# The Delta Lake transaction log

A key to understanding how Delta Lake provides all these capabilities is the transaction log. The Delta Lake transaction log is the common thread that runs through many of Delta Lake's most notable features, including ACID transactions, scalable metadata handling, time travel and more. The Delta Lake transaction log is an ordered record of every transaction that has ever been performed on a Delta Lake table since its inception.

Delta Lake is built on top of Spark to allow multiple readers and writers of a given table to work on a table at the same time. To always show users correct views of the data, the transaction log serves as a single source of truth: the central repository that tracks all changes that users make to the table.

When a user reads a Delta Lake table for the first time or runs a new query on an open table that has been modified since the last time it was read, Spark checks the transaction log to see what new transactions are posted to the table. Then, Spark updates the table with those recent changes. This ensures that a user's version of a table is always synchronized with the master record as of the most recent query, and that users cannot make divergent, conflicting changes to a table.

# Flexibility and broad industry support

Delta Lake is an open source project, with an engaged community of contributors building and growing the Delta Lake ecosystem atop a set of open APIs and is part of the Linux Foundation. With the growing adoption of Delta Lake as an open storage standard in different environments and use cases, comes a broad set of integration with industry–leading tools, technologies and formats.

Organizations leveraging Delta Lake on the Databricks Data Intelligence Platform gain flexibility in how they ingest, store and query data. They are not limited in storing data in a single cloud provider and can implement a true multicloud approach to data storage.

Connectors to tools, such as Fivetran, allow you to leverage Databricks' ecosystem of partner solutions, so organizations have full control of building the right ingestion pipelines for their use cases. Finally, consuming data via queries for exploration or business intelligence (BI) is also flexible and open.

databricks

# Delta Lake integrates with all major analytics tools

Eliminates unnecessary data movement and duplication

In addition to a wide ecosystem of tools and technologies, Delta Lake supports a broad set of data formats for structured, semi-structured and unstructured data. These formats include image binary data that can be stored in Delta Tables, graph data format, geospatial data types and key-value stores.

**Learn more**

Delta Lake on the Databricks Data Intelligence Platform
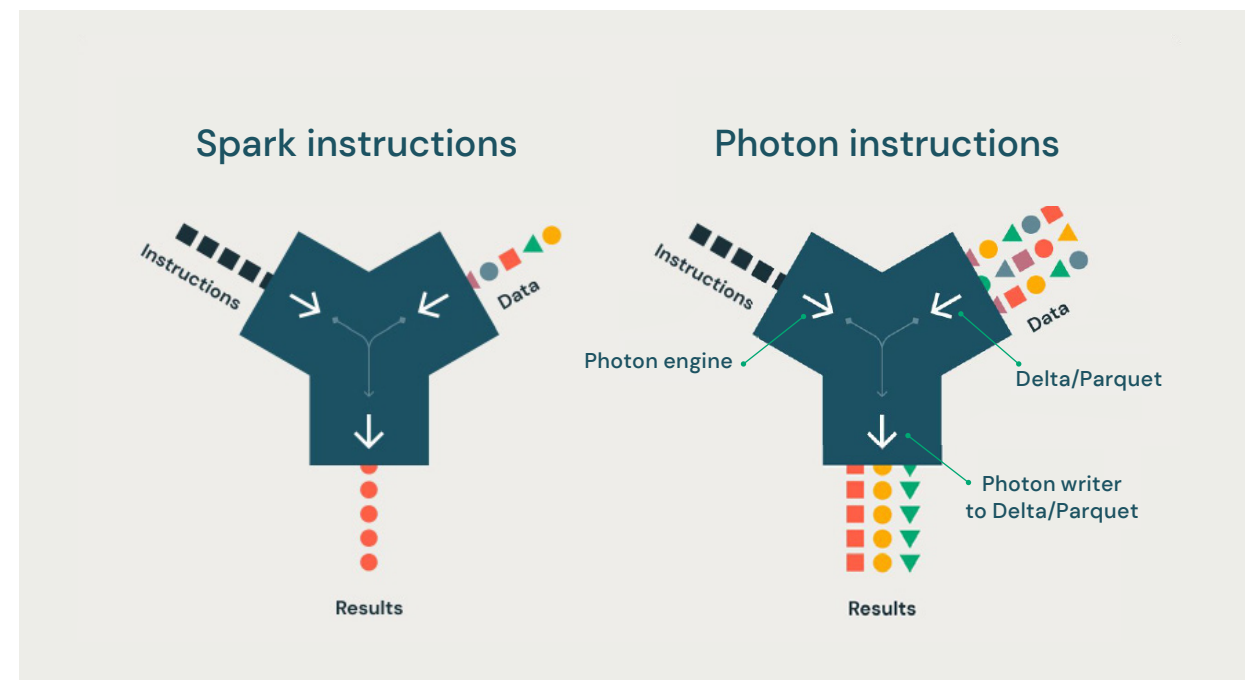
Data Intelligence Platform Documentation

Delta Lake Open Source Project

eBooks: The Delta Lake Series

## What is Photon?

As many organizations standardize on the lakehouse paradigm, this new architecture poses challenges with the underlying query execution engine for accessing and processing structured and unstructured data. The execution engine needs to provide the performance of a data warehouse and the scalability of data lakes.
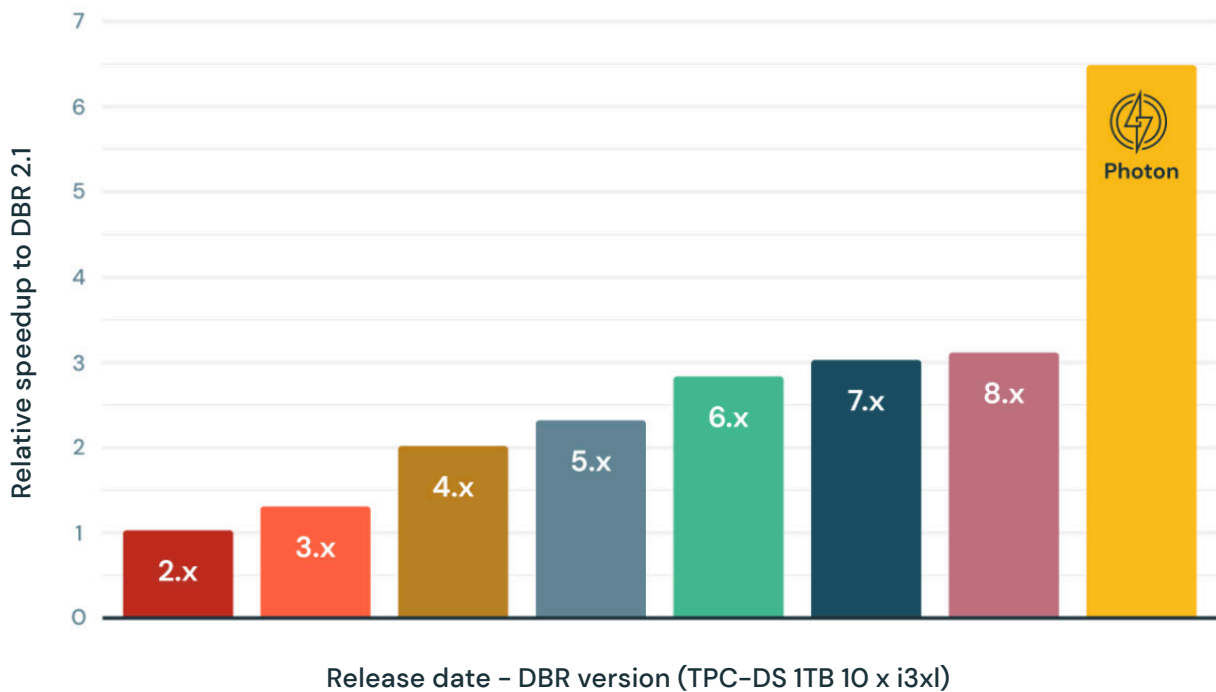
Photon is the next-generation query engine on the Databricks Data Intelligence Platform that provides dramatic infrastructure cost savings and speedups for all use cases — from data ingestion, ETL, streaming, data science and interactive queries — directly on your data lake. Photon is compatible with Spark APIs and implements a more general execution framework that allows efficient processing of data with support of the Spark API. This means getting started is as easy as turning it on — no code change and no lock-in. With Photon, typical customers are seeing up to 80% TCO savings over traditional Databricks Runtime (Spark) and up to 85% reduction in VM compute hours.



databricks

## Why process queries with Photon?

Query performance on Databricks has steadily increased over the years, powered by Spark and thousands of optimizations packaged as part of the Databricks Runtime (DBR). Photon provides an additional 2x speedup per the TPC-DS 1TB benchmark compared to the latest DBR versions.

### Relative speedup to DBR 2.1 by DBR version
Higher is better



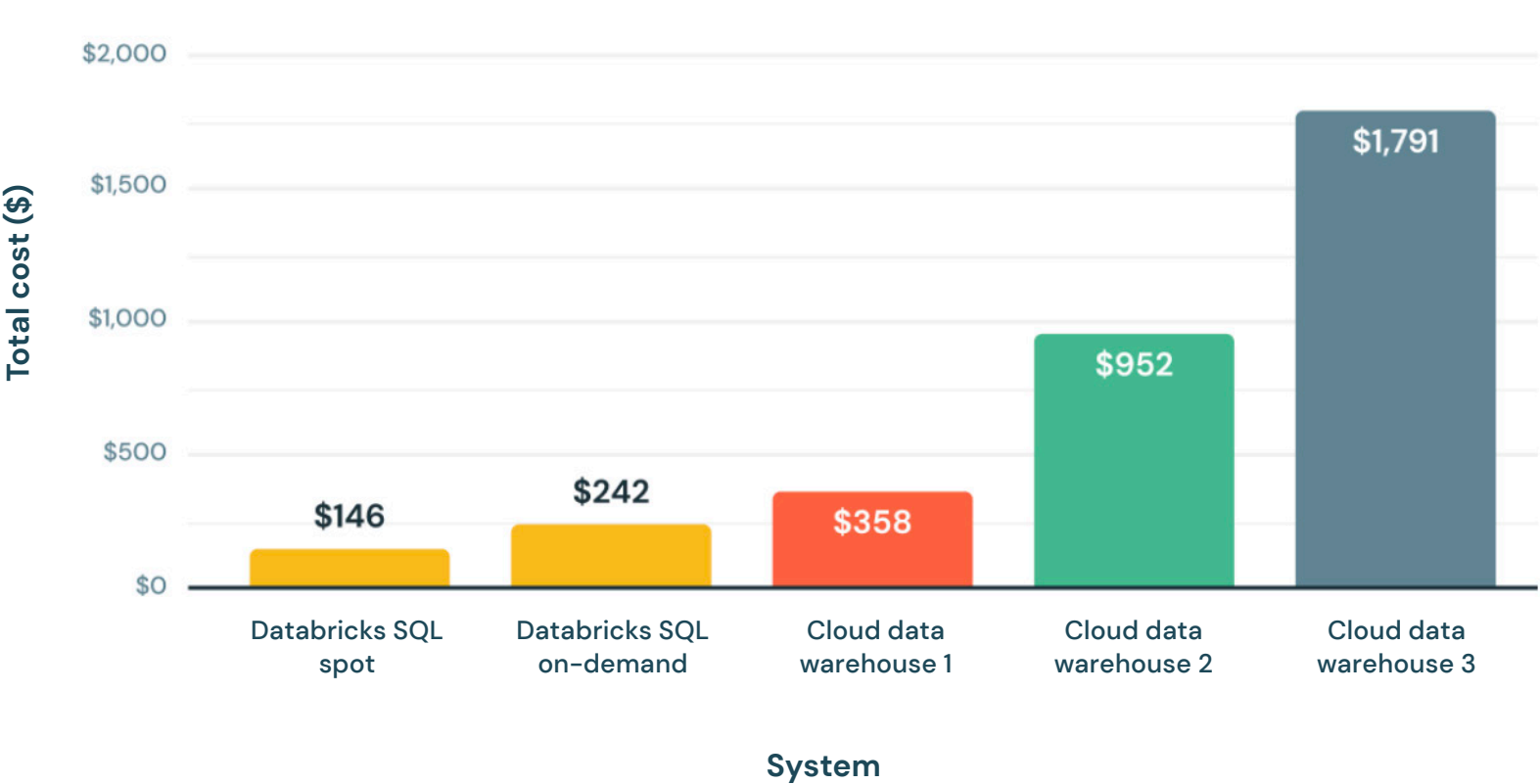Release date – DBR version (TPC-DS 1TB 10 x i3xl)

**Customers have observed significant speedups using Photon on workloads such as:**

- **SQL-based jobs:** Accelerate large-scale production jobs on SQL and Spark DataFrames

- **IoT use cases:** Faster time-series analysis using Photon compared to Spark and traditional Databricks Runtime

- **Data privacy and compliance:** Query petabytes-scale data sets to identify and delete records without duplicating data with Delta Lake, production jobs and Photon

- **Loading data into Delta and Parquet:** Vectorized I/O speeds up data loads for Delta and Parquet tables, lowering overall runtime and costs of data engineering jobs

databricks

## Best price/performance for analytics in the cloud

Written from the ground up in C++, Photon takes advantage of modern hardware for faster queries, providing up to 12x better price/performance compared to other cloud data warehouses — all natively on your data lake.
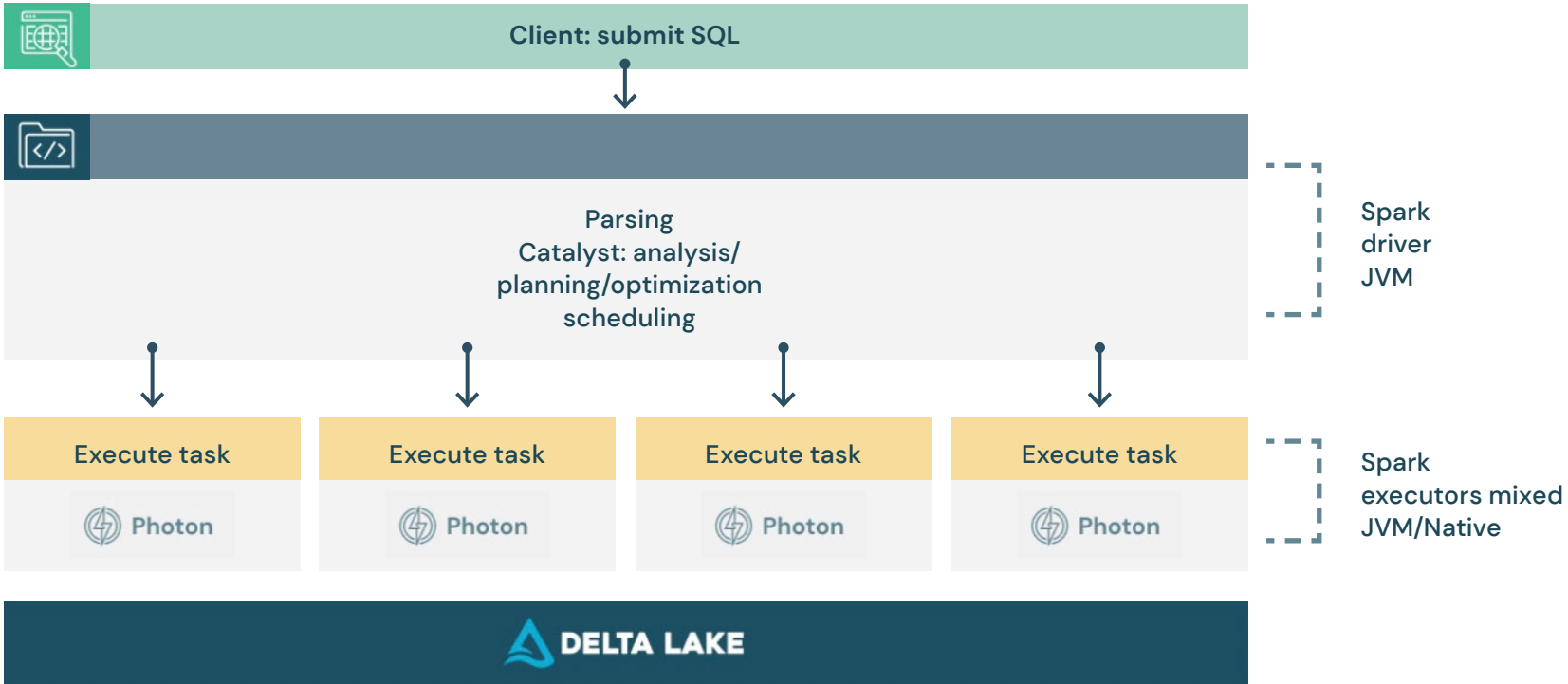
### 100TB TPC-DS price/performance
Lower is better



**databricks**

## Works with your existing code and avoids vendor lock–in

Photon is designed to be compatible with the Apache Spark DataFrame and SQL APIs to ensure workloads run seamlessly without code changes. All you do is turn it on. Photon will seamlessly coordinate work and resources and transparently accelerate portions of your SQL and Spark queries. No tuning or user intervention required.
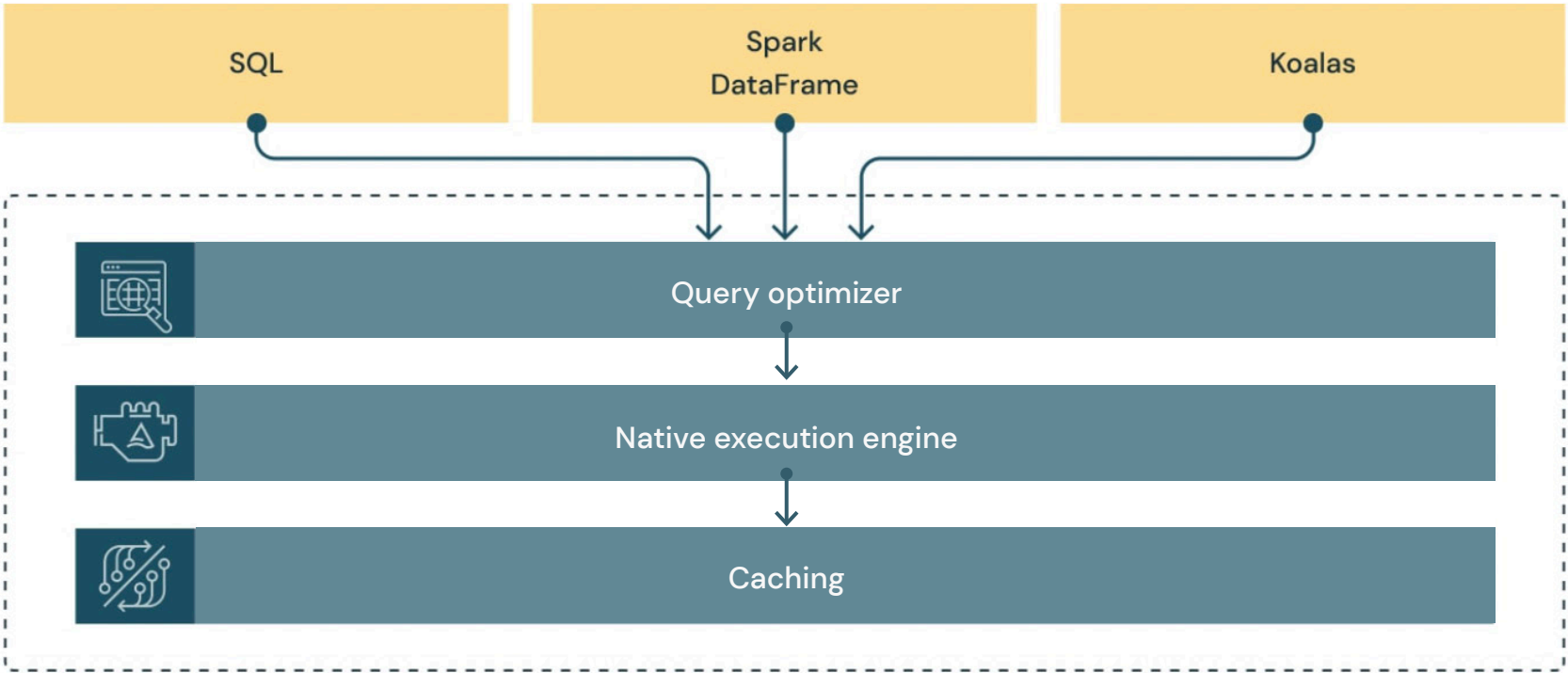
## Photon in the Databricks Data Intelligence Platform



*Lifecycle of a Photon query*

## Optimizing for all data use cases and workloads

Photon is the first purpose-built lakehouse engine designed to accelerate all data and analytics workloads: data ingestion, ETL, streaming, data science, and interactive queries. While we started Photon primarily focused on SQL to provide customers with world-class data warehousing performance on their data lakes, we've significantly increased the scope of ingestion sources, formats, APIs and methods supported by Photon since then. As a result, customers have seen dramatic infrastructure cost savings and speedups on Photon across all their modern Spark (e.g., Spark SQL and DataFrame) workloads.



*Accelerating all workloads on the lakehouse*

**Learn more**

Announcing Photon Public Preview: The Next-Generation Query Engine on the Databricks Data Intelligence Platform

Databricks Sets Official Data Warehousing Performance Record

# 04

## Unified governance and sharing for data, analytics and AI

Today, more and more organizations recognize the importance of making high-quality data readily available to data teams across business units, or externally with customers and partners to drive actionable insights and business value. At the same time, organizations also understand the risks of data breaches which negatively impact brand value and inevitably lead to erosion of customer trust. Governance is one of the most critical components of a lakehouse architecture; it helps ensure that data assets are securely managed throughout the enterprise. However, many companies are using different incompatible governance models leading to complex and expensive solutions.

databricks

# Key challenges with data and AI governance

## Diversity of data and AI assets

The increased use of data and the added complexity of the data landscape have left organizations with a difficult time managing and governing all types of their data-related assets. No longer is data stored in files or tables. Data assets today take many forms, including dashboards, machine learning models and unstructured data like video and images that legacy data governance solutions simply are not built to govern and manage.

## Two disparate and incompatible data platforms

Organizations today use two different platforms for their data analytics and AI efforts — data warehouses for BI and data lakes for AI. This results in data replication across two platforms, presenting a major governance challenge. With no unified view of the data landscape, it is difficult to see where data is stored, who has access to what data, and consistently define and enforce data access policies across the two platforms with different governance models.

## Rising multicloud adoption

More and more organizations now leverage a multicloud strategy to optimize costs, avoid vendor lock-in, and meet compliance and privacy regulations. With nonstandard, cloud-specific governance models, data governance across clouds is complex and requires familiarity with cloud-specific security and governance concepts, such as identity and access management (IAM).

## Disjointed tools for data governance on the lakehouse

Today, data teams must deal with a myriad of fragmented tools and services for their data governance requirements, such as data discovery, cataloging, auditing, sharing, access controls, etc. This inevitably leads to operational inefficiencies and poor performance due to multiple integration points and network latency between the services.

databricks

# One security and governance approach

Lakehouse systems provide a uniform way to manage access control, data quality and compliance across all of an organization's data using standard interfaces similar to those in data warehouses by adding a management interface on top of data lake storage.

Modern lakehouse systems support fine-grained (row, column and view level) access control via SQL, query auditing, attribute-based access control, data versioning and data quality constraints and monitoring. These features are generally provided using standard interfaces familiar to database administrators (for example, SQL GRANT commands) to allow existing personnel to manage all the data in an organization in a uniform way. Centralizing all the data in a lakehouse system with a single management interface also reduces the administrative burden and potential for error that comes with managing multiple separate systems.

# What is Unity Catalog?

Databricks Unity Catalog offers a unified governance layer for data and AI within the Databricks Data Intelligence Platform. Unity Catalog simplifies governance by empowering data teams with a common governance model based on ANSI-SQL to define and enforce fine-grained access controls. With attribute-based access controls, data administrators can enable fine-grained access controls on rows and columns using tags (attributes). Built-in data search and discovery allows data teams to quickly find and reference relevant data for any use case. Unity Catalog offers automated data lineage for all workloads in SQL, R, Scala and Python, to build a better understanding of the data and its flow in the lakehouse. Unity Catalog also allows data sharing across or within organizations and seamless integrations with your existing data governance tools.

With Unity Catalog, data teams can simplify governance for all data and AI assets with one consistent model to discover, access and share data, giving you much better native performance, management and security across clouds.

databricks

## Key benefits

### Catalog, secure and audit access to all data assets on any cloud

Unity Catalog provides centralized metadata, enabling data teams to create a single source of truth for all data assets ranging from files, tables, dashboards to machine learning models in one place.
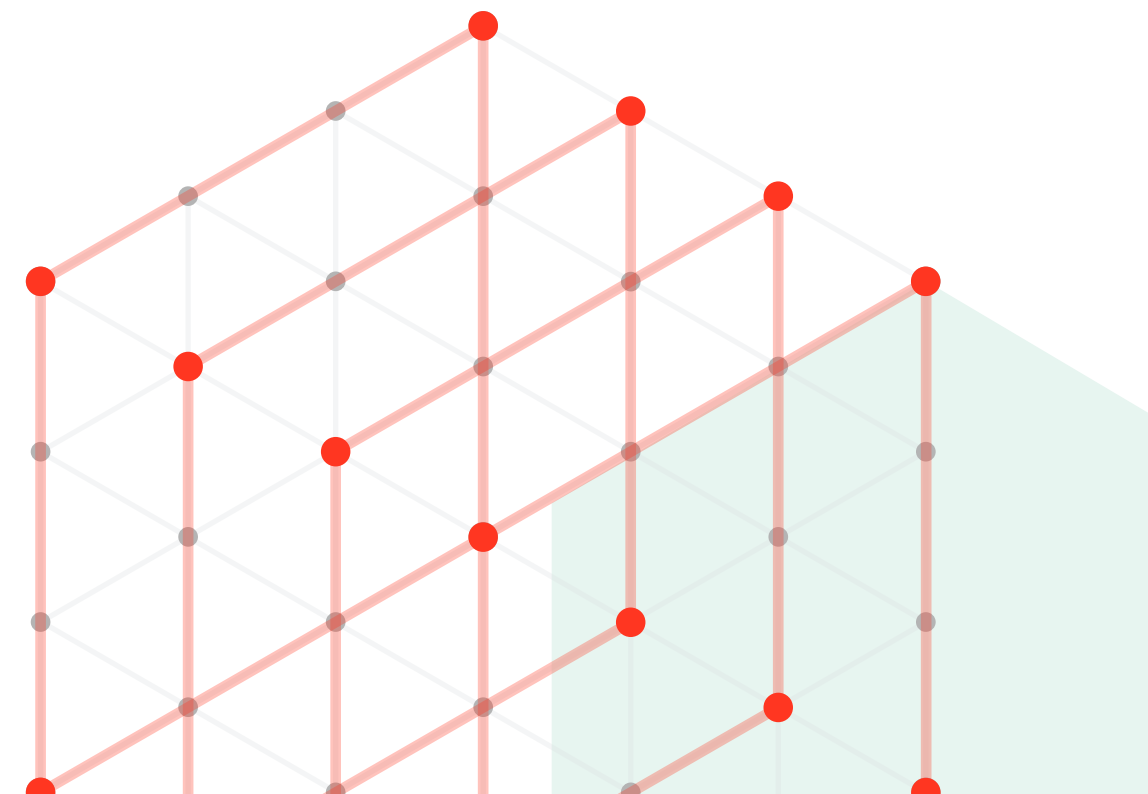
The common metadata layer for cross-workspace metadata is at the account level and eases collaboration by allowing different workspaces to access Unity Catalog metadata through a common interface and break down data silos. Further, the data permissions in Unity Catalog are applied to account-level identities, rather than identities that are local to a workspace, allowing a consistent view of users and groups across all workspaces.

**Without Unity Catalog**

Databricks workspace 1
- Access controls
- User management
- Metastore
- Clusters SQL endpoints

Databricks workspace 2
- Access controls
- User management
- Metastore
- Clusters SQL endpoints

**With Unity Catalog**

Unity Catalog
- Access controls
- User management
- Metastore

Databricks workspace 1
- Clusters SQL endpoints

Databricks workspace 2
- Clusters SQL endpoints

databricks

Unity Catalog offers a unified data access layer that provides a simple and streamlined way to define and connect to your data through managed tables, external tables, or files, while managing their access controls. Unity Catalog centralizes access controls for files, tables and views.
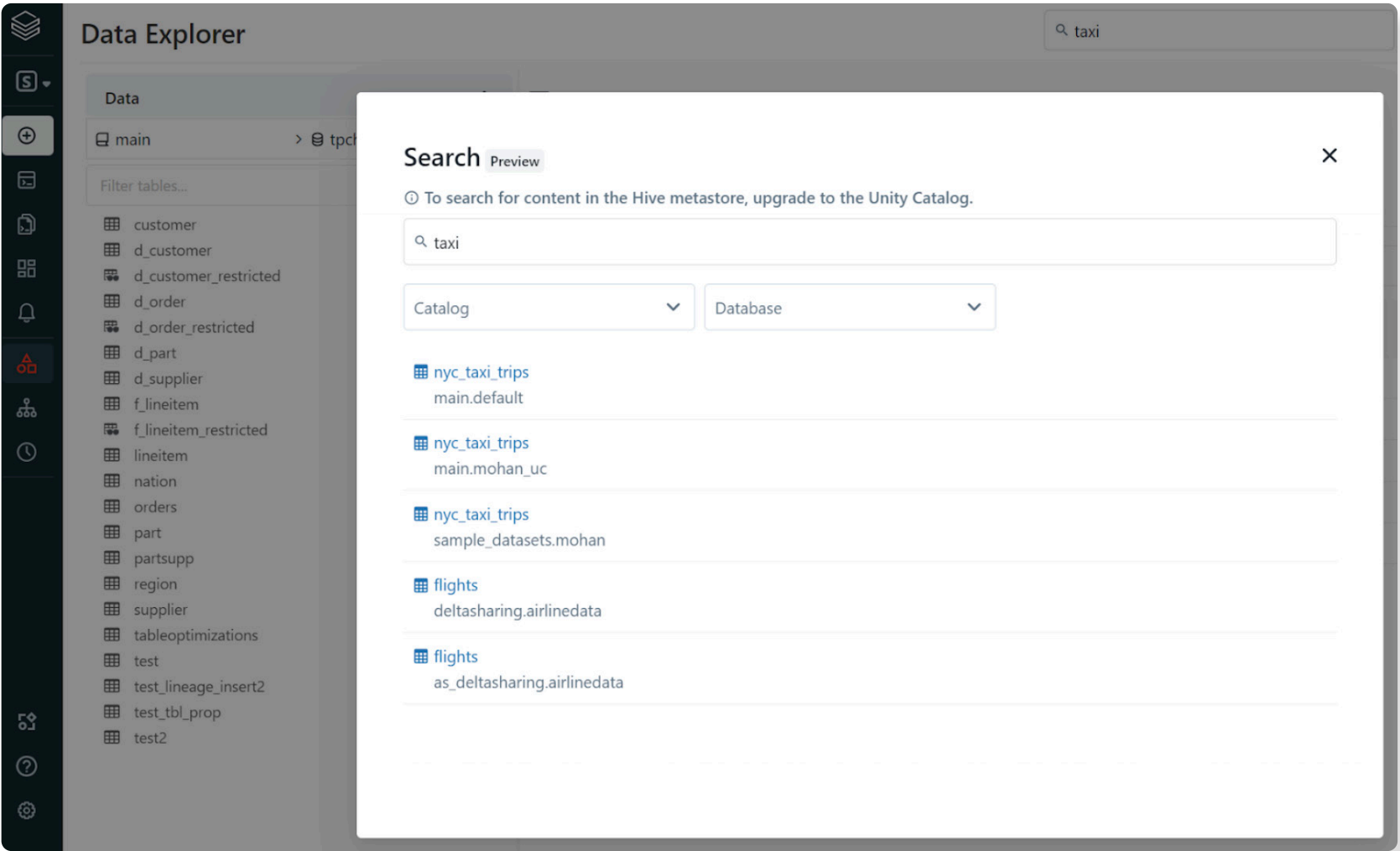
It allows fine-grained access controls for restricting access to certain rows and columns to the users and groups who are authorized to query them. With Attribute-Based Access Controls (ABAC), you can control access to multiple data items at once based on user and data attributes, further simplifying governance at scale. For example, you will be able to tag multiple columns as personally identifiable information (PII) and manage access to all columns tagged as PII in a single rule.

Today, organizations are dealing with an increased burden of regulatory compliance, and data access auditing is a critical component to ensure your organization is set up for success while meeting compliance requirements. Unity Catalog also provides centralized fine-grained auditing by capturing an audit log of operations such as create, read, update and delete (CRUD) that have been performed against the data. This allows a fine-grained audit trail showing who accessed a given data set and helps you meet your compliance and business requirements.

databricks

## Built-in data search and discovery

Data discovery is a critical component to break down data silos and democratize data across your organization to make data-driven decisions. Unity Catalog provides a rich user interface for data search and discovery, enabling data teams to quickly search relevant data assets across the data landscape and reference them for all use cases — BI, analytics and machine learning — accelerating time-to-value and boosting productivity.

## Automated data lineage for all workloads

Data lineage describes the transformations and refinements of data from source to insight. Lineage includes capturing all the relevant metadata and events associated with the data in its lifecycle, including the source of the data set, what other data sets were used to create it, who created it and when, what transformations were performed, which other data sets leverage it, and many other events and attributes. Unity Catalog offers automated data lineage down to table and column level, enabling data teams to get an end-to-end view of where data is coming from, what transformations were performed on the data and how data is consumed by end applications such as notebooks, workflows, dashboards, machine learning models, etc.

With automated data lineage for all workloads — SQL, R, Python and Scala, data teams can quickly identify and perform root cause analysis of any errors in the data pipelines or end applications. Second, data teams can perform impact analysis to see dependencies of any data changes on downstream consumers and notify them about the potential impact. Finally, data lineage also empowers data teams with increased understanding of their data and reduces tribal knowledge. Unity Catalog can also capture lineage associated with non-data entities, such as notebooks, workflows and dashboards. Lineage can be



*Data lineage with Unity Catalog*

retrieved via REST APIs to support integrations with other catalogs.

## Integrated with your existing tools

Unity Catalog helps you to future-proof your data and AI governance with the flexibility to leverage your existing data catalogs and governance solutions — Collibra, Alation, Immuta, Privacera, Microsoft Purview and AWS Lakeformation.

## Resources

Learn more about Unity Catalog

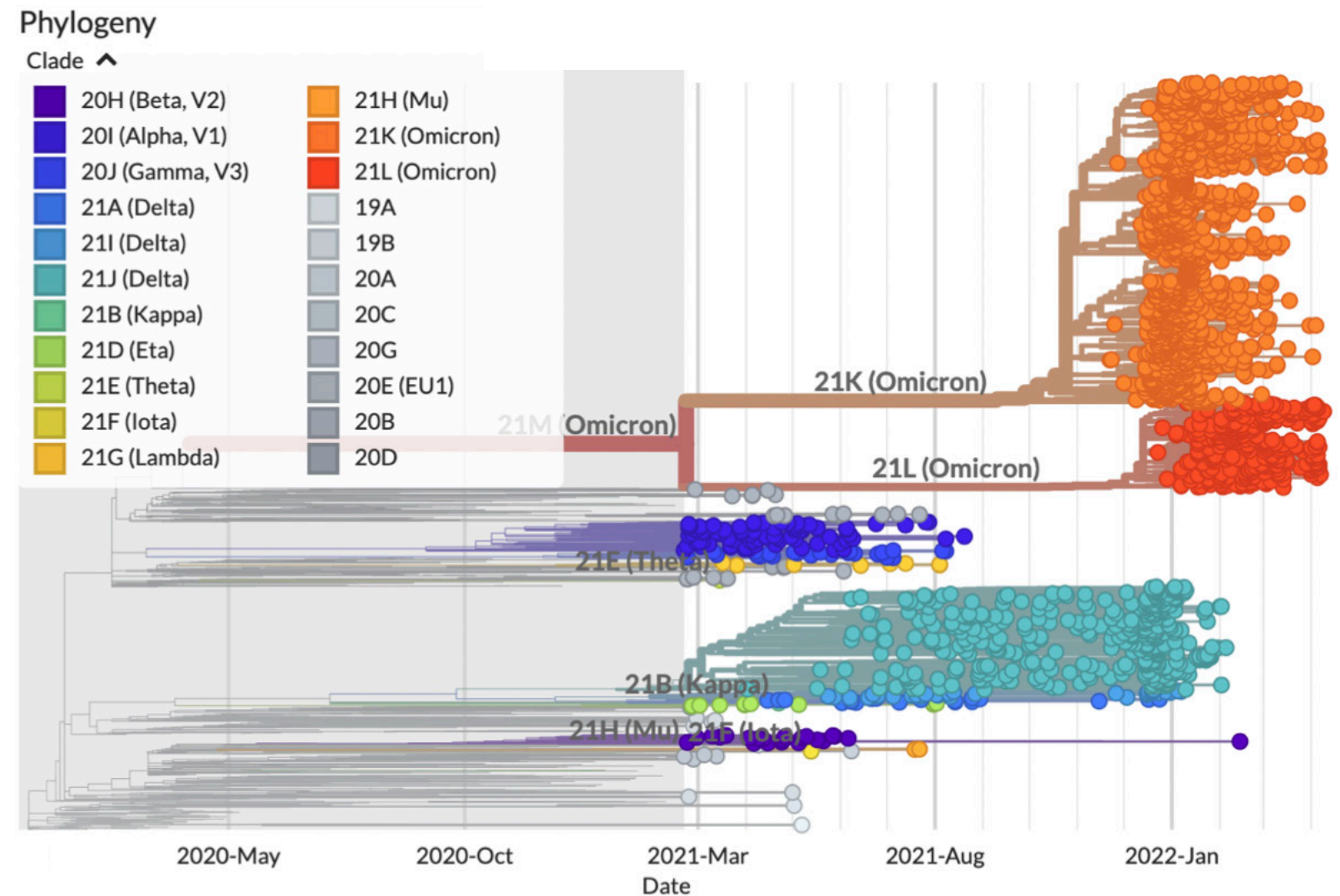AWS Documentation

Azure Documentation

# Open data sharing and collaboration

Data sharing has become important in the digital economy as enterprises wish to exchange data easily and securely with their customers, partners, suppliers and internal lines of business to better collaborate and unlock value from that data. But to date, a lack of standards-based data sharing protocol has resulted in data sharing solutions tied to a single vendor or commercial product, introducing vendor lock-in risks. What the industry deserves is an open approach to data sharing.

## Why data sharing is hard

Data sharing has evolved from an optional feature of a few data platforms to a business necessity and success factor for organizations. Our solution architects encounter daily the classic scenarios of a retailer looking to publish sales data to their suppliers in real time or a supplier that wants to share real-time inventory.

As a reminder, data sharing recently triggered the most impressive scientific development that humankind has ever seen. On January 5, 2021, the first sample of the genome of the coronavirus was uploaded to the internet. It wasn't a lung biopsy



from a patient in Wuhan, but a shared digital genomic data set that triggered the development of the first batch of COVID vaccines worldwide.

Since then, coronavirus experts have daily exchanged public data sets, looking for better treatments, tests and tracking mutations as they

are passed down through a lineage, a branch of the coronavirus family tree. The above graphic shows such a publicly shared mutation data set.

Sharing data, as well as consuming data from external sources, allows you to collaborate with partners, establish new partnerships, enable research and can generate new revenue streams with data monetization.

Despite those promising examples, existing data sharing technologies come with several limitations:

- Traditional data sharing technologies, such as Secure File Transfer Protocol (SFTP), do not scale well and only serve files offloaded to a server

- Cloud object stores operate on an object level and are cloud–specific

- Commercial data sharing offerings baked into vendor products often share tables instead of files, but scaling them is expensive and they are not open and, therefore, do not permit data sharing with a different platform

The following table compares proprietary vendor solutions with SFTP, cloud object stores and Delta Sharing.

| | Proprietary vendor solutions | SFTP | Cloud object store | Delta Sharing |
|---|---|---|---|---|
| Secure | ✓ | ✓ | ✓ | ✓ |
| Cheap | | ✓ | ✓ | ✓ |
| Vendor agnostic | | ✓ | | ✓ |
| Multicloud | | ✓ | | ✓ |
| Open source | | ✓ | | ✓ |
| Table/DataFrame abstraction | ✓ | | | ✓ |
| Live data | ✓ | | | ✓ |
| Predicate pushdown | ✓ | | | ✓ |
| Object store bandwidth | | | ✓ | ✓ |
| Zero compute cost | | | ✓ | ✓ |
| Scalability | | | ✓ | ✓ |

databricks

# Open source data sharing and Databricks

To address the limitations of existing data sharing solutions, Databricks developed Delta Sharing, with various contributions from the OSS community, and donated it to the Linux Foundation.

An open source–based solution, such as Delta Sharing, eliminates the lock–in of commercial solutions and brings a number of additional benefits such as community–developed integrations with popular, open source data processing frameworks. In addition, open protocols allow the easy integration of commercial clients, such as BI tools.

## What is Databricks Delta Sharing?

Databricks Delta Sharing provides an open solution to securely share live data from your lakehouse to any computing platform. Recipients don't have to be on the Databricks Platform or on the same cloud or a cloud at all. Data providers can share live data, without replicating or moving it to another system. Recipients benefit from always having access to the latest version of data and can quickly query shared data using tools of their choice for BI, analytics and machine learning, reducing time–to–value. Data providers can centrally manage, govern, audit and track usage of the shared data on one platform.

Unity Catalog natively supports Delta Sharing, the world's first open protocol for data sharing, enabling organizations to share live, large-scale data without replication and make data easily and quickly accessible from tools of your choice, with enterprise–grade security.

# Key benefits

## Open cross–platform sharing

Easily share existing data in Delta Lake and Apache Parquet formats between different vendors. Consumers don't have to be on the Databricks Platform, same cloud or a cloud at all. Native integration with Power BI, Tableau, Spark, pandas and Java allow recipients to consume shared data directly from the tools of their choice. Delta Sharing eliminates the need to set up a new ingestion process to consume data. Data recipients can directly access the fresh data and query it using tools of their choice. Recipients can also enrich data with data sets from popular data providers.

## Sharing live data without copying it

Share live ready–to–query data, without replicating or moving it to another system. Most enterprise data today is stored in cloud data lakes. Any of the existing data sets on the provider's data lake can easily be shared across clouds, regions or data platforms without any data replication or physical movement of data. Data providers can update their data sets reliably in real time and provide a fresh and consistent view of their data to recipients.
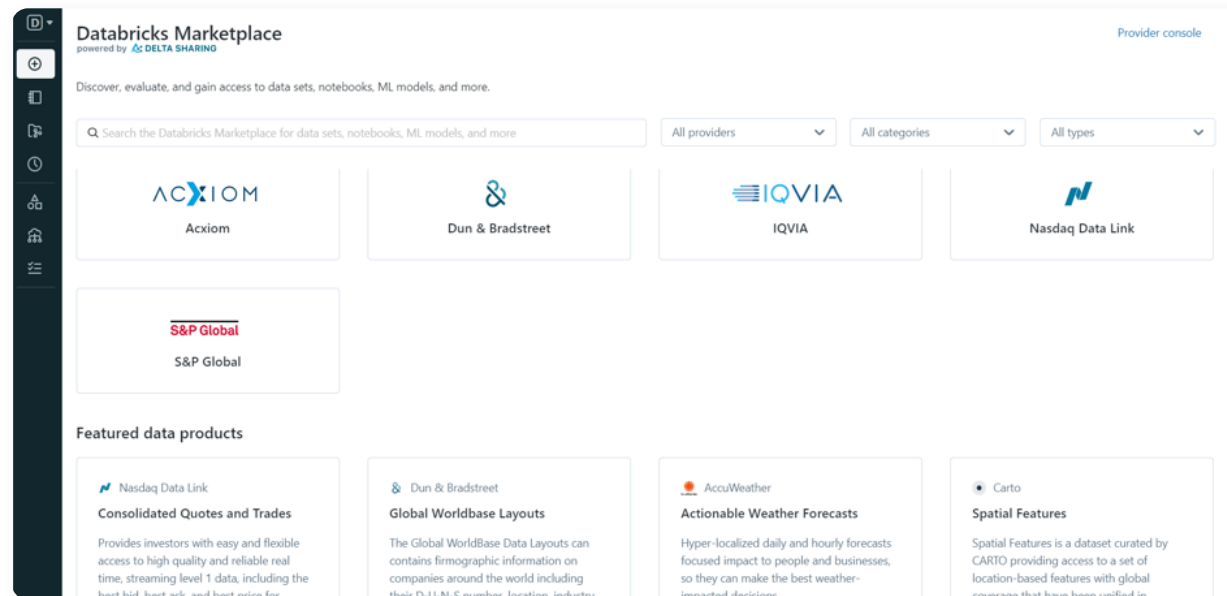
## Centralized administration and governance

You can centrally govern, track and audit access to the shared data from a single point of enforcement to meet compliance requirements. Detailed user–access audit logs are kept to know who is accessing the data and monitor usage of the shared data down to table, partition and version level.

databricks

## An open Marketplace for data solutions

The demand for third-party data to make data-driven innovations is greater than ever, and data marketplaces act as a bridge between data providers and data consumers to help facilitate the discovery and distribution of data sets.

Databricks Marketplace provides an open marketplace for exchanging data products such as data sets, notebooks, dashboards and machine learning models. To accelerate insights, data consumers can discover, evaluate and access more data products from third-party vendors than ever before. Providers can now commercialize new offerings and shorten sales cycles by providing value-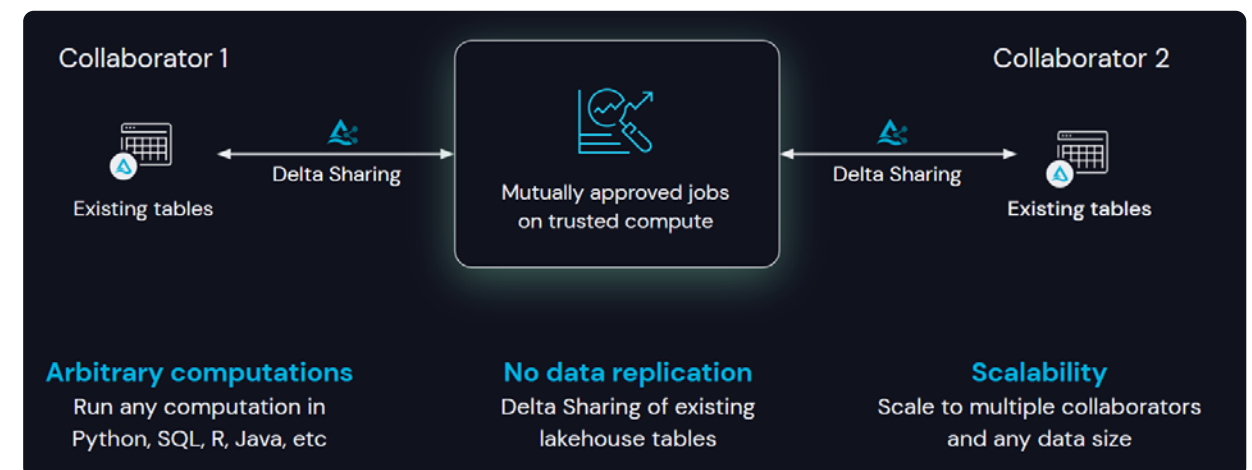added services on top of their data. Databricks Marketplace is powered by Delta Sharing, allowing consumers to access data products without having to be on the Databricks Platform. This open approach allows data providers to broaden their addressable market without forcing consumers into vendor lock-in.



*Databricks Marketplace*

## Privacy-safe data cleanrooms

Powered by open source Delta Sharing, the Databricks Data Intelligence Platform provides a flexible data cleanroom solution allowing businesses to easily collaborate with their customers and partners on any cloud in a privacy-safe way. Participants in the data cleanrooms can share and join their existing data, and run complex workloads in any language — Python, R, SQL, Java and Scala — on the data while maintaining data privacy. Additionally, data cleanroom participants don't have to do cost-intensive data replication across clouds or regions with other participants, which simplifies data operations and reduces cost.
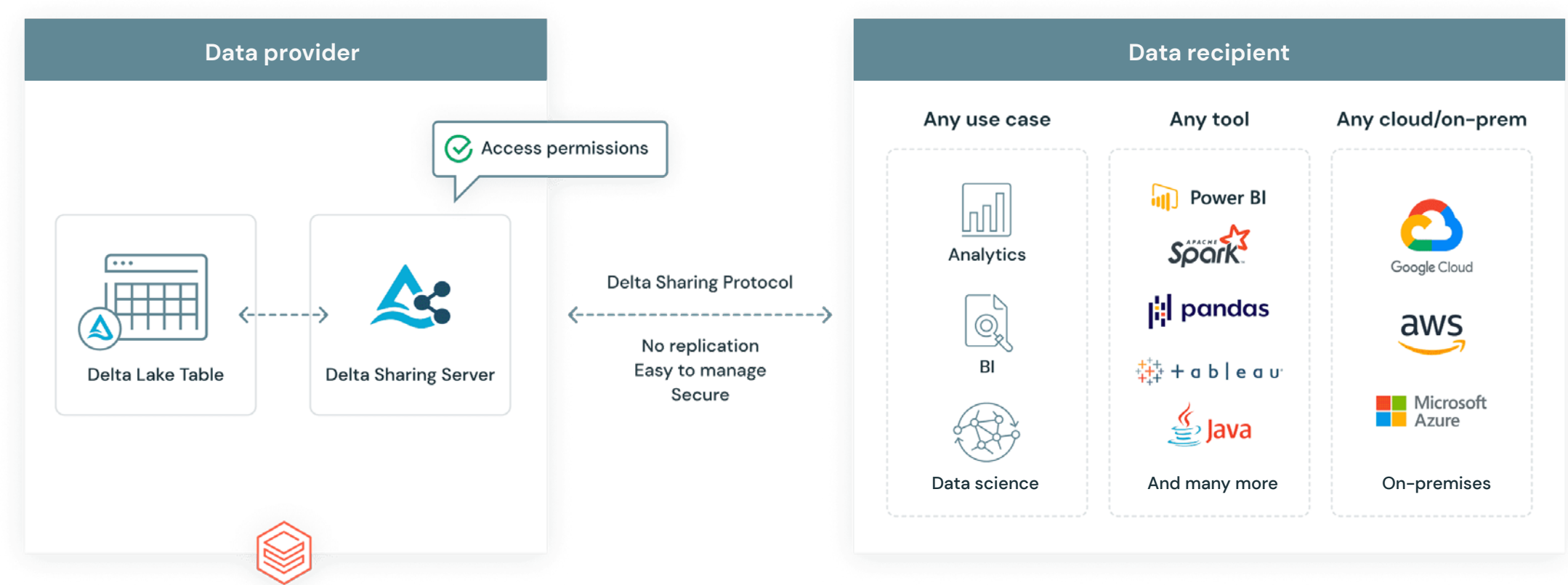


*Data cleanrooms with the Databricks Data Intelligence Platform*

databricks

# How it works

Delta Sharing is designed to be simple, scalable, non-proprietary and cost-effective for organizations that are serious about getting more from their data. Delta Sharing is natively integrated with Unity Catalog, which allows customers to add fine-grained governance and security controls, making it easy and safe to share data internally or externally.

Delta Sharing is a simple REST protocol that securely shares access to part of a cloud data set. It leverages modern cloud storage systems — such as AWS S3, Azure ADLS or Google's GCS — to reliably transfer large data sets. Here's how it works for data providers and data recipients.



The data provider shares existing tables or parts thereof (such as specific table versions or partitions) stored on the cloud data lake in Delta Lake format. The provider decides what data they want to share and runs a sharing server in front of it that implements the Delta Sharing protocol and manages access for recipients. To manage shares and recipients, you can use SQL commands or the Unity Catalog CLI or the intuitive user interface.

The data recipient only needs one of the many Delta Sharing clients that supports the protocol. Databricks has released open source connectors for pandas, Apache Spark, Java and Python, and is working with partners on many more.

databricks