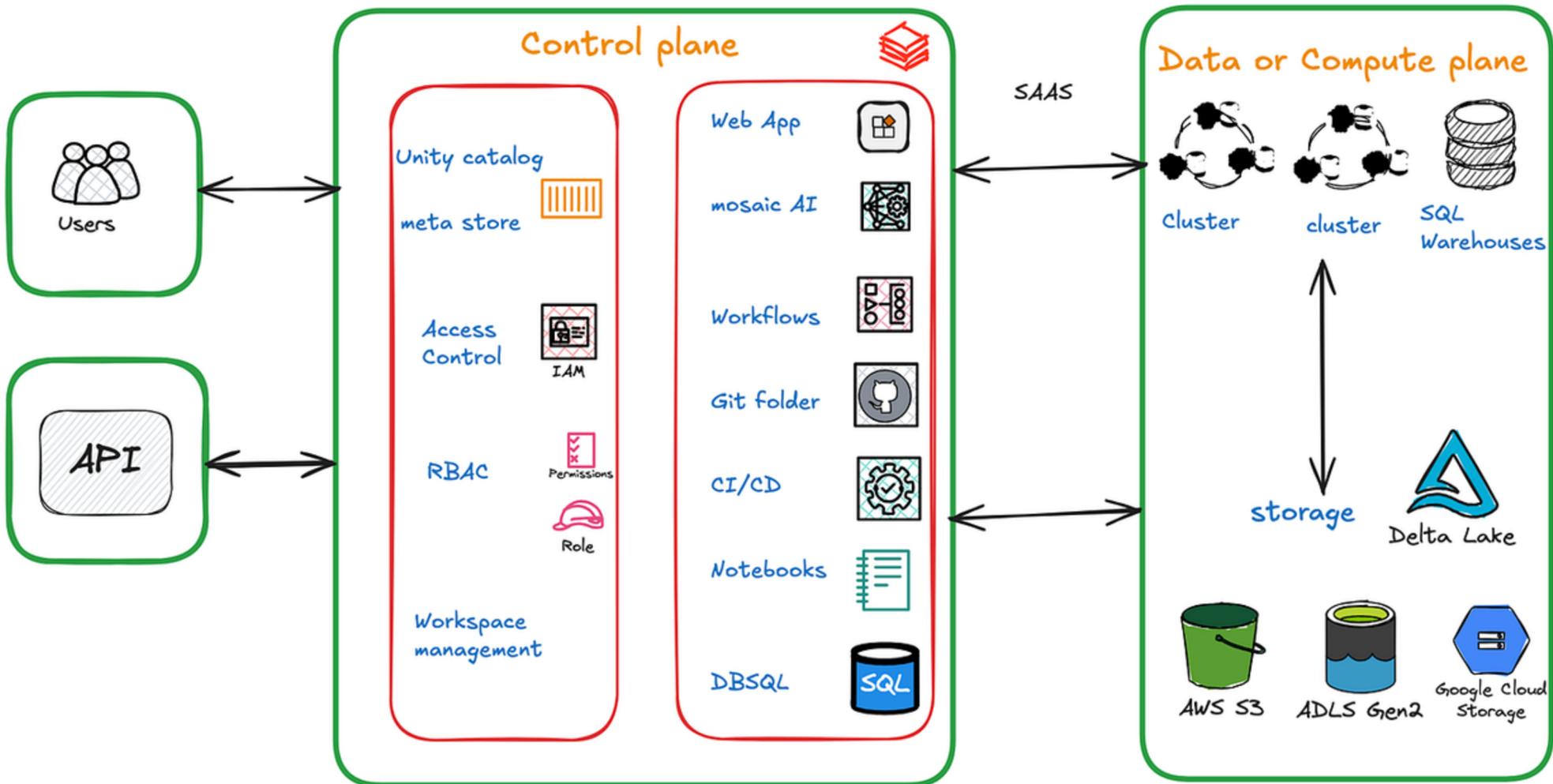


Databricks Architecture

In simple words..



Databricks Architecture



Databricks Architecture Overview

- Databricks is a cloud-based data platform built on Apache Spark that helps you handle big data, machine learning, and analytics in one place.
- It sits on top of cloud providers like Azure, AWS, or GCP and provides a unified workspace.



Main Components of Databricks Architecture

- 1. Control Plane**
- 2. DataPlane**
- 3. Workspace**
- 4. Clusters**
- 5. Delta Lake (Storage Layer)**

1. Control Plane

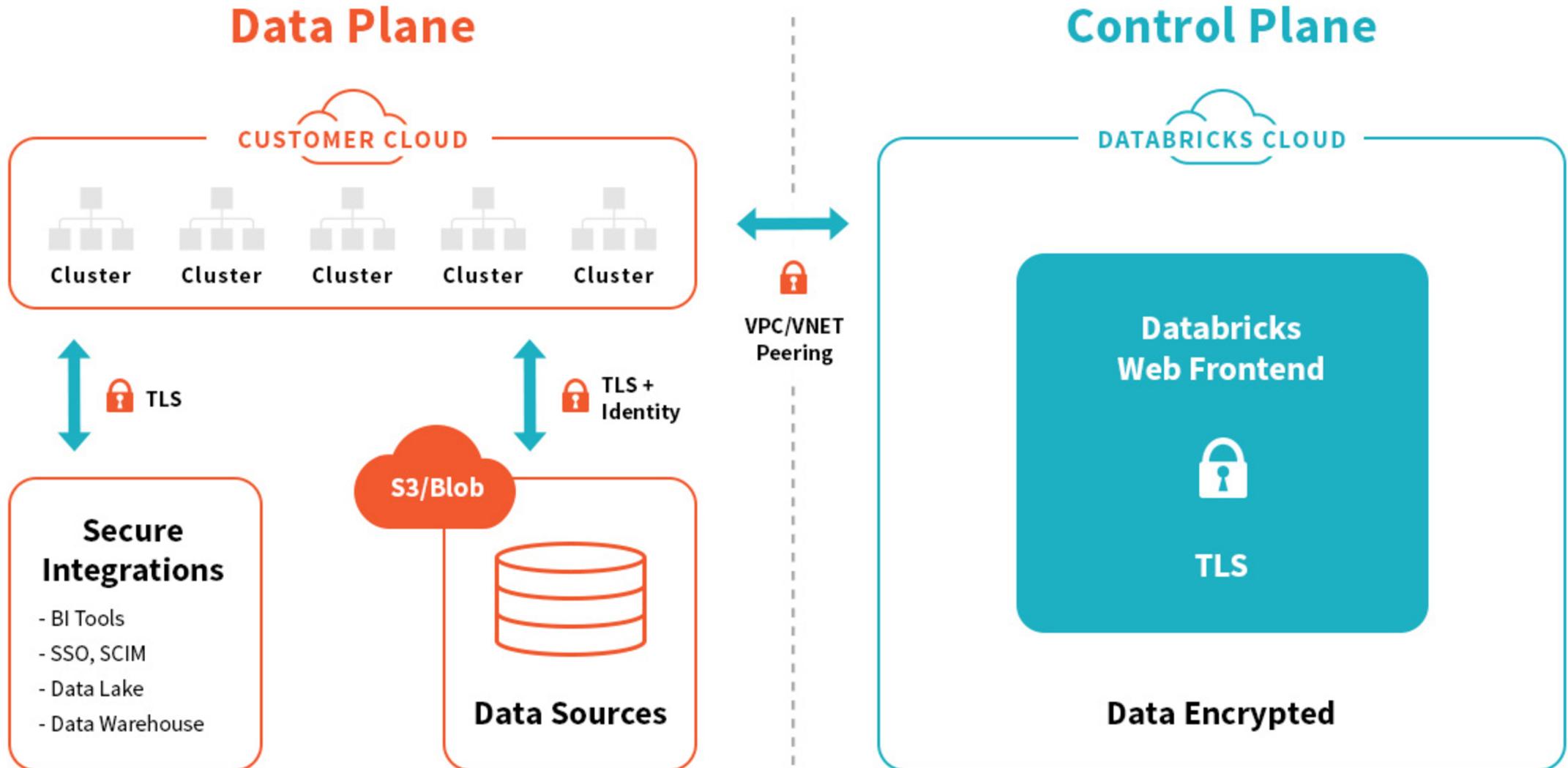
- Managed by Databricks(you don't manage it).
- Responsible for authentication, jobs scheduling, notebooks,
- REST APIs, UI, cluster management.
- Stores metadata (but not your data).



Think of this as the “brain” that manages everything.

2. Data Plane

- This runs inside your cloud account(Azure, AWS, GCP).
 - Your data and compute clusters live here.
 - Databricks spins up Spark clusters on-demand in your cloud environment.
-  Think of this as the “muscle” where actual data processing happens.



3. Workspace

- A collaborative **UI/Notebook environment** for developers, data scientists, and analysts.
 - Supports **Python, SQL, R, Scala** in the same notebook.
-  Example: A Data Engineer writes ETL in PySpark, while a Data Analyst writes SQL in the same workspace.

Databricks Workspace



Clusters



Notebooks



Jobs



Data

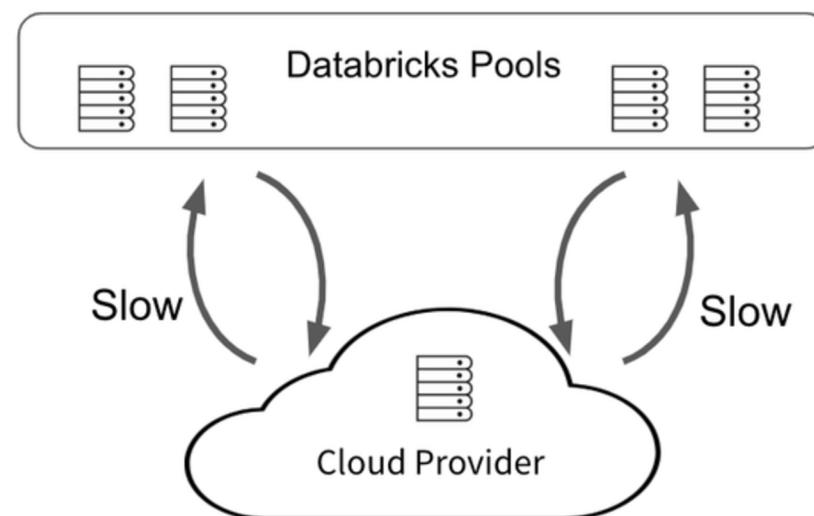
4. Clusters

- Virtual machines(Spark Runtime) that run in your cloud environment.
- **Types:**
 - Interactive clusters:** used for exploration/analysis.
 - Job clusters:** auto-created for scheduled ETL/ML jobs.

Automated Job

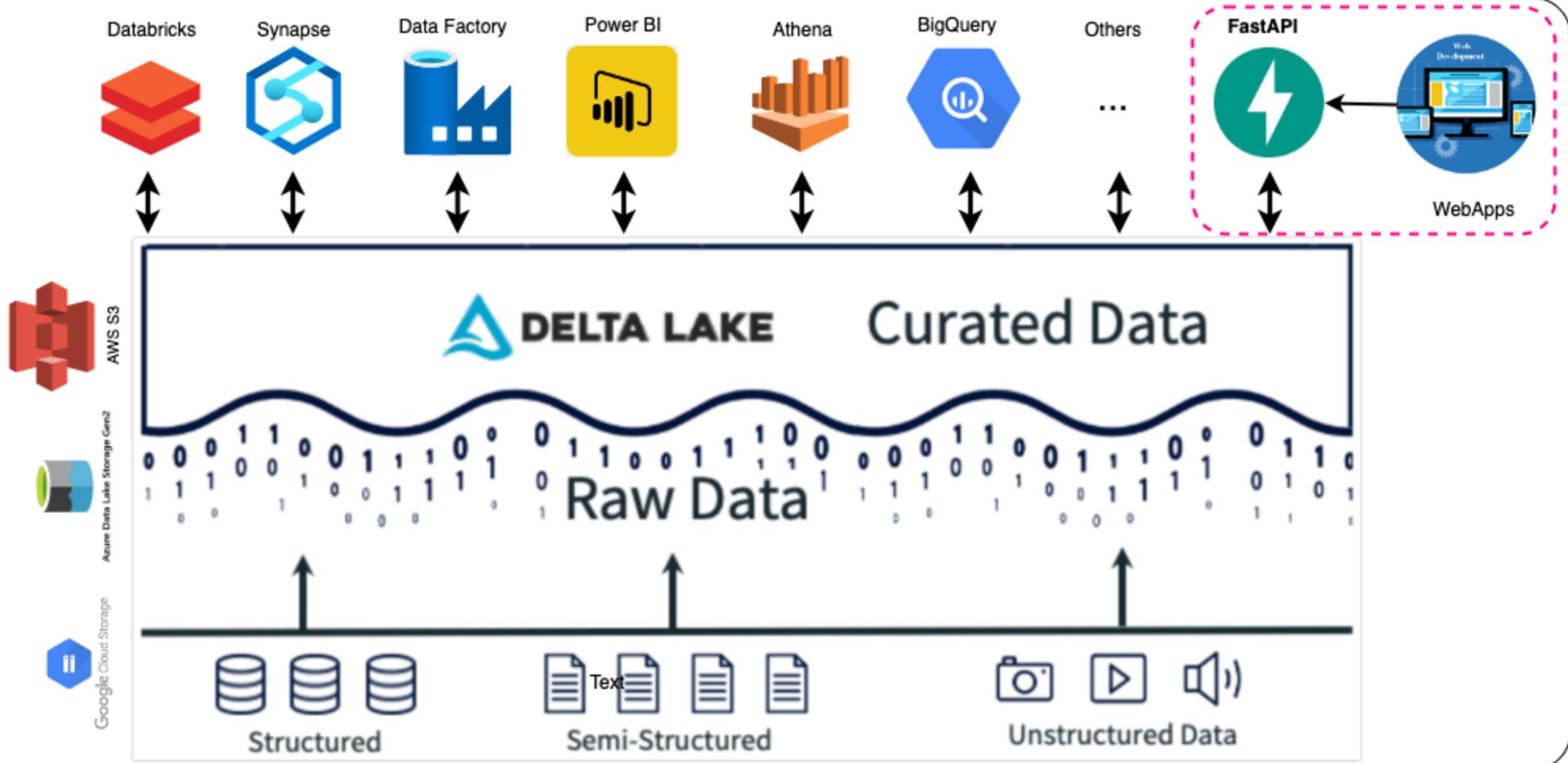


Interactive Scale Up



5. Delta Lake(Storage Layer)

- An open-source storage format on top of cloud storage (S3, ADLS, GCS).
- Provides ACID transactions, schema enforcement, time travel.
- Makes data lakes reliable→ called a Lakehouse.



6. Databricks Runtime

- Pre-packaged Spark environment optimized by Databricks.
- Includes ML flow, Koalas(pandas API),Delta Lake, libraries for ML/AI.

Databricks runtime

The screenshot shows a dropdown menu titled "Databricks runtime" with the selected option "17.1" highlighted. Below the dropdown is a search bar with the placeholder "Search". The list of runtime versions is as follows:

Version	Scala / Spark Version
17.1	Scala 2.13, Spark 4.0.0
17.0	Scala 2.13, Spark 4.0.0
16.4 LTS (Scala 2.13)	Scala 2.13, Spark 3.5.2
16.4 LTS (Scala 2.12)	Scala 2.12, Spark 3.5.2
16.3	Scala 2.12, Spark 3.5.2
16.2	Scala 2.12, Spark 3.5.2
15.4 LTS	Scala 2.12, Spark 3.5.0
14.3 LTS	Scala 2.12, Spark 3.5.0
13.3 LTS	Scala 2.12, Spark 3.4.1
12.2 LTS	Scala 2.12, Spark 3.3.2
11.3 LTS	Scala 2.12, Spark 3.3.0

Picking a Databricks Runtime Version



Key Benefit of Databricks Architecture

- Unified platform (ETL +Analytics +ML)
- Secure separation(control plane vs data plane)
- Scalable(auto-scaling Spark clusters)
- Reliable storage (Delta Lake)

Integrated Data Services

-  Azure Data Factory
-  Azure Data Lake Storage
-  Azure Blob Storage
-  Azure Event Hubs
-  Azure Cosmos DB

Azure Databricks

End-to-End Analytics and ML

-  Azure Synapse Analytics
-  PowerBI
-  Azure Machine Learning

Integrated Management



Azure Security
Azure Active Directory
Single Sign-On, Identity
Passthrough, Network



Azure Portal
1-Click Setup
Unified Billing



Azure DevOps
Notebook integration