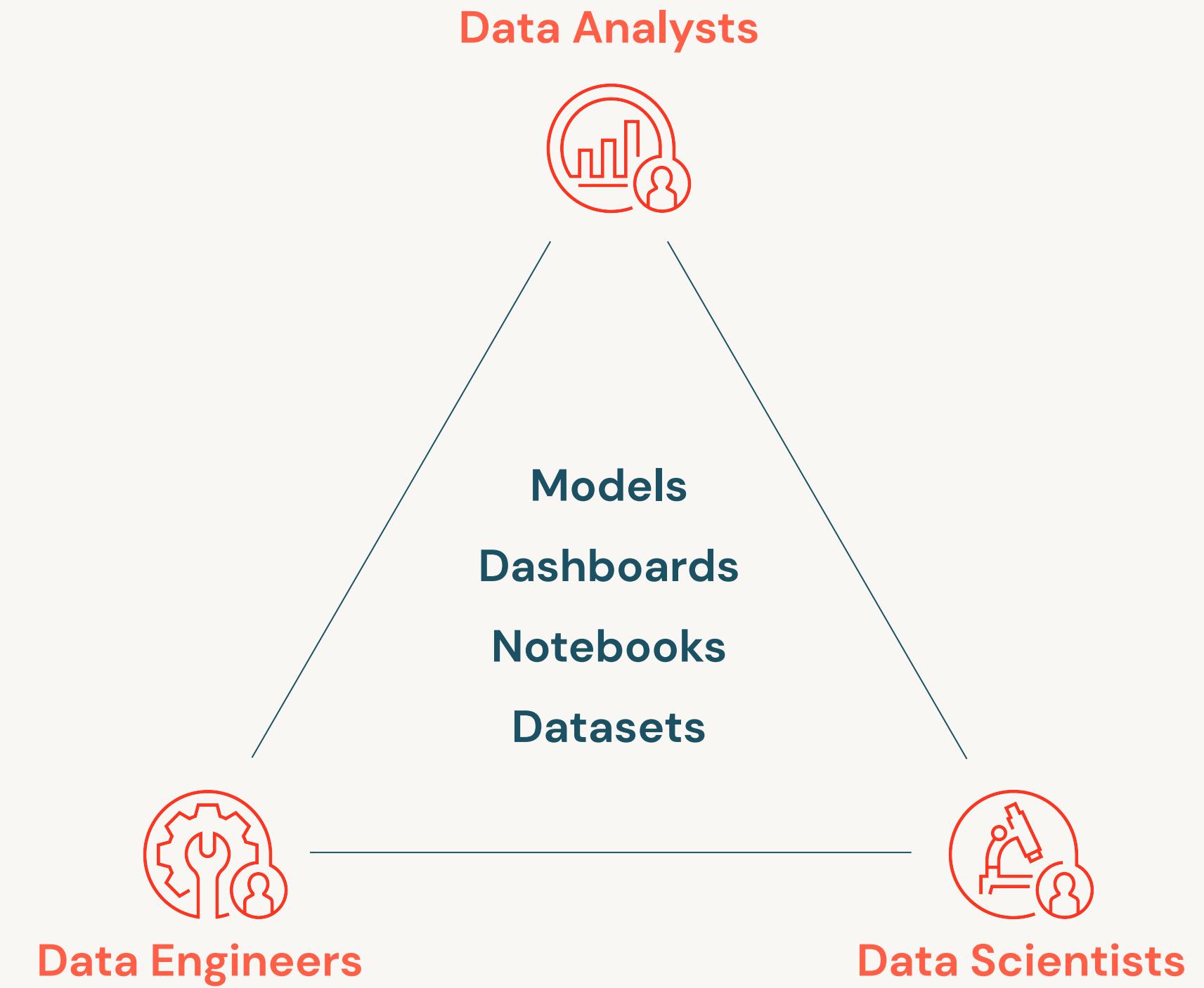
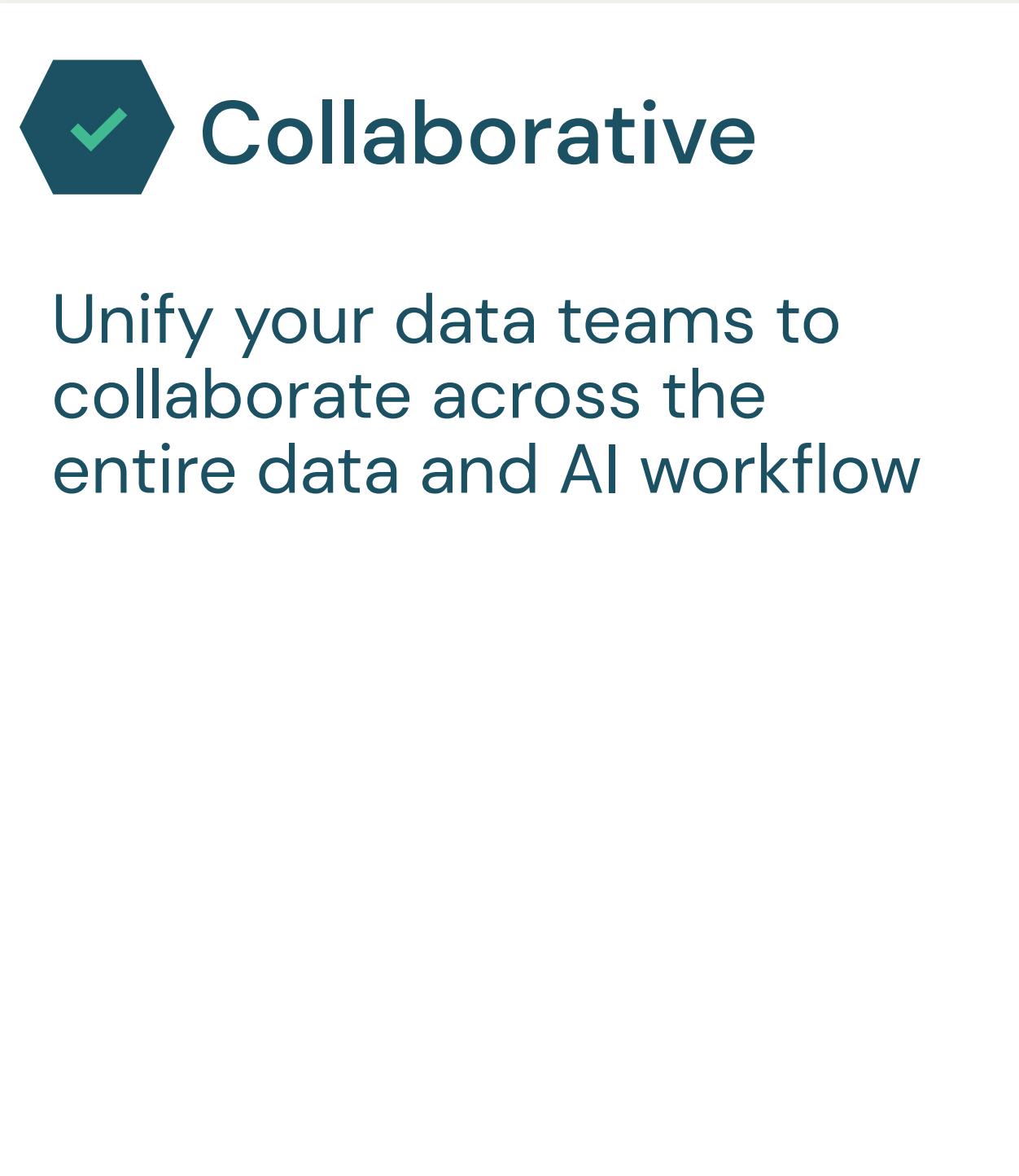
Machine Learning



Centralized Governance



The Databricks Lakehouse Platform

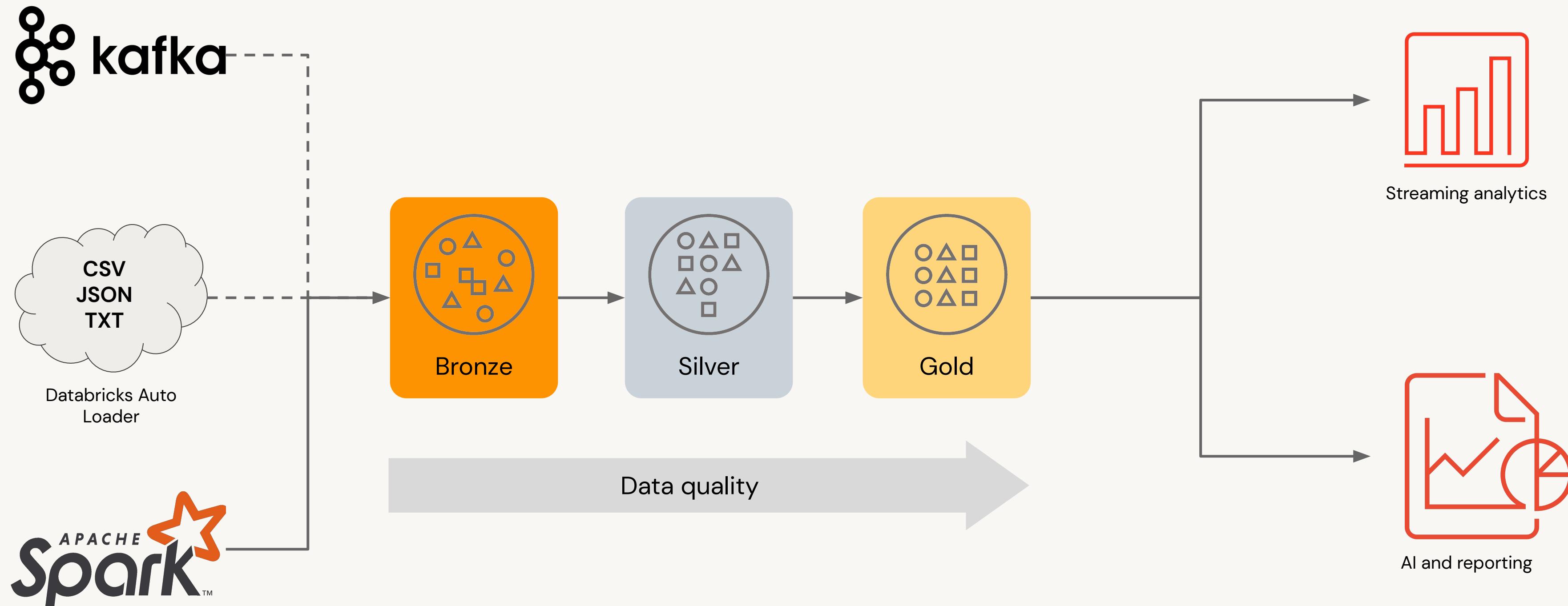




The Medallion Architecture



Multi-Hop in the Lakehouse



Multi-Hop in the Lakehouse

Bronze Layer

Typically just a raw copy of ingested data

Replaces traditional data lake

Provides efficient storage and querying of full, unprocessed history of data



Multi-Hop in the Lakehouse

Silver Layer

Reduces data storage complexity, latency, and redundancy

Optimizes ETL throughput and analytic query performance

Preserves grain of original data (without aggregations)

Eliminates duplicate records

Production schema enforced

Data quality checks, corrupt data quarantined



Multi-Hop in the Lakehouse

Gold Layer

Powers ML applications, reporting, dashboards, ad hoc analytics

Refined views of data, typically with aggregations

Reduces strain on production systems

Optimizes query performance for business-critical data

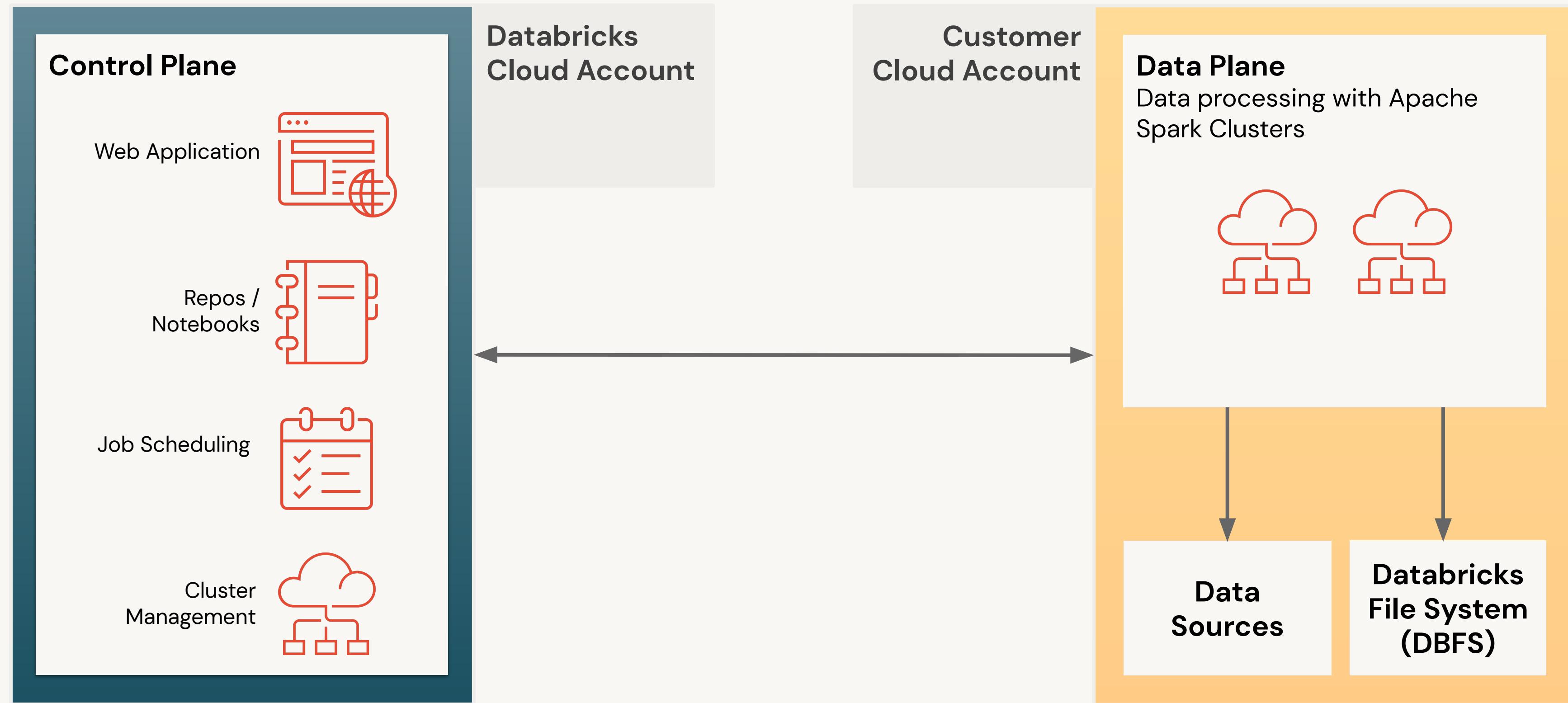




Databricks Architecture and Services



Databricks Architecture



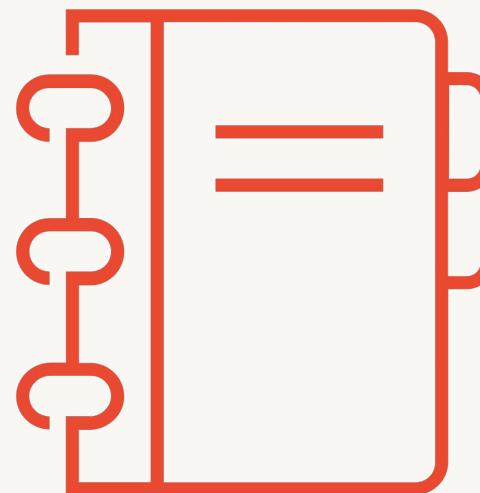
Databricks Services

Control Plane in Databricks

Manage customer accounts, datasets, and clusters



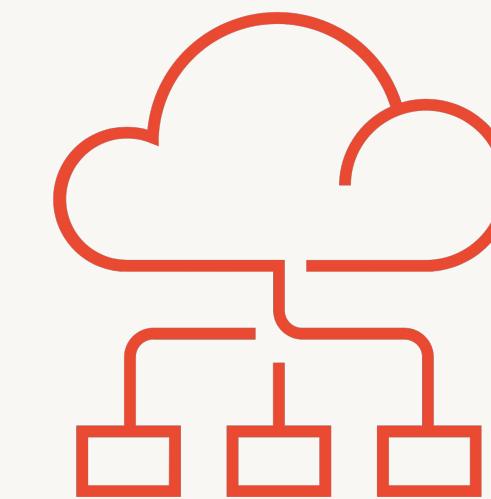
Databricks Web
Application



Repos /
Notebooks

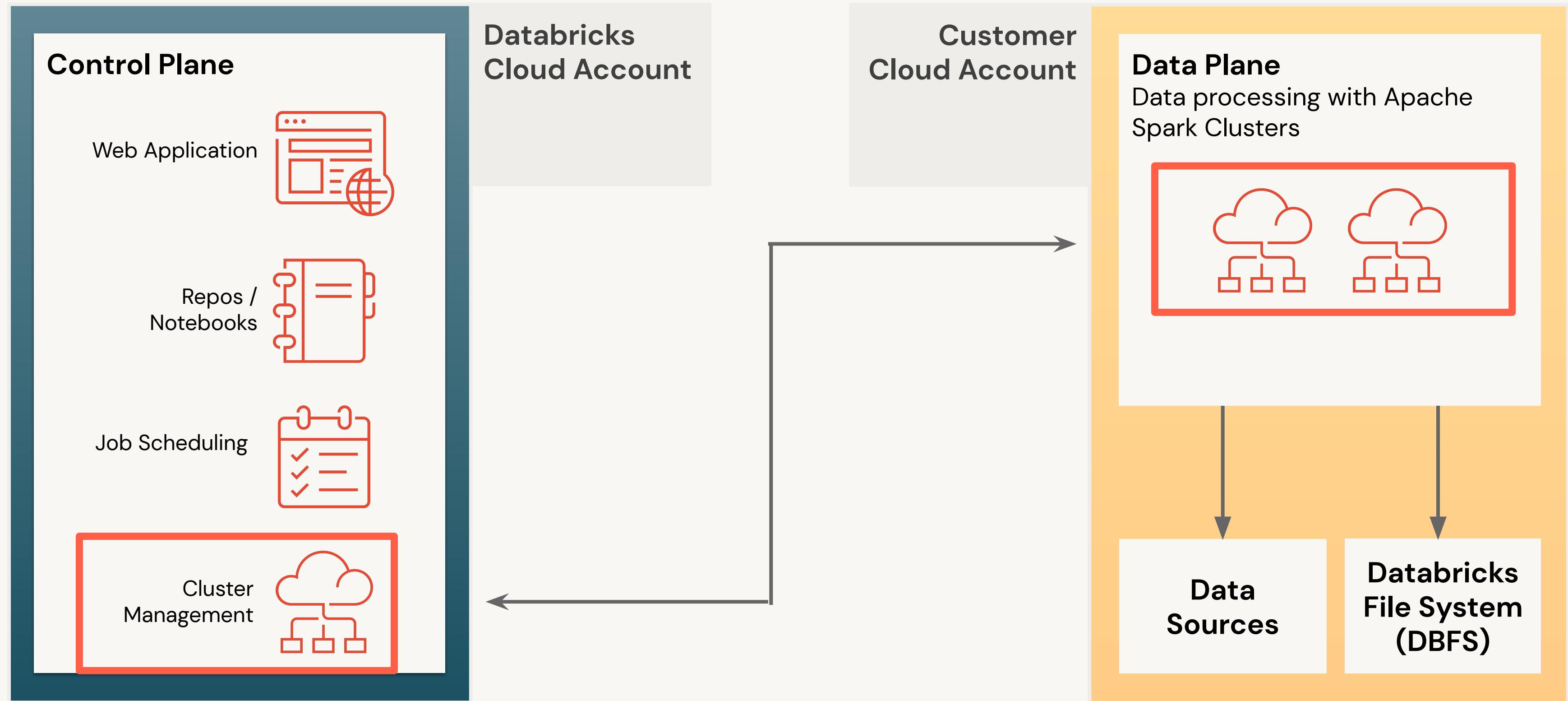


Jobs



Cluster
Management

Clusters



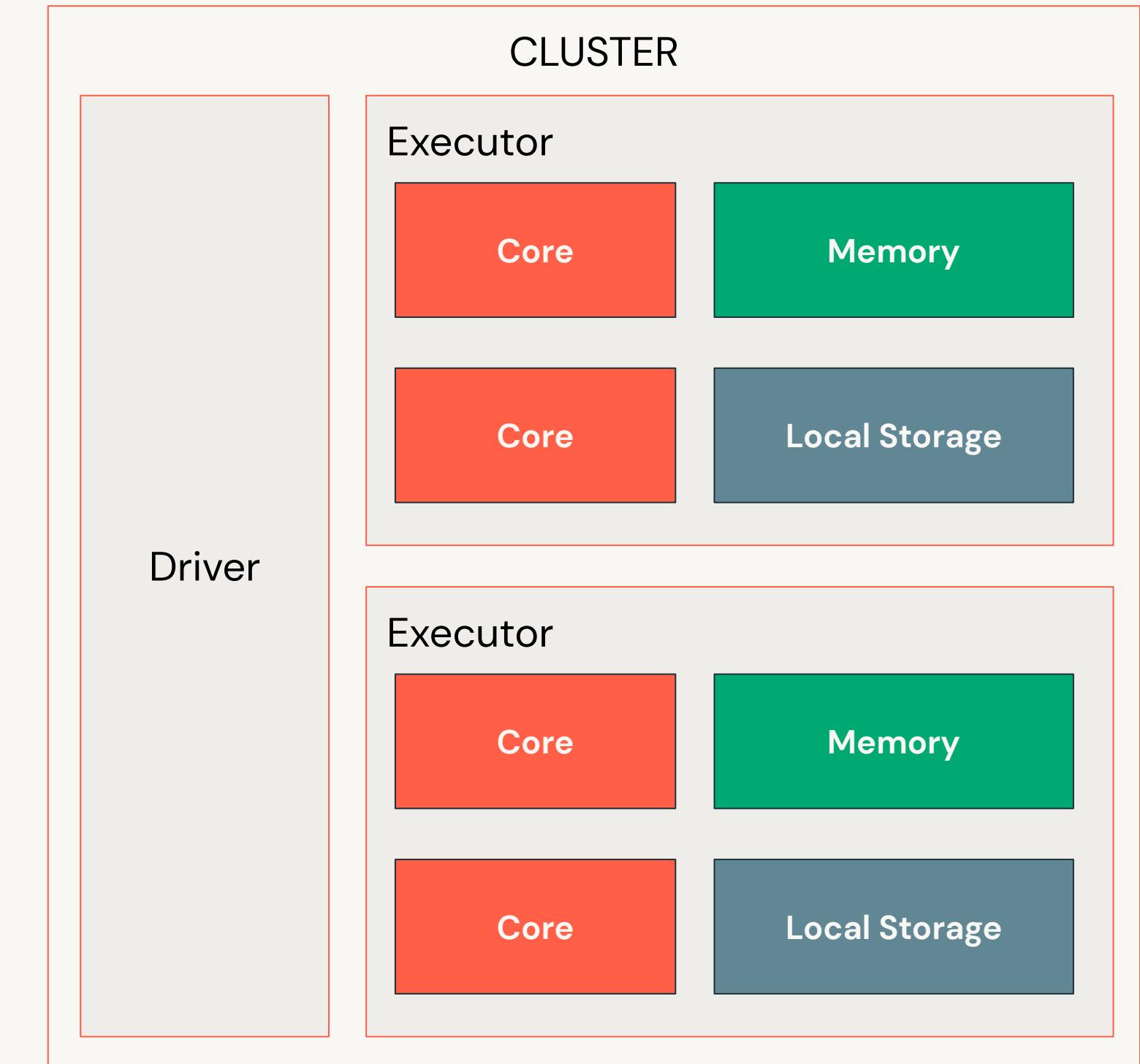
Clusters

Overview

Clusters are made up of one or more virtual machine (VM) instances

Driver coordinates activities of executors

Executors run tasks composing a Spark job



Clusters

Types

All-purpose Clusters

Analyze data collaboratively using interactive notebooks

Create clusters from the Workspace or API

Retains up to 70 clusters for up to 30 days.

Job Clusters

Run automated jobs

The Databricks job scheduler creates job clusters when running jobs.

Retains up to 30 clusters.





Git Versioning with Databricks Repos



Databricks Repos

Overview

Git Versioning

Native integration with Github, Gitlab, Bitbucket and Azure Devops

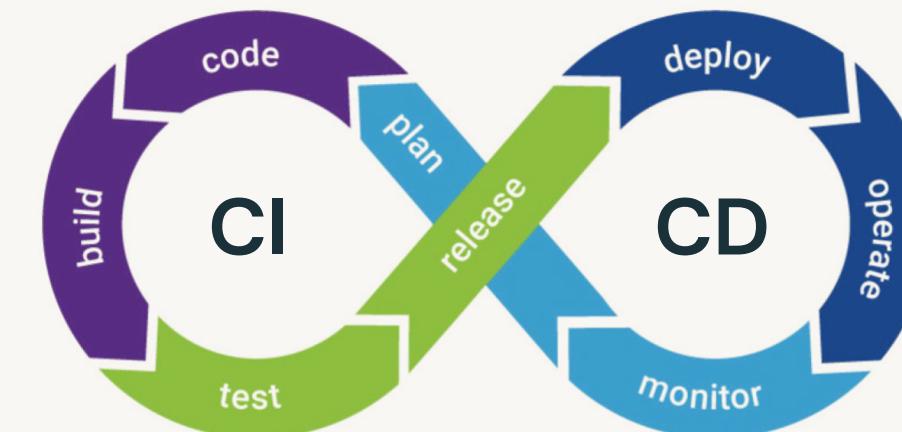
UI-based workflows



CI/CD Integration

API surface to integrate with automation

Simplifies the dev/staging/prod multi-workspace story



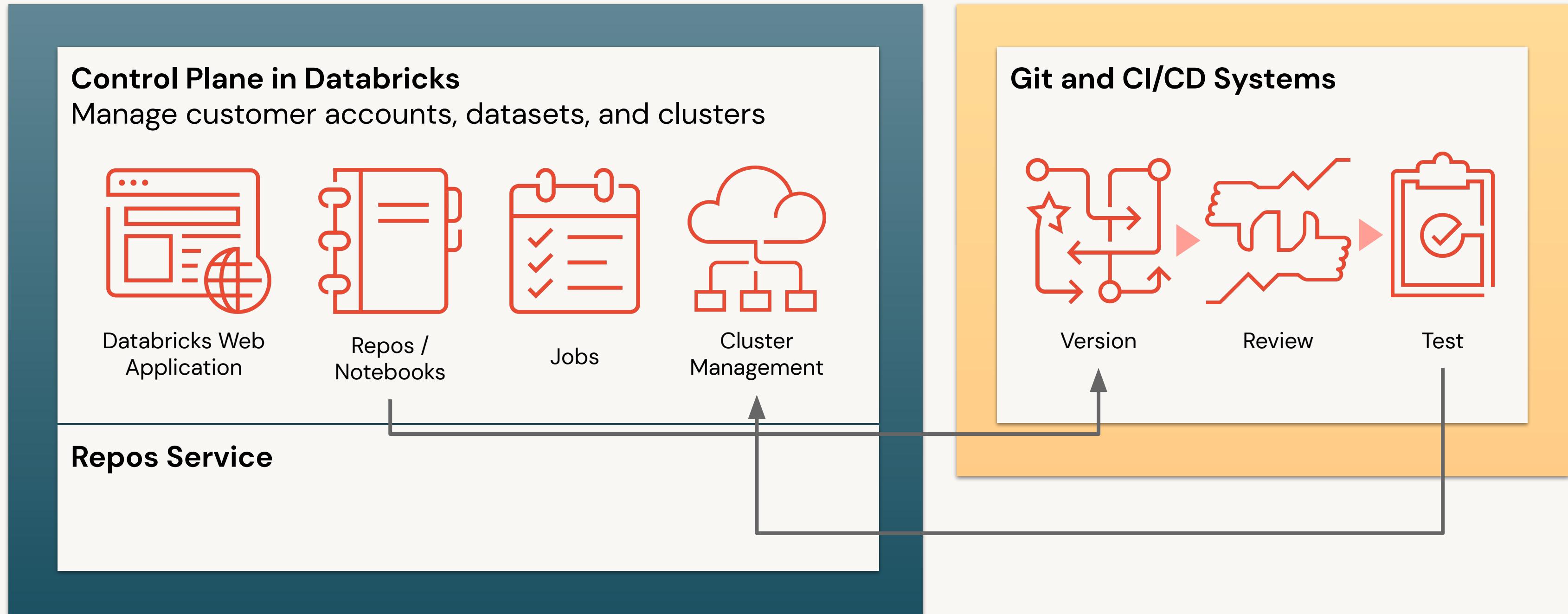
Enterprise ready

Allow lists to avoid exfiltration

Secret detection to avoid leaking keys

Databricks Repos

CI/CD Integration



Databricks Repos

Best practices for CI/CD workflows

