# Deloitte.

databricks

# The power of an innovative data and analytics platform

Intelligence and insights at scale

MAKING AN IMPACT THAT MATTERS
*since 1845*

# Making the complex simple

## Unlocking the power of data-driven insights

Unprecedented levels of data and an increasing need for businesses to quickly identify meaningful trends have propelled the demand for **more effective, flexible, and amplified data and analytics platform (DAP)** capabilities. Today's DAP requires unified and cloud platforms capable of handling massive data growth, advanced analytics, enterprise-scale machine learning, and artificial intelligence workloads as well as the ability to oversee entire data oceans.

To accelerate the digital transformation journey and realize tangible business outcomes, platforms must provide enterprises with an accessible data product to "**do more with less**," operating seamlessly across their organization's ecosystem of services and technologies, incorporating several tools with various capabilities for actionable insights.

The objective of this white paper is to impartially explain how Databricks capabilities align to a modern data and analytics platform (DAP).

# Key Tenets

## Harnessing insights to drive impact

At the core of data and analytics platforms is the capacity to transform raw data into actionable insights. **The platform must enable secure, scalable, and cost-effective data analysis in order to maximize the potential of your business**. By having a well-architected DAP, you can unlock real-time insights to make informed decisions and drive better results. To unleash your business' potential with a data and analytics platform, it is important to carefully address the key tenets of data and analytics.

In this whitepaper, we'll break down how Databricks aligns to the core principles and capabilities of DAP.

### SaaS / PaaS platform
On-demand access to a complete, ready-to-use, cloud-hosted platform/ cloud-hosted software at scale
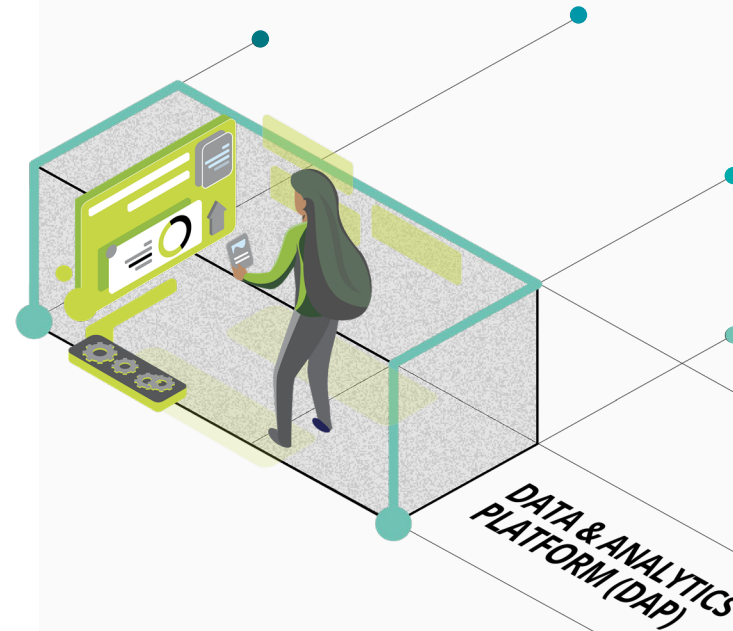
### Data engineering
Designing and building systems for collecting, transforming, and storing data for analysis

### Data warehouse & advanced analytics
Techniques, tools, and technologies for centralized repositories of traditional business intelligence (BI) and beyond to gain deeper insights, make predictions, or make data-driven decisions

### Security & governance
Overseeing the data and analytics platform to mitigate business risks

**DATA & ANALYTICS PLATFORM (DAP)**

# What About Generative AI?

## Setting the stage for technology advancements

Generative AI is the focus of many enterprise technology conversations right now. AI-driven text generation models can simulate reality and generate novel, previously unseen data. This can be used to enhance AI-driven decisions or chatbots, or forecast outcomes such as customer sentiment analysis, and to improve the performance of applications through the generation of more powerful and efficient programming code. With these capabilities, it is possible to achieve outcomes beyond those achieved with existing technologies.

**But in order to maximize the full power of Generative AI and unlock its potential, it is imperative to establish an unparalleled data platform that serves as the catalyst for accelerated, cost-effective, and superior advancements in all facets of technology and innovation.**

Databricks offers a unified platform for the storage of LLM-based vector search, feature stores, and the life cycle of AI and LLM model development. This encompasses model management, training, evaluation, serving, and monitoring, whether utilizing existing LLMs or constructing new ones. The platform also incorporates several innovative Gen AI features, including but not limited to Lakehouse IQ, Databricks Assistant, and Dolly 2.0, among others.

# Generative AI in action

Generative AI features on The Data Intelligence  Platform
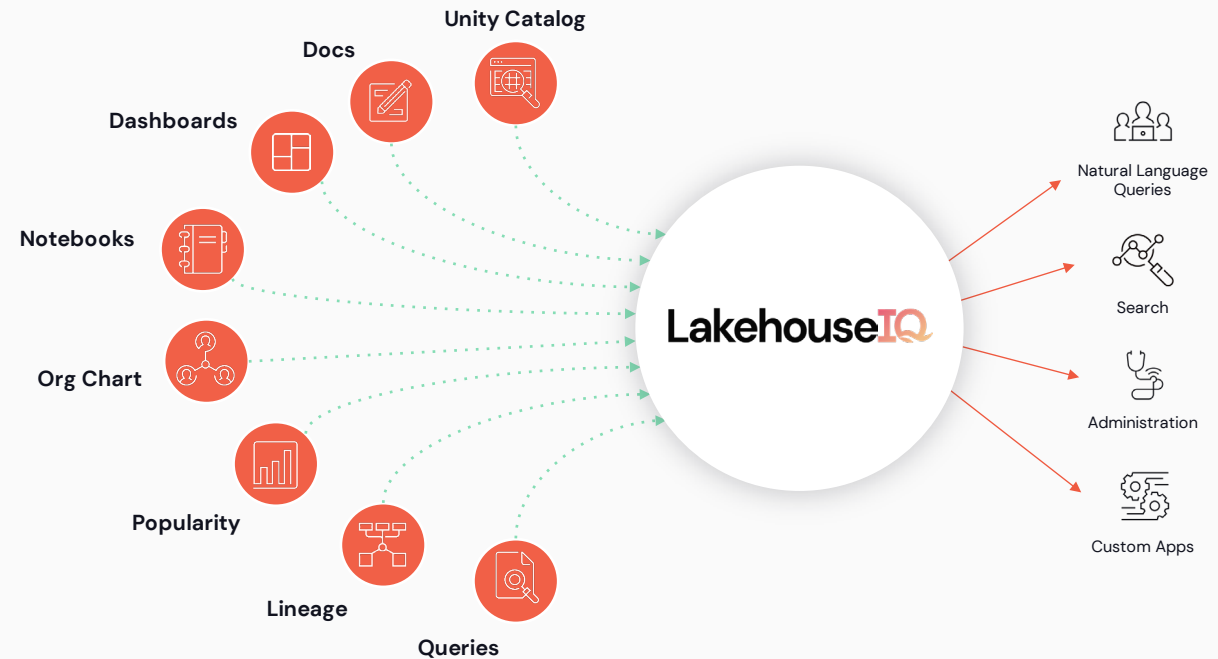
## Dolly 2.0

An open-source community contribution that demonstrates how organizations can create, own, and customize their own Language Models (LLMs) to communicate with people, without having to pay for API access or share their data with third-parties. Although Databricks does not expect Dolly to be the most cutting-edge in terms of performance yet, its open-source dataset is meant to act as the foundation for further development.

## Lakehouse IQ (Preview)

Utilizes a range of corporate assets, encompassing schemas, popularity metrics, lineage data, descriptions, organizational charts, and data usage insights, to tackle user queries. This includes generating meaningful data by comprehending enterprise assets and terminology, searching for organization-specific information, and intelligently crafting code through NLP-driven interactions.

## Databricks Assistant

Facilitates the conversion of English natural language queries into SQL queries, offering code explanations and error corrections to enhance the quality and productivity of the development community..
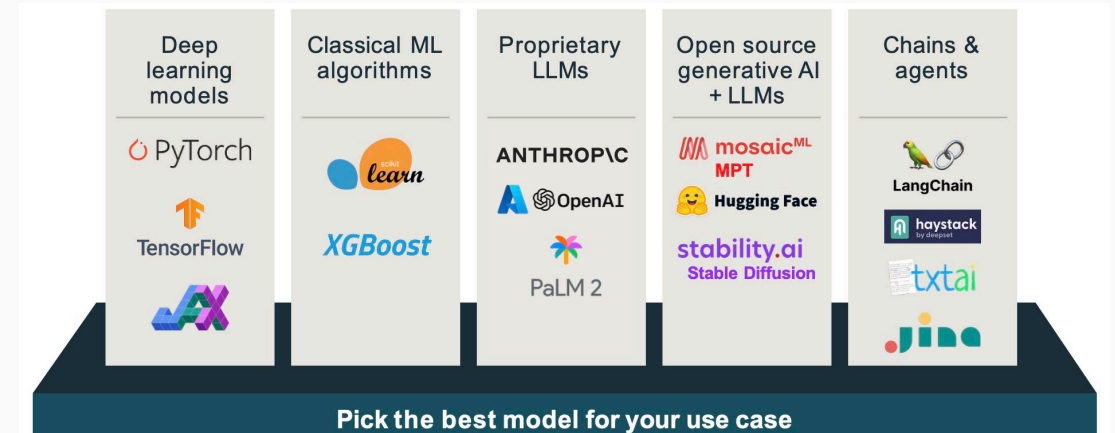
# Generative AI in action

## Lakehouse AI works for both proprietary and open-source AI models

Using Lakehouse AI, users are afforded the flexibility to access a diverse array of proprietary models, including but not limited to offerings like Anthrop\C, OpenAI and Palm 2.

Additionally, Lakehouse AI empowers users to harness the capabilities of open-source models, such as those provided by Hugging Face, Llama 2, Mosaic ML MPT, Stability.ai, and more.

This versatility allows for a range of AI model choices and ensures users can select the most suitable options for their specific needs.



| Deep learning models | Classical ML algorithms | Proprietary LLMs | Open source generative AI + LLMs | Chains & agents |
|---|---|---|---|---|
| PyTorch TensorFlow JAX | learn XGBoost | ANTHROP\C OpenAI PaLM 2 | mosaicML MPT Hugging Face stability.ai Stable Diffusion | LangChain haystack by deepset txtai jina |

**Pick the best model for your use case**

# SaaS / PaaS platform

Choose a flexible and reliable platform for your business needs

Software as a Service (SaaS) & Platform as a Service (PaaS)

Enterprises are transitioning to platforms that **reduce or eliminate the need for procuring, installing, upgrading, patching, safeguarding, backing up, restoring, and managing their environment, with minimal upfront costs**. These platforms are designed to provide enterprises with the latest features and functionalities with no hardware constraints, as well as the ability to run multiple versions of the same software.

**A Databricks workspace is a unified and collaborative environment that provides access to various Databricks assets, such as a web application, customer notebooks, jobs and queues, and cluster management. The Databricks Data Intelligence Platform as a Service includes Databricks Runtime, which provides optimized Apache Spark™, Unity Catalog, Delta Lake, and MLflow. Databricks simplifies code and data management by allowing easy version control and dependency management.**
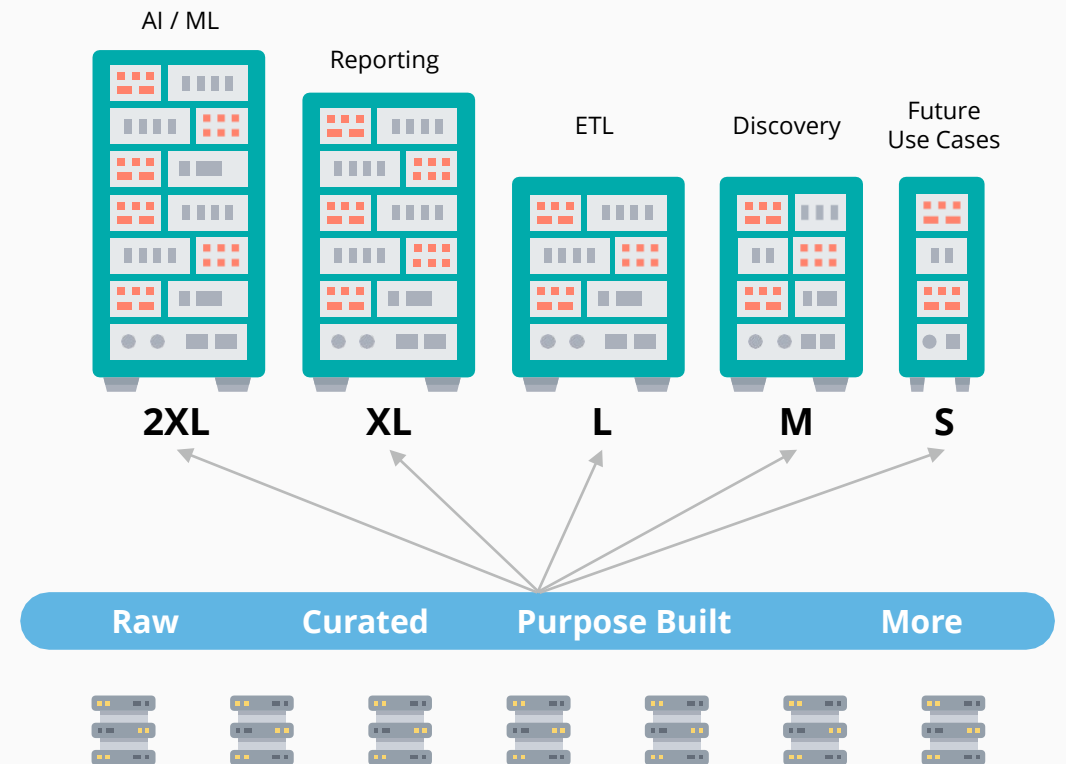
# SaaS / PaaS Platform

## Break free from dependence on compute and storage

The cost of ownership and sunken cost can be drastically reduced when systems are loosely coupled since storage and compute can be expanded independently and on demand. Cloud technologies **enabled through software provide virtually unlimited expansion of storage and compute power**. This allows companies to spend less on these resources, leaving more room for innovation and discovery. The software enables efficient functioning of the system and is optimized to work within and across a multi-cloud environment without sacrificing performance.

Databricks utilizes a decoupled storage and compute architecture to improve cost efficiency and granular scalability. This allows users to bring their own compute power of their desired type (memory optimized, compute optimized, or balanced instances) to scale with no code changes and without needing to make additional data copies or move data.

**SHARED STORAGE**

Consuming organization has choice of computer power



AI / ML

Reporting

ETL

Discovery

Future Use Cases

**2XL**　　**XL**　　**L**　　**M**　　**S**

| Raw | Curated | Purpose Built | More |
|-----|---------|---------------|------|

Compute with decoupled storage

# SaaS / PaaS Platform

## Maximize elasticity and scalability

DAP that **"automatically" scales (out or in/up or down)** resources to meet analytics demands, from small data sets to large, with no effort is critical for controlling cloud costs. It will also meet service-level agreements, even when data volume grows during seasonal demand, outages, or large volumes of inflow of data. While cloud provides the scalability/elasticity, configuring to achieve "economies of scale," it is an involved, manual effort. **Modern DAP SaaS software makes it seamless, with minimal or no manual effort, to reduce costs and gain efficiency** when operating at a larger or smaller scale.

### Serverless Compute
Instant compute available to process user queries within seconds reduces latency, improves scalability, reduces management overhead, accelerates innovation, and reduces the cost.

### Auto-Scaling
Databricks clusters can be quickly spun up when there is a need for processing large amounts of data, and then shut down when no longer required, allowing for efficient use of resources.
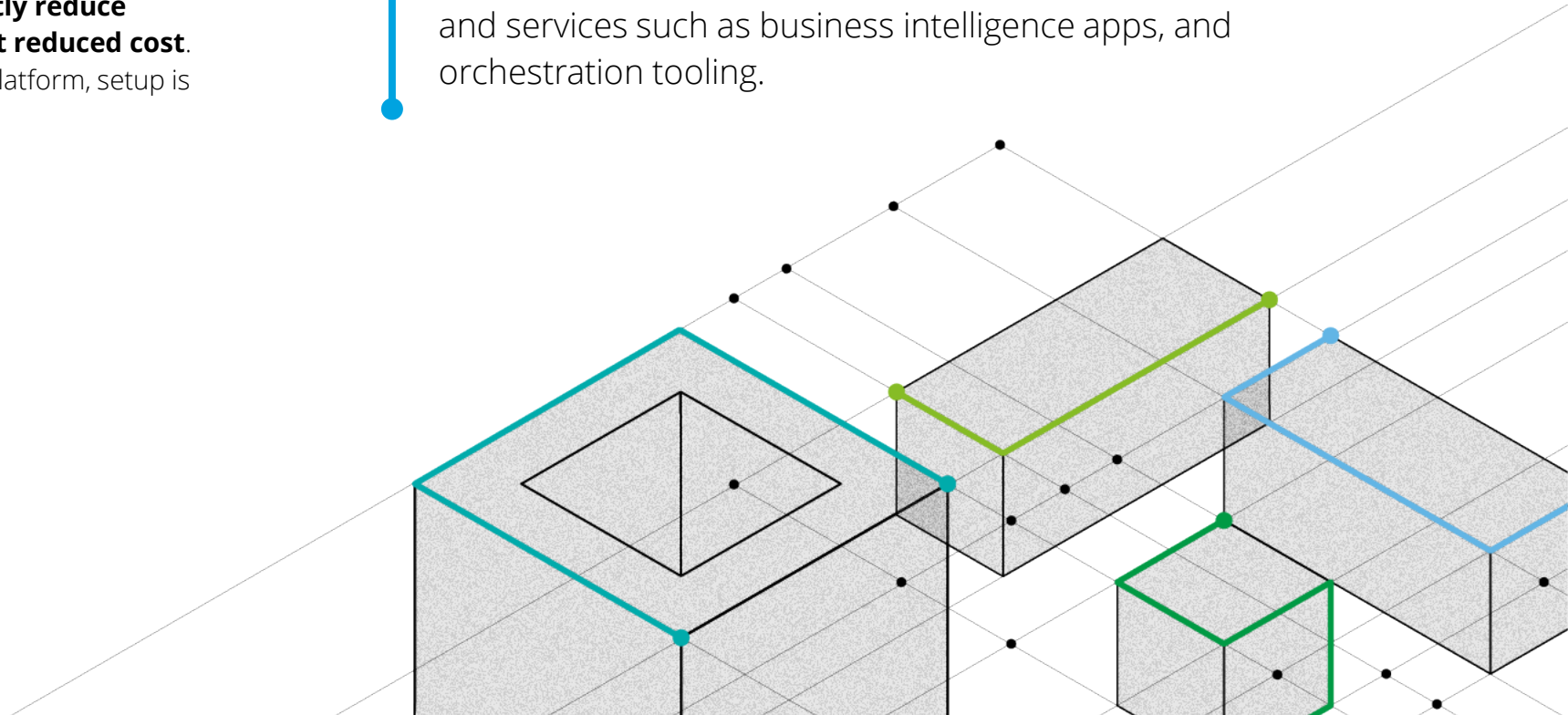
### Pools
Enables faster cluster start-up and scaling by providing a managed cache of pre-created virtual machine instances that can be quickly acquired when needed.

# SaaS / PaaS Platform

Be prepared and plan ahead for disaster recovery

Enabling cloud-based "**managed disaster recovery**" (DR) service(s) helps to quickly recover an organization's critical systems after a disaster and provides remote access to the systems in a secure environment. Leveraging native sync across DR environments/regions is critical to reducing downtime. SaaS- and PaaS-based DAP software services reduce complex Infrastructure as a Service-based cloud DR setup, provide **faster spin-ups and failovers of a particular service, significantly reduce downtime, and provide faster recovery times, all at reduced cost**. However, due to the number of services used in a DAP platform, setup is still complex.

Databricks is available for disaster recovery in various regions around the world on AWS, Azure, and GCP. This involves setting up upstream data ingestion services for both batch and streaming data, cloud native storage, libraries and other dependencies, downstream tools and services such as business intelligence apps, and orchestration tooling.
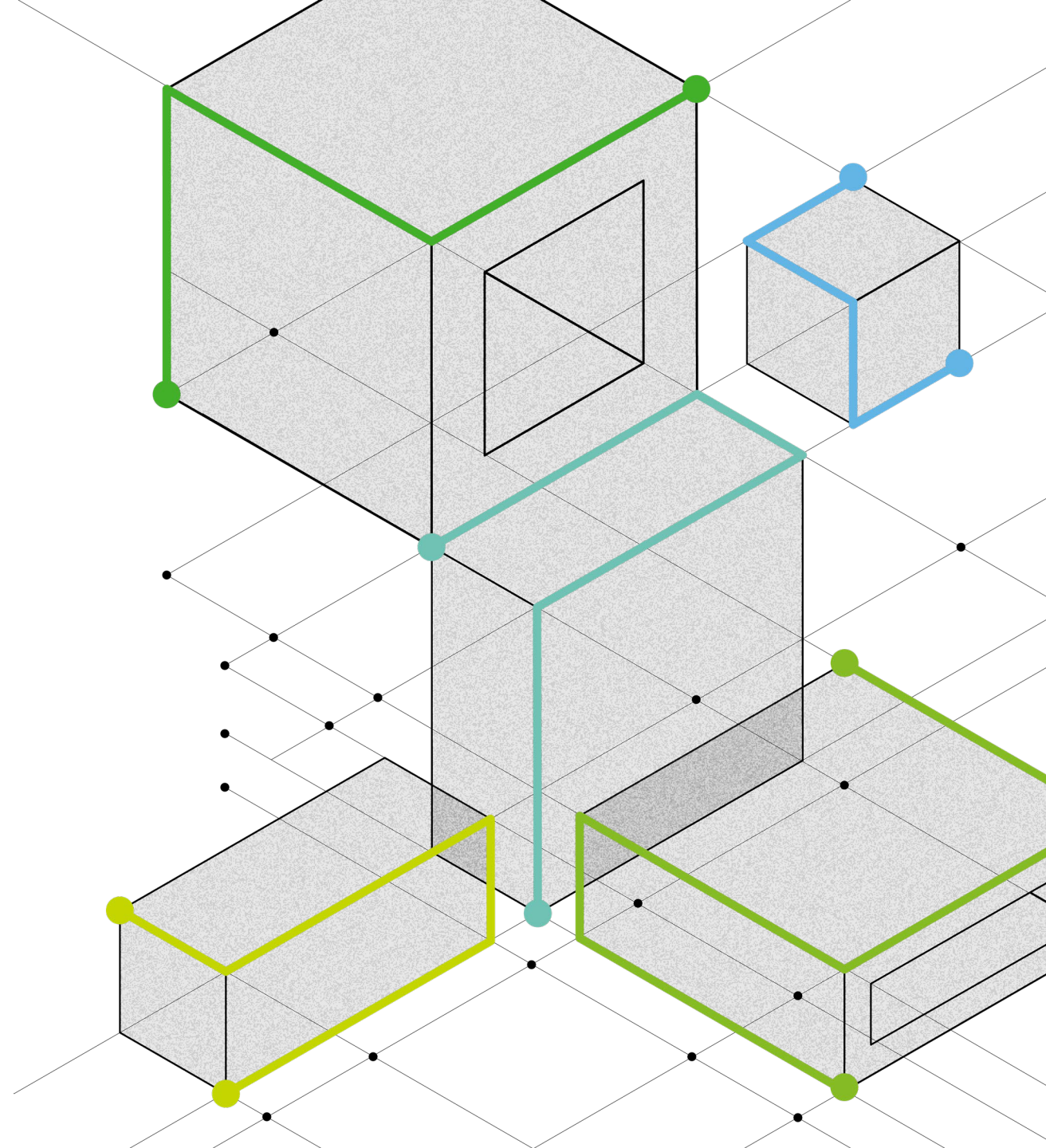
# SaaS / PaaS Platform

## Get the most value for your money

Despite the cloud's benefits, many organizations tend to make cloud product decisions based solely on the costs associated with performance benchmarking results. However, organizations that make product decisions based on the actual **cost of the overall services and whether they can deliver realize the best economic value and returns on their investment**.

The recent TPC-DS benchmark of Databricks' price performance demonstrated higher performance with lower investments. Databricks helps to reduce cost by optimizing the use of cloud resources. It provides an all-in-one offering of data engineering, SQL analytics, cataloging, data sharing, and AI/ML capabilities integrated with DevOps and MLOps. Databricks allows users to choose their preferred programming language without having to leave the platform.
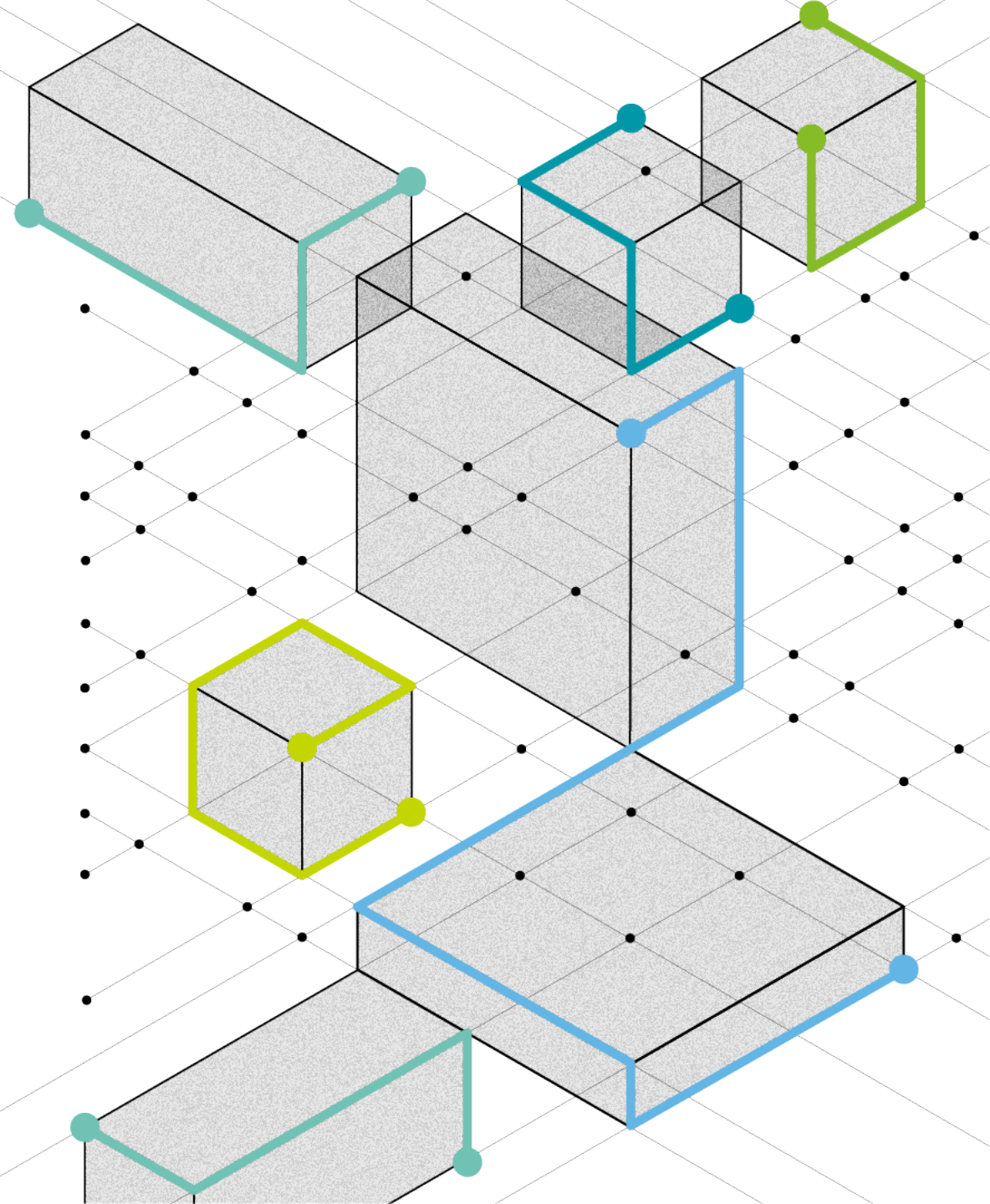
# SaaS / PaaS Platform

## Powering the vendor-supported, open-source revolution

Vendor-backed, open-source enables a more long-term approach to product development, making it easier to integrate with other services and **solve bugs and performance issues faster**. It also promotes healthy competition between vendors to continue to develop best-in-class products for consumers, without a vendor lock-in.

The Databricks foundation is based on open-source Apache Spark™, but is improved through the performance optimizations and functionality provided by Delta Lake, such as ZOrder, Change Data Feed, Dynamic Partition Overwrites, and Dropped Columns. This enables performance capabilities for all lakehouse workloads **across hyperscalers (AWS, Azure and GCP)**, from streaming to batch processing, compared to other storage layers. Databricks contributes to the open-source community, most recently with the Data Lake 2.0 series of innovations, which enable the ability to leverage the power of Spark advancements.

# Data Engineering

## Engineering data for enhanced analysis and decision-making

**Unify streaming and batch processing**

Streaming and batch processing use case demands are increasing and edge analytics are gaining momentum due to high-powered, modern edge hardware (sensor, phone, etc.) capable of collecting, analyzing, and creating actionable insights in real time, directly from the IoT devices generating the data. While batch processing is still utilized, more and more consumers are eager to get their data in real time and in combination with batch data to explore, analyze, and discover insights. **Maintaining separate platforms reduces operational efficiency, requires different skill sets and a lengthier process to fix failures, and increases costs**.

**Databricks provides data processing capabilities for streaming and batch data on AWS, Azure, and GCP. By using Delta Live Tables for ETL/ELT, data engineers can unify batch and streaming workloads, eliminating the need to design and code two different applications. Autoloader and Delta Live Tables (DLT) can make streaming and batching applications faster and more efficient. For more information, please see the <u>DLT whitepaper</u>. DLT offers Automated Infrastructure Management and a new optimization layer (Enzyme) designed specifically to speed up the ETL processing.**
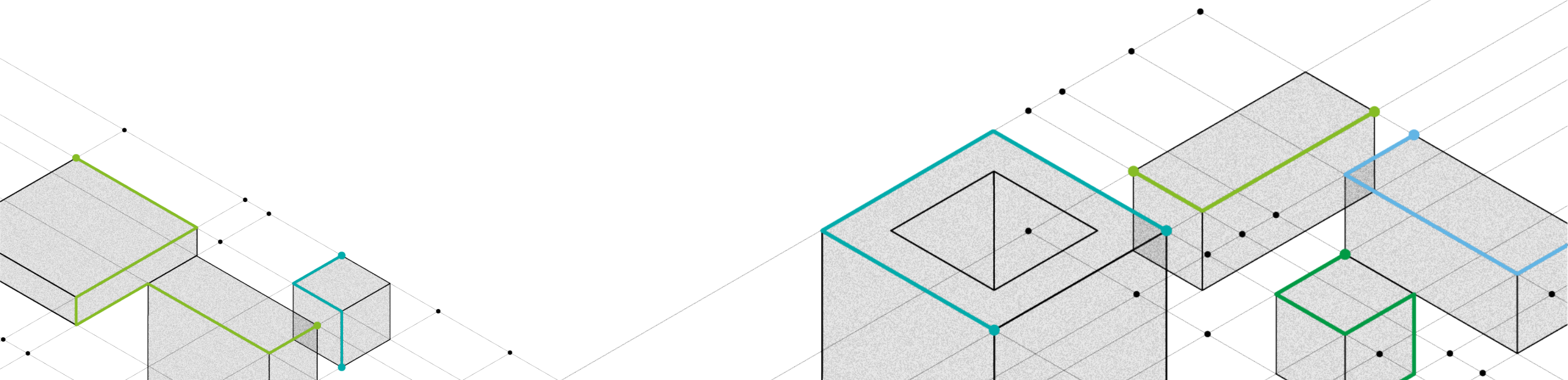
# Data Engineering

## Embrace automation with CI/CD & DevOps

Manual or semi-automated release management leads to lack of scalability and automation, slower dev cycles, and lower software quality. The integration of CI/CD and DevOps best practices throughout the software development life cycle helps organizations **develop higher-quality applications, reduce human error, and facilitate faster release management processes**. By automating testing, production isolation, and monitoring, manual release management is replaced with a streamlined, consistent process that requires fewer human resources and reduces the potential for errors and delays.

Databricks can utilize the leading Continuous Integration/Continuous Deployment (CI/CD) services, such as Azure DevOps, CircleCI, GitLab, Jenkins, and Octopus Deploy to automate and streamline the development, testing, and deployment of applications and infrastructure as code. Databricks Repos provide git integration and can be integrated into CI/CD pipelines. Databricks Workflows allow users to orchestrate complex data pipelines, and the Terraform provider can be leveraged to provision and manage resources using Infrastructure as code principles.
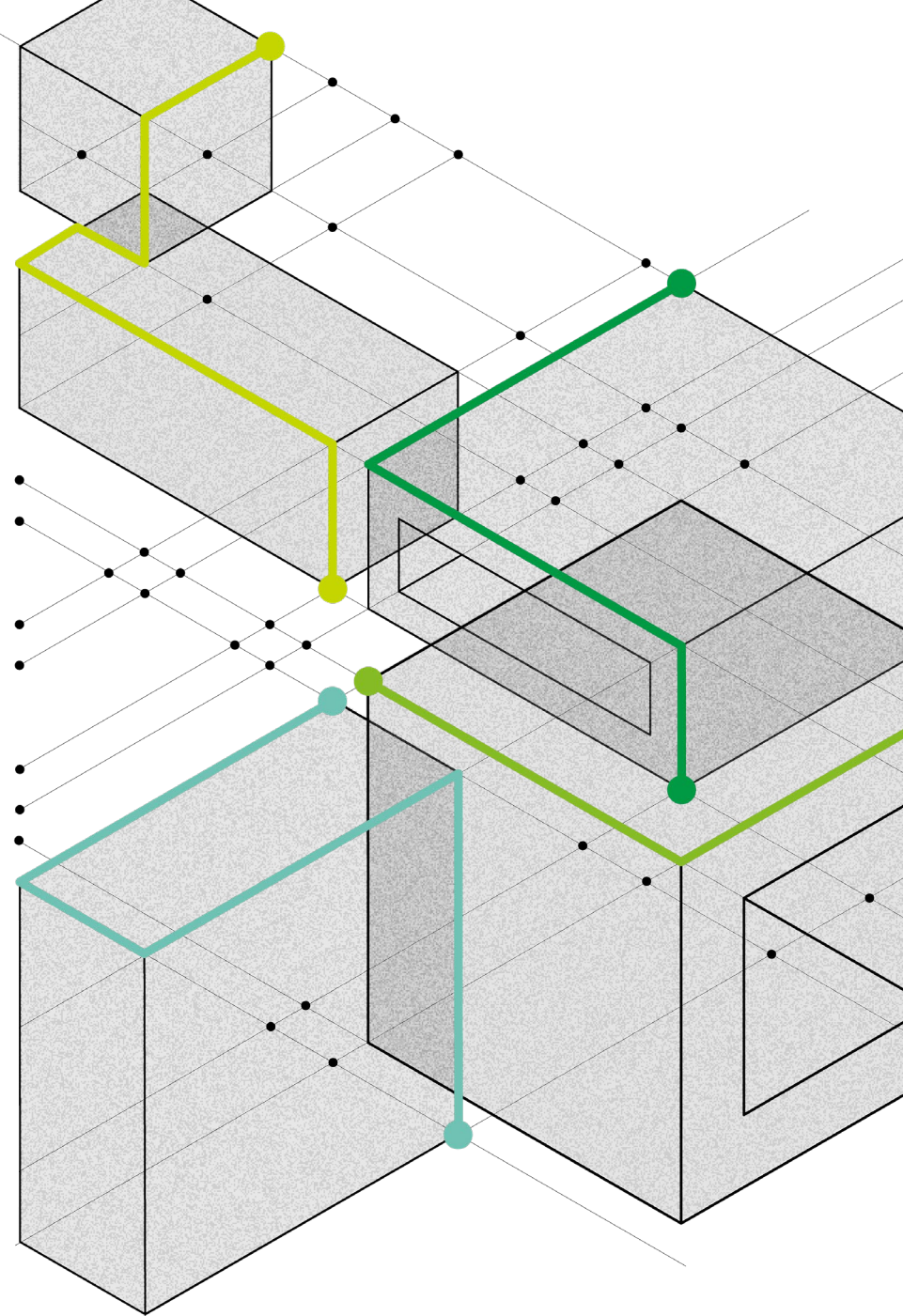
# Data Engineering

## Unlock cloud freedom with cloud-agnostic software

Cloud-agnostic software is constantly evolving. Enterprises are encouraged to leverage tools and technologies that are compatible with any cloud infrastructure and can be moved to and from different cloud environments without any operational issues. It is not mandated to select a cloud-agnostic platform but **very effective when organizations embark on a multi-cloud strategy, given enterprises are moving into multi-cloud**.

Databricks runs on AWS, Azure, and GCP platforms as a multi-cloud platform, allowing applications to be ported without making any changes to their code. However, minimal connectivity and configuration changes may be necessary to integrate with hyperscaler-specific technologies, such as S3 and private networks. Databricks provides a unified user interface across all supported clouds, which reduces the learning curve for users when switching between clouds. The underlying open-source foundations of the platform reduce switching costs. Additionally, the Databricks Data Intelligence Platform provides capabilities that enable data sharing and collaboration across clouds.

Tools and technologies equipped with query engines tailored for Iceberg or Hudi data in multi-cloud environments can access Delta tables via UniForm (preview), eliminating the need for data duplication or conversion.
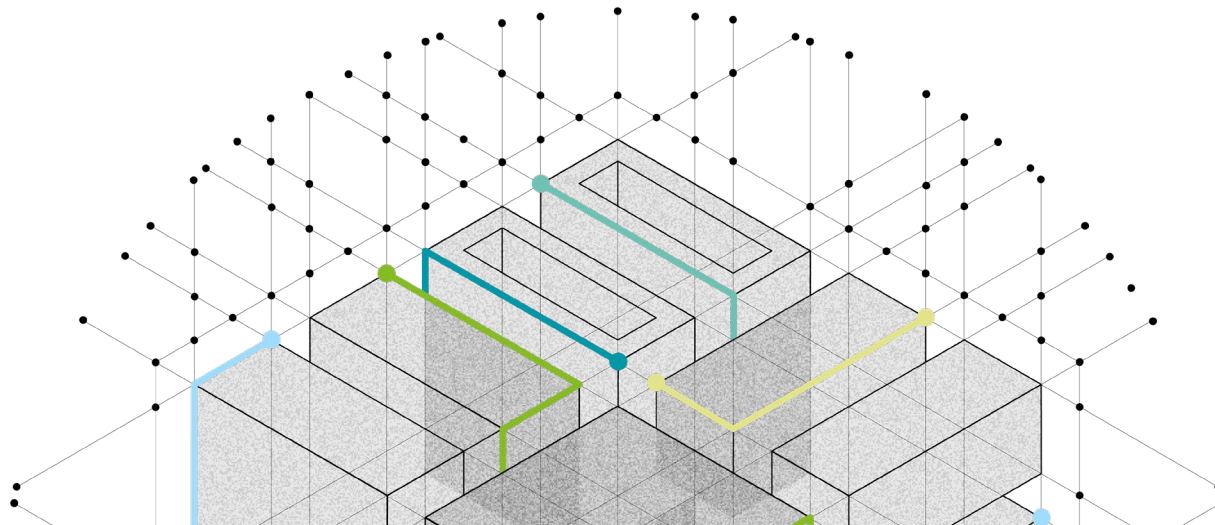
# Data Engineering

Copy data without copying any data!

As data continues to grow into the hundreds of terabytes and petabytes, duplicating it to lower environments and sharing it with consumers becomes increasingly time consuming and costly. Furthermore, there are significant overhead costs associated with verifying the integrity, accuracy, and consistency of the data.

Cloning helps resolve several use cases when organizations struggle to configure multiple environments (Dev, test, UAT and Prod) and manage H/W and S/W, **refreshing the data from production "without duplicating the data," which was taking days to weeks**, and "securely" sharing point-in-time data to another department or external users and have them to **use their own compute power**.

*Clone* is a Databricks-exclusive feature enabled in the Databricks Runtime by default. Databricks offers two types of cloning. Deep clone is a clone that physically copies the source table data to the clone target, including streaming metadata that allows loads to be resumed on the target.

Shallow clones avoid data duplications, and these clones are cheaper to create. Clones are useful for data archiving, ML model reproduction, short-term experiments and data sharing use cases. Cloning is a useful complement to data ingestion capabilities that get data into the lakehouse architecture, such as Partner Connect, Auto Loader and Delta Live Tables.
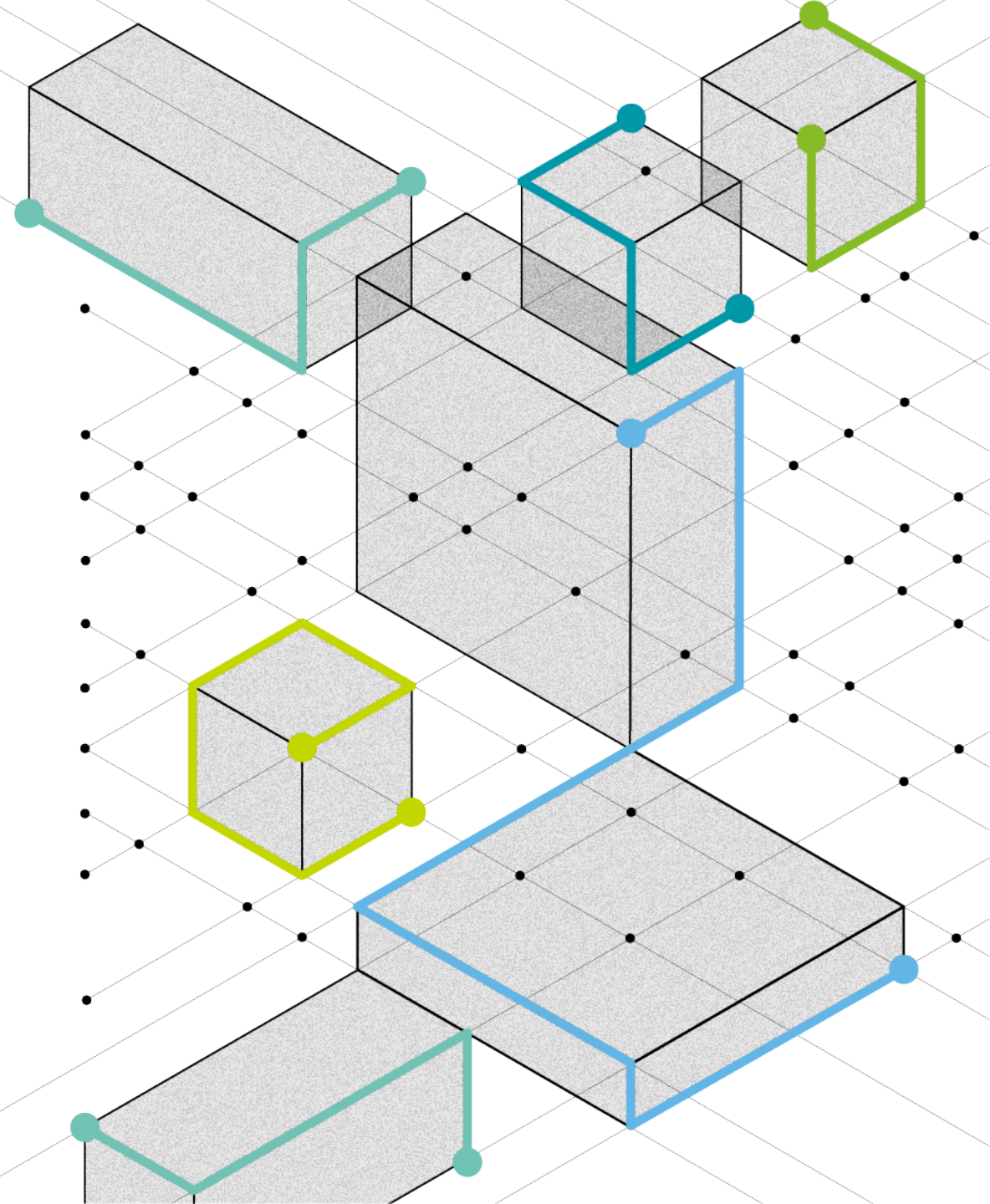
# Data Engineering

## Automate, streamline and schedule

Workflows define the sequence of steps that must be taken to complete a data process, and scheduling ensures that those processes take place at the appropriate times. Hyperscalers are investing and innovating in order to manage the workflows and scheduling software, ensuring that data is transformed efficiently and is available when needed. For this to be successful, these tasks must be repeatable and carried out in the correct order at the right time with the right resources, as well as provide observability and the ability to identify and correct errors.

Databricks allows users to schedule and run notebooks, jobs and other tasks using built-in Workflows. Delta Live Tables provides additional capabilities to build entire ETL pipelines and manage their reliable execution using auto-scaling cloud resources. Both Workflows and Delta Live Tables provide monitoring capabilities and allow for dynamic error-handling. Furthermore, Delta Live Tables gives users the ability to declare data quality rules and specify what should happen when those rules are violated.
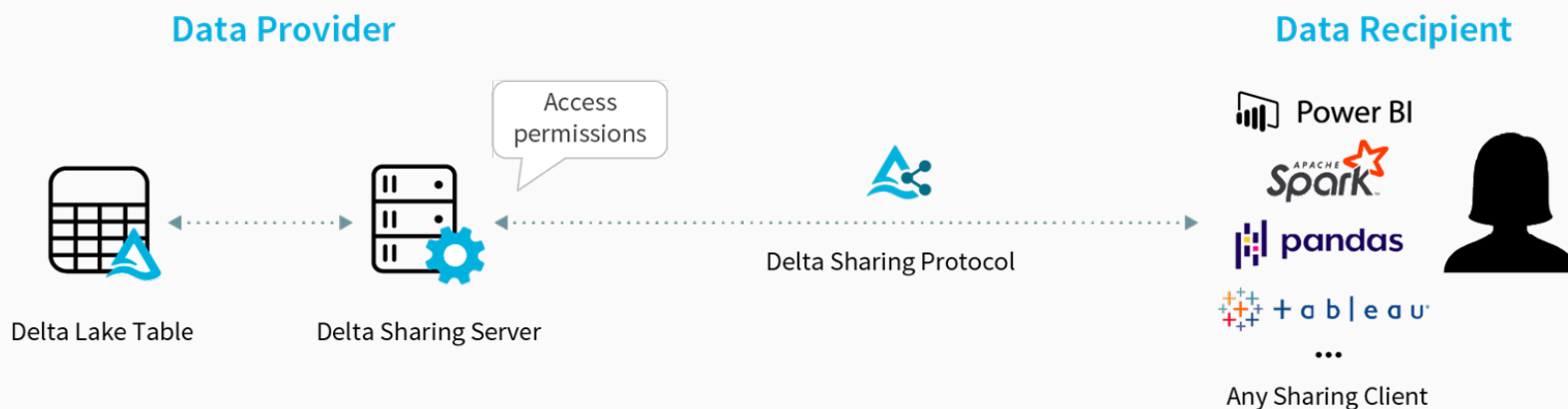
# Data Engineering

## Lead the way with data sharing and clean rooms

Data sharing and clean rooms, the new era of data monetization, is a change in how data is securely exchanged both internally and externally with customers, partners, and producers or applications while maintaining data fidelity across all entities consuming the data. Recent innovations in secure data sharing and clean rooms, without data copy, enable organizations to share data not only with high confidence but to monetize the data as well. Data monetization is expected to grow exponentially. In a recent survey, Forrester Research found that **more than 70% of global data and analytics decision-makers are expanding their ability to use external data, and another 17% plan to do so within the next 12 months**.

*Delta Sharing* is an open protocol that allows for the secure and real- time exchange of large datasets between different products and services. With this protocol, users can securely share data across multiple products without vendor lock-in. Data is shared in place, and data movement is minimized. Additionally, the data marketplace enables producers to monetize the data they produce while also allowing consumers to purchase high quality data at competitive prices.

Databricks clean rooms provide secure, privacy preserving environment that allows arbitrary computations of running in Python, SQL, R, Java, etc., and scale to multiple collaborators and any data size. (Some capabilities are in private preview, and more are expected to be added).



DELTA SHARING, COURTESY DATABRICKS

# Data warehouse and advanced analytics

## Unlock insights to drive business success

By leveraging data and applying advanced analytics techniques, such as predictive analytics, descriptive analytics, and prescriptive analytics, businesses can gain valuable insights to create more efficient processes, improve customer service, and develop new products and services.

**The Databricks Data Intelligence Platform is a cloud-native platform that helps organizations to maximize the potential of their data for all analytics use cases, including business intelligence and AI/ML. Databricks provides a high-performance data warehouse with built-in BI capabilities and integrations with popular BI tools. The platform offers end-to-end capabilities to enable efficient development, deployment, and management of analytics applications.**

# Advanced analytics

## Unleash data insights: data warehouse

In addition to data lakes, a data warehouse is a crucial component of analytics transformation for most organizations. It cleanses, applies business rules, and models data to accurately represent the business while providing efficient access for traditional descriptive business intelligence analytics to enable their users to run their business effectively with reliable data. By unifying data lakes and data warehouses into one platform, numerous advantages can be gained, such as a centralized managed service; increased data security access; and faster, cheaper, and better insights.

Databricks provides a set of data warehouse and BI capabilities that enable users to efficiently store, manage and analyze large volumes of data. *Delta Lake* combines the cost-effectiveness of a data lake with the performance and reliability of a data warehouse. Databricks includes *Photon*, a vectorized engine that accelerates SQL and DataFrame operations to support faster performance for data processing and queries. Databricks has built-in capabilities that enable users to query data and build dashboards.
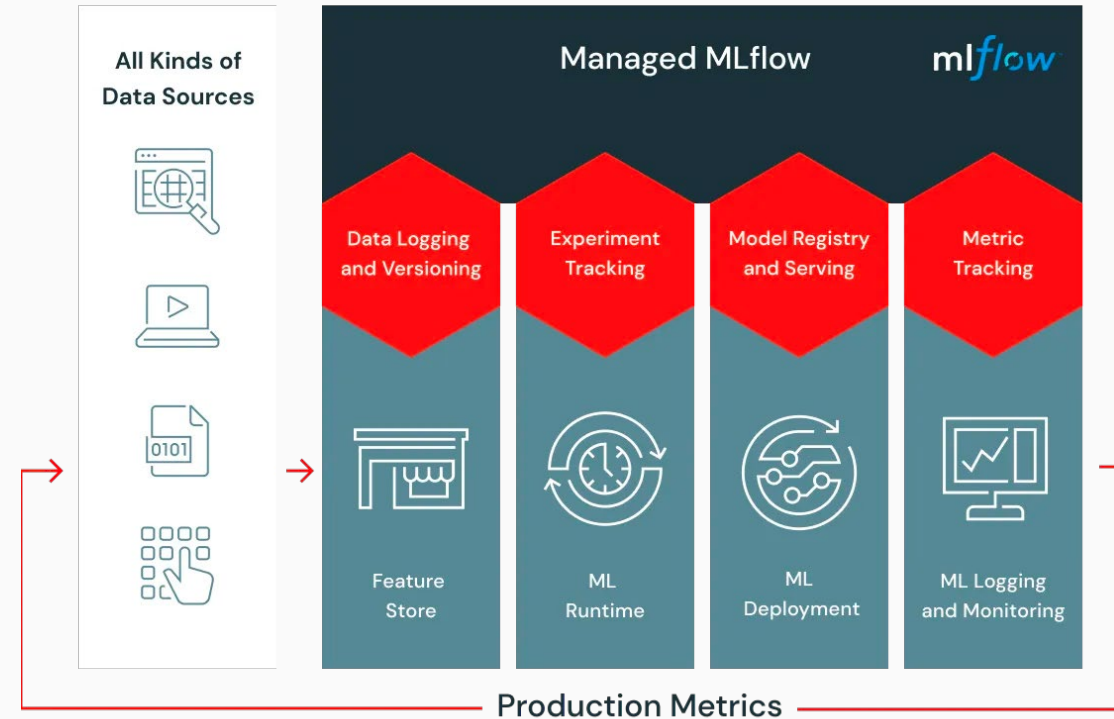
# Advanced analytics

## Embrace AI/ML potential: insights at scale

Advanced analytics, such as predictive, prescriptive, and cognitive analytics, allows organizations to identify patterns, trends, and correlations in their data, enabling them to make more informed decisions and take timely, effective action. Enterprises must evaluate data management software's inherent AI/ML capabilities, **native MLOps (like DevOps), integration with market-leading AI/ML tools, and incorporating multi-compute instances** (GPUs, memory optimized, compute optimized) for the model development and execution. SaaS products can auto-optimize to reduce model training time without specialized skills.

Databricks Machine Learning (ML) enables ML teams to streamline data preparation and processing, facilitate cross-team collaboration, and standardize the ML lifecycle from experimentation to production. Data scientists can easily integrate 3rd party libraries without limitation, including proprietary and open-source large language models (LLMs), manage models through MLOps, and containerize them to run on any platform at scale. Models can be deployed as streaming services for near real-time analytics, enabling faster insights and the ability to react to events such as potential machine failures, fraud, and health alerts. Additionally, Databricks provides the capability to deploy models on other cloud AI services using the MLflow API.



*DATA FOUNDATION FOR THE FULL ML LIFECYCLE, COURTESY DATABRICKS*

# Security and governance

Ensure data protection and regulatory compliance

# Security

## Secure your platform, secure your future

Security of the cloud and on the cloud is getting better and better every day. End-to-end data-security and compliance for the entire enterprise are increasingly important in selecting a data management platform. Key items to be considered are **E2E encryption, role-based access control, single sign-on, audit logs, secure connection for third-party tool integration, no public IP and compliance such as SOC 2 Type 1 & 2 certified, HIPAA-compliant deployment available, ICD-503 INT-A, INT-B, NIST 800-53 rev 4, ARS 3.0 and FedRAMP** (in process), etc.
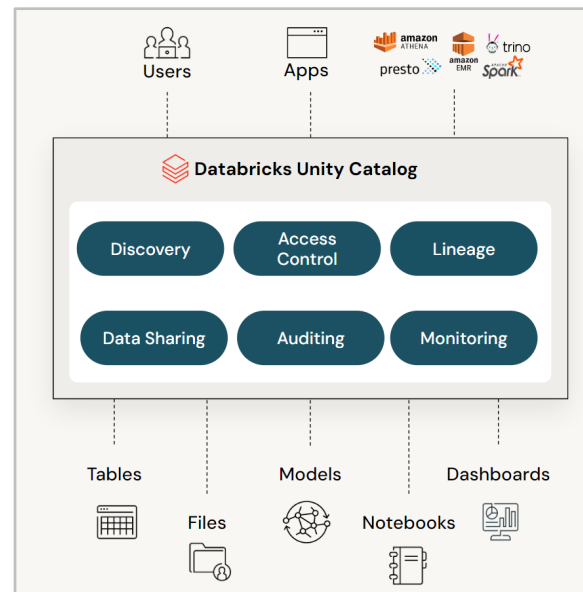
Databricks meets the above security requirements, certifications, and standards. The platform provides a variety of tools for safeguarding the network infrastructure, managing sensitive data and ensuring compliance with regulatory requirements. This includes support for customer-managed keys, private link and robust penetration testing for vulnerabilities.

# Governance

## Set the foundation for success

Governance, observability, data quality, logging, and monitoring is crucial for sensitive data, such as protected and personally identifiable information that is subject to privacy regulations. As organizations struggle with consolidating, locating, and analyzing vastly distributed and diverse data sets, modern enterprise data catalog and data quality tools solve the problem through their augmented ML capabilities. These capabilities enable **faster information access, ensure collaboration and trust, and the analytics** everyone can agree on.

Being able to track and catalog the sources and characteristics of the data sets used to build analytics models, data versions, operational metrics, search, and access management helps to ensure that data is used properly by data scientists, data engineers, and businesses.



## DATABRICKS UNITY CATALOG

Stores operational data, ML models, analytics artifacts, and metadata from other catalogs (e.g. Hive metastore).

### Unified visibility into data and AI
One catalog for all of your data and AI assets making discovery easy, with auto-generated data insights

### Single permission model
Unified and intuitive SQL-based interface to manage access policies for all of your data and AI and audit access platforms. Define fine-grained access controls on rows and columns.

### Data sharing
Share data, notebooks and AI assets across multiple clouds without duplication

### Data lineage
Real-time lineage to see how data flows down to the column level

### AI-powered monitoring and observability
Monitor quality of ML models and data with autogenerated alerts and dashboards
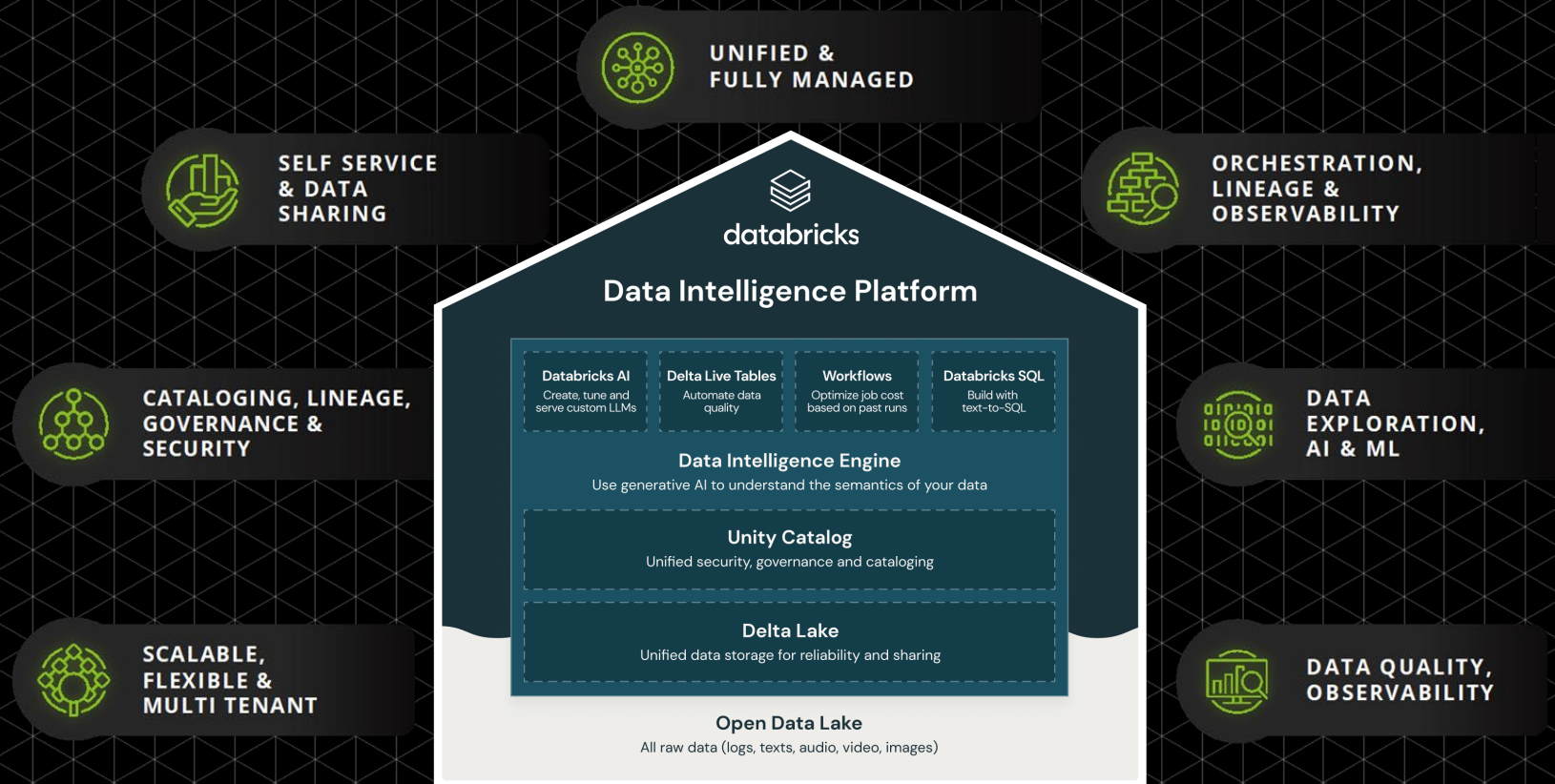
# Additional features

## Expand your horizons

When designing a Data and Analytics Platform (DAP), organizations should take into account the key tenets and capabilities, as well as other sub capabilities that will make the platform complete. **The requirements and priorities of the enterprise sub capabilities can differ from one organization to another**. Such capabilities may include support for multiple programming languages, third- party visualization and virtualization integration, and integration with various data cataloging tools.

**Databricks Federation, currently in public preview, establishes secure connections to data sources, enforces policies in accordance with Unity Catalog access control, and optimally routes queries to the relevant sources as required.**

**Databricks is a Data and Analytics Platform (DAP) that provides the necessary building blocks for enterprises to design a modern DAP. Additionally, Databricks provides a library of solution accelerators that can be utilized to tackle common industry use cases.**

# The Databricks Data Intelligence Platform

The Databricks Data Intelligence Platform is a unified platform that can enable enterprises to achieve AI-driven decision intelligence, rapid innovation, and scalable insights. It is continuously evolving its enhancements to enable users to achieve greater productivity while optimizing resources.



UNIFIED & FULLY MANAGED

SELF SERVICE & DATA SHARING

ORCHESTRATION, LINEAGE & OBSERVABILITY

CATALOGING, LINEAGE, GOVERNANCE & SECURITY

DATA EXPLORATION, AI & ML

SCALABLE, FLEXIBLE & MULTI TENANT

DATA QUALITY, OBSERVABILITY

**databricks**

**Data Intelligence Platform**

**Databricks AI**
Create, tune and serve custom LLMs

**Delta Live Tables**
Automate data quality

**Workflows**
Optimize job cost based on past runs

**Databricks SQL**
Build with text-to-SQL

**Data Intelligence Engine**
Use generative AI to understand the semantics of your data

**Unity Catalog**
Unified security, governance and cataloging

**Delta Lake**
Unified data storage for reliability and sharing

**Open Data Lake**
All raw data (logs, texts, audio, video, images)

# So what? Unlock your potential

## Making a difference with every conversation

Insight into the data leads to better decision-making and problem solving. Thus, the platform should be focused on **solving problems, not just asking questions**.

**Streaming and batch processing**
Improved data processing leads to improved decision-making.

**CI/CD & DevOps**
Enhances user experience by expediting the deployment of new features and updates.

**Cloud-agnostic software**
Simplifies the process of switching cloud providers and accessing data and analytics services.

**Data cloning**
Accelerates testing and development by cloning.

**Data sharing and clean rooms**
Allows secure collaboration and exchange of information.

**AI/ML**
Enables users to use machine learning and artificial intelligence to gain deeper insights into data.

**Self-service capabilities**
Empowers users to quickly access data and analytics without IT dependency.

**Data security**
Ensures that data is secure and only accessible by authorized users.

**Governance**
Enables users to easily enforce data governance policies to ensure compliance with regulations.

**Data warehousing and BI**
Enables rapid analysis and visualization of large-scale data.

**Generative AI**
Exceeding human capabilities and understanding, artificial intelligence advances.

# Conclusion

## Leverage the power of data and analytics for improved decision-making

Organizations should **choose the right data and analytics platform based on their strategic goals, cloud investment, and internally available skills**, rather than just individual tool capabilities, which can be available in every technology sooner or later.

Once a strategy has been finalized, constructing a data and analytics platform necessitates an evaluation of the organization's **data and analytics needs; assessment of the existing infrastructure; identification of any gaps or areas for improvement; and consideration of scalability, security, data governance, and accessibility**. Subsequent to this, businesses should choose technology solutions and vendors that will provide the desired features and capabilities.

# Deloitte.

### Mani Kandasamy

Mani is a technology fellow with Deloitte Consulting LLP AI & Data Engineering Offering and leads @scale cloud data modernization and analytics solutions for a global portfolio of Deloitte's clients.

### Vijay Balasubramaniam

Vijay is a Sr. Partner Solutions Architect at Databricks. He leverages his expertise in data management to help partners and customers be successful in large-scale analytics initiatives.