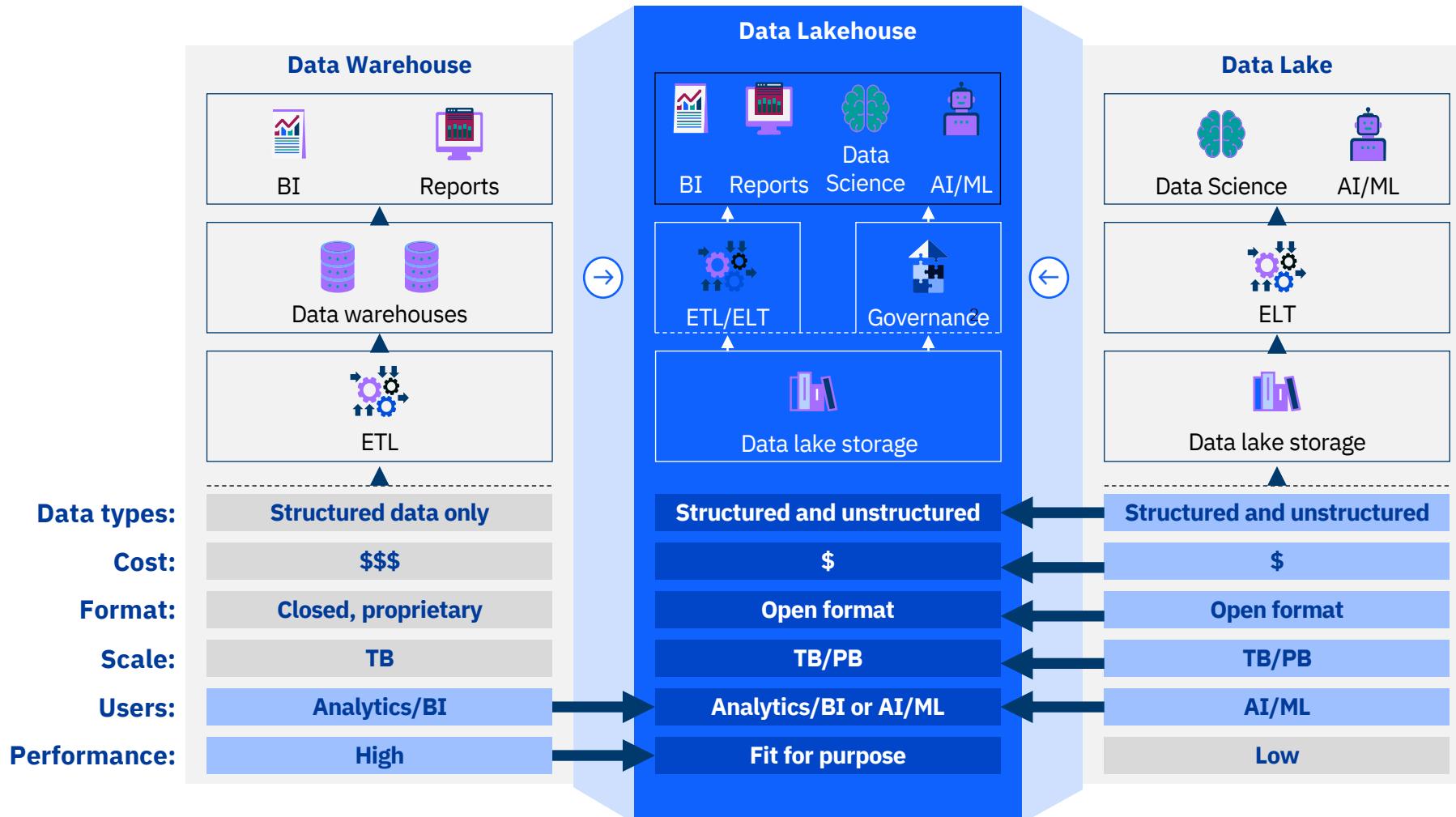


watsonx.data™

IBM Ecosystem Engineering – SI Lab

Lakehouses are meant to be a new class of data store that combines the best of data warehouses and data lakes

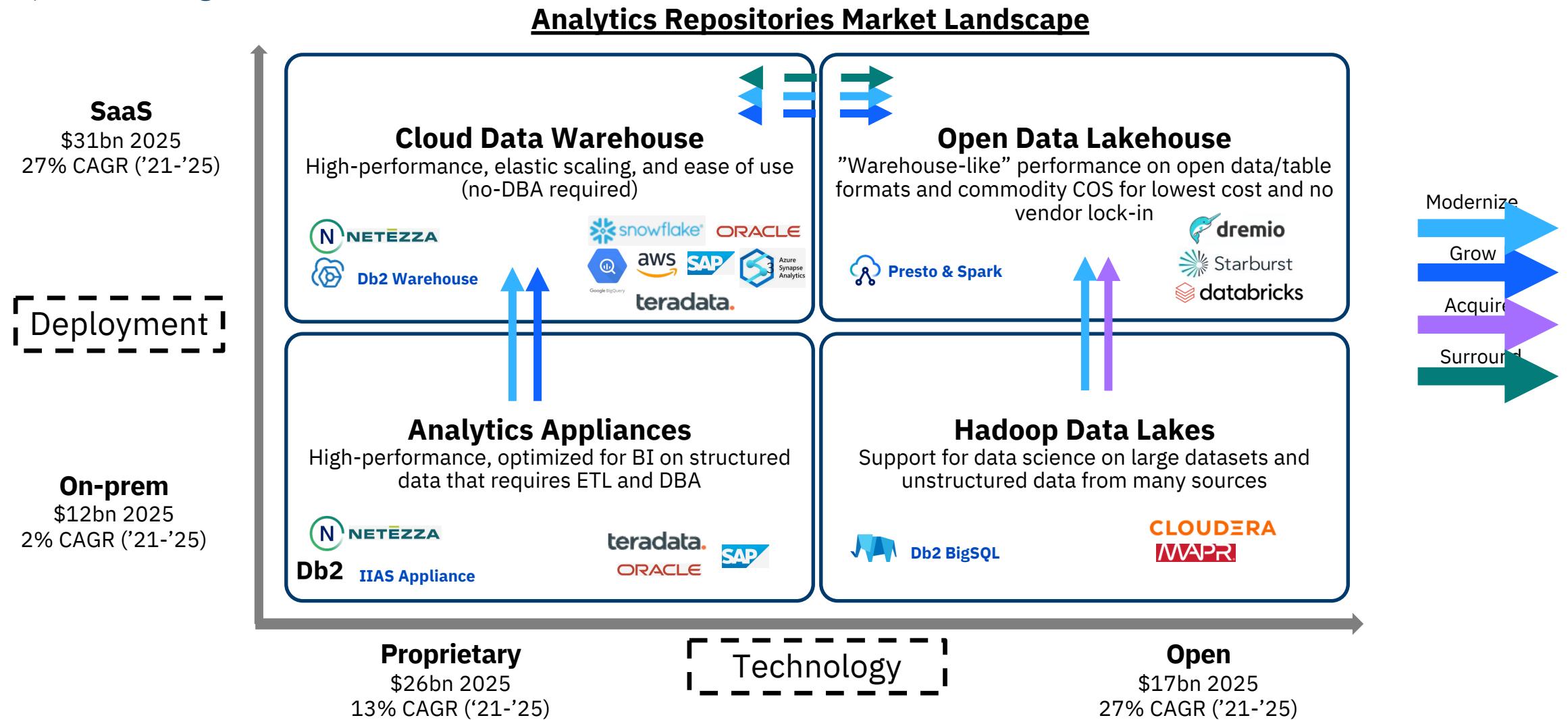


However, first generation lakehouses still have key constraints that limit their ability to address cost and complexity challenges:

- ① Single query engines set up to support limited workloads –typically just BI or ML
- ② Typically deployed over cloud only with no support for multi-/hybrid cloud deployments
- ③ Minimal governance and metadata capabilities to deploy across the entire ecosystem

Market dynamics

Major disruptions are driving the growth in the analytics repositories market from on-prem to SaaS and from proprietary to open technologies



The platform for AI and data

watsonx

Scale and
accelerate the
impact of AI with
trusted data.

- [watsonx.ai](#)
- Train, validate, tune and
deploy AI models

A next generation enterprise studio for AI builders to train, validate, tune, and deploy both traditional machine learning and new generative AI capabilities powered by foundation models. It enables clients to build AI applications in a fraction of the time with a fraction of the data.

watsonx.data

Scale AI workloads, for all your data, anywhere

Fit-for-purpose data store, built on an open lakehouse architecture, supported by querying, governance and open data formats to access and share data.

watsonx.governance

Accelerate responsible,
transparent and explainable
AI workflows

End-to-end toolkit for AI governance across the entire model lifecycle to accelerate responsible, transparent, and explainable AI workflows.

Purpose optimized engines



Data
Engineering



Data
Exploration



Predictive
Analytics



Business
Intelligence

Watson Query (Orchestrator)
VISION/Not-yet-reality



Apache Spark/Flink/Other

- Large scale batch analytics.
- Exabyte scale.
- Data engineering+transformation.
- Low volume batch transactions.



Presto

- Interactive queries+Adhoc analytics.
- Lightweight scalable engines.
- Low volume batch transactions.

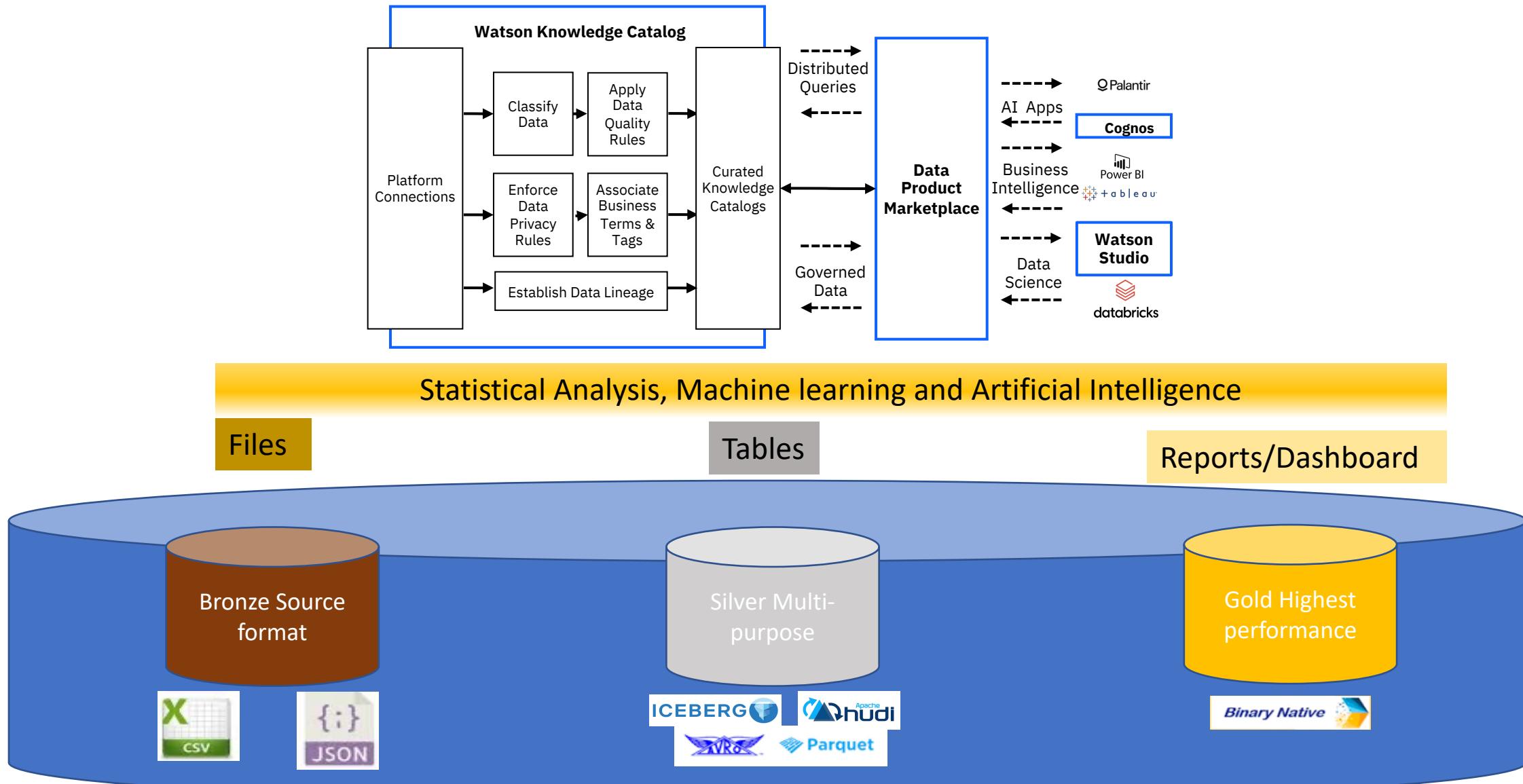


Data Warehouse(s)

- High performance BI+Analytics.
- High Concurrency.
- High volume transactions.

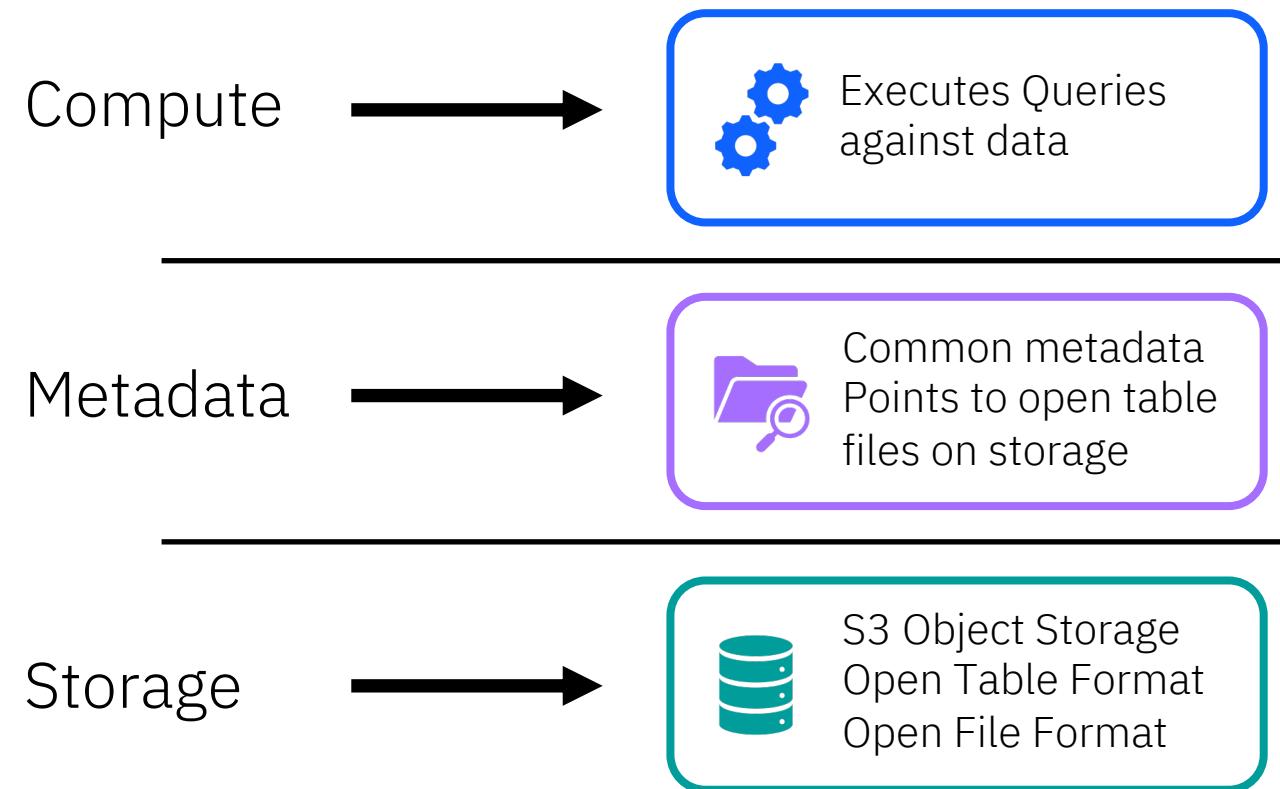


Universal Data Governance for Hybrid/Multi-Cloud



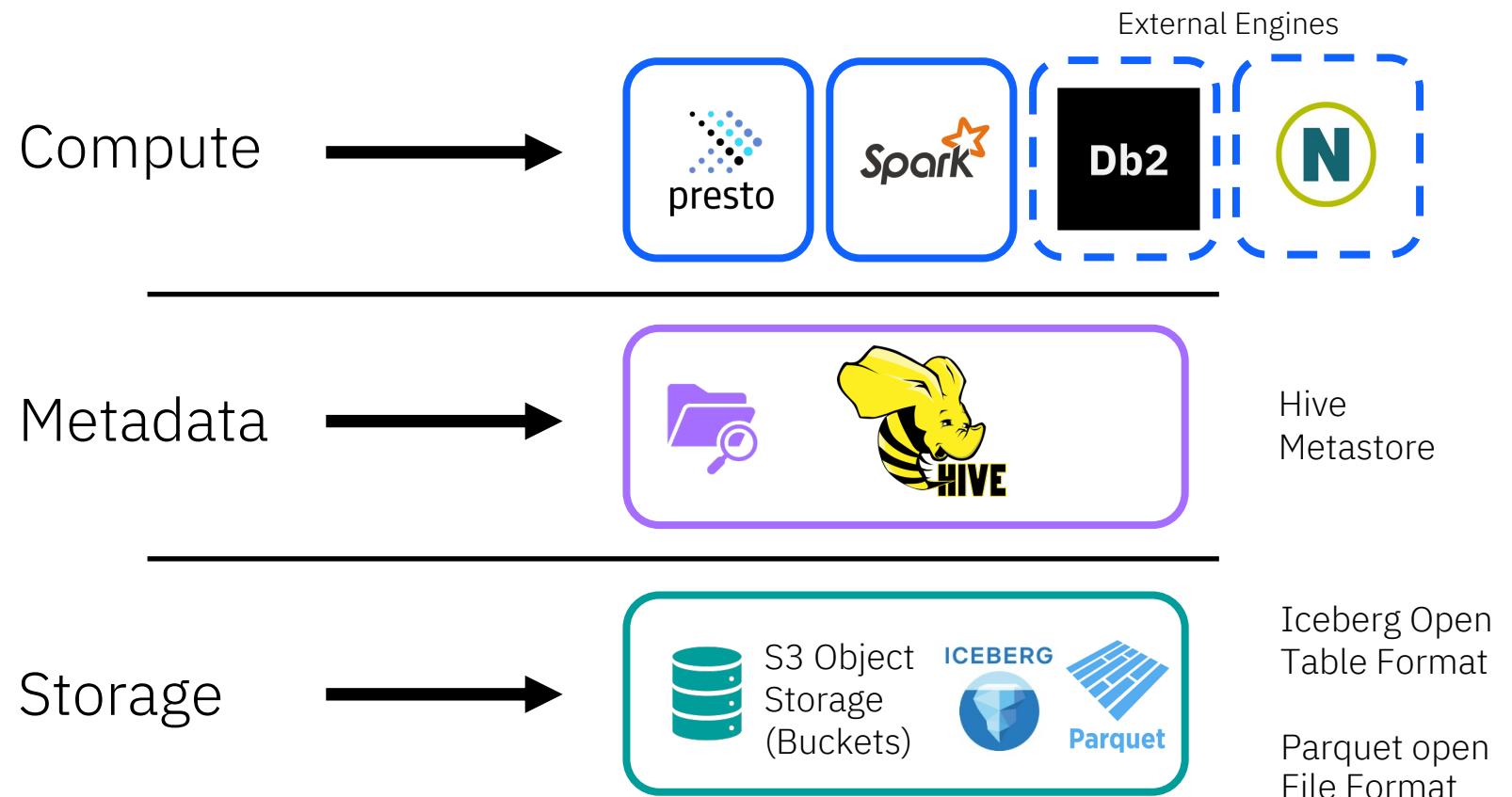
General Lakehouse Architecture

What is a lakehouse architecture?



watsonx.data Architecture

What is the
watsonx.data
architecture?



watsonx.data Architecture

The Presto Engine adds functionality:

Connectors to non-iceberg data files

Connectors to databases (some allowing CRUD operations)

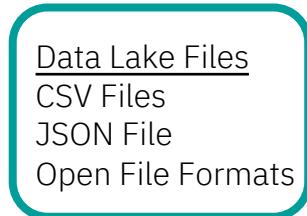
Compute



Metadata



Storage



Hive Buckets



Iceberg Buckets



Databases

watsonx.data Architecture

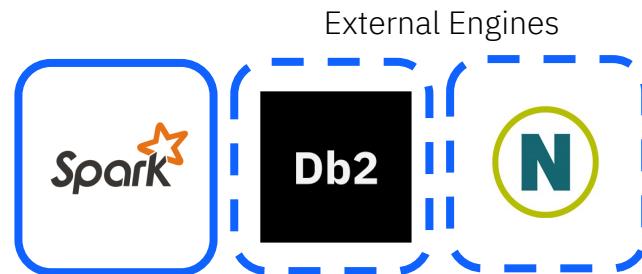
Spark, Db2, Netezza
Engines

Can connect to
Iceberg tables

Can access their own
data objects/table
within their
environments

Cannot access Presto
connectors

Compute



Metadata



Storage



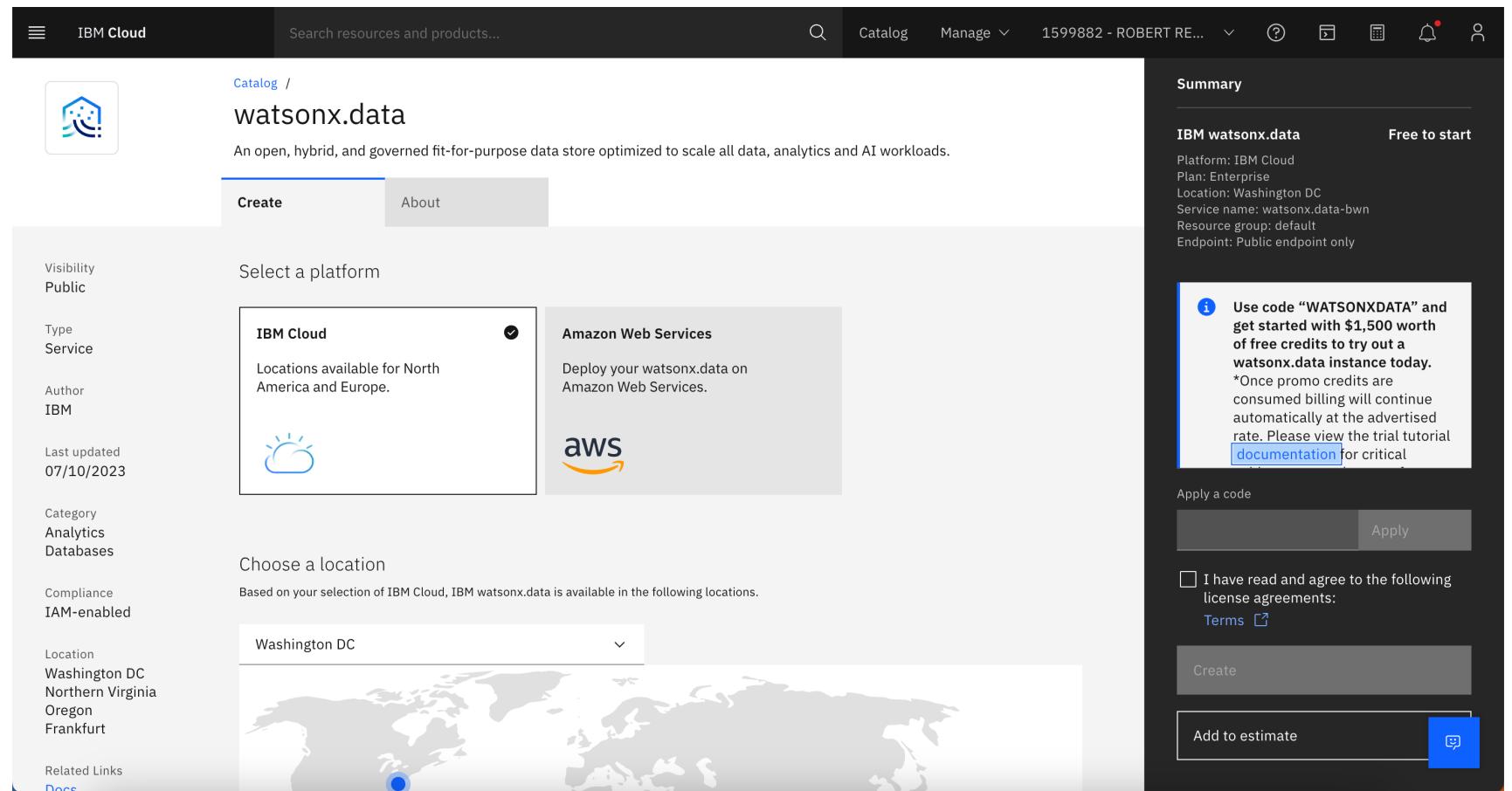
<https://cloud.ibm.com/lakehouse>

Deployment options

SaaS



IBM Cloud
AWS



The screenshot shows the IBM Cloud Catalog interface. At the top, there's a navigation bar with 'IBM Cloud' and a search bar. Below the navigation, the page title is 'Catalog / watsonx.data'. A brief description follows: 'An open, hybrid, and governed fit-for-purpose data store optimized to scale all data, analytics and AI workloads.' There are two main tabs: 'Create' (which is selected) and 'About'. On the left, there's a sidebar with various metadata: Visibility (Public), Type (Service), Author (IBM), Last updated (07/10/2023), Category (Analytics, Databases), Compliance (IAM-enabled), Location (Washington DC, Northern Virginia, Oregon, Frankfurt), and Related Links (Docs). The main content area has two sections: 'Select a platform' where 'IBM Cloud' is chosen (with a note about locations in North America and Europe) and 'Amazon Web Services' (with a note about deployment and the AWS logo); and 'Choose a location' which shows a map with a blue dot indicating the selected location. On the right side, there's a 'Summary' section with details like Platform: IBM Cloud, Plan: Enterprise, Location: Washington DC, Service name: watsonx.data-bwn, Resource group: default, and Endpoint: Public endpoint only. It also features a promotional message about using code 'WATSONXDATA' for credits, a checkbox for license agreements, and buttons for 'Create' and 'Add to estimate'.

- Engine Scaling
 - Multi-Engine Support
 - Caching Support
 - For Enterprise Workloads
-
- Runs on Red Hat OpenShift

Deployment options

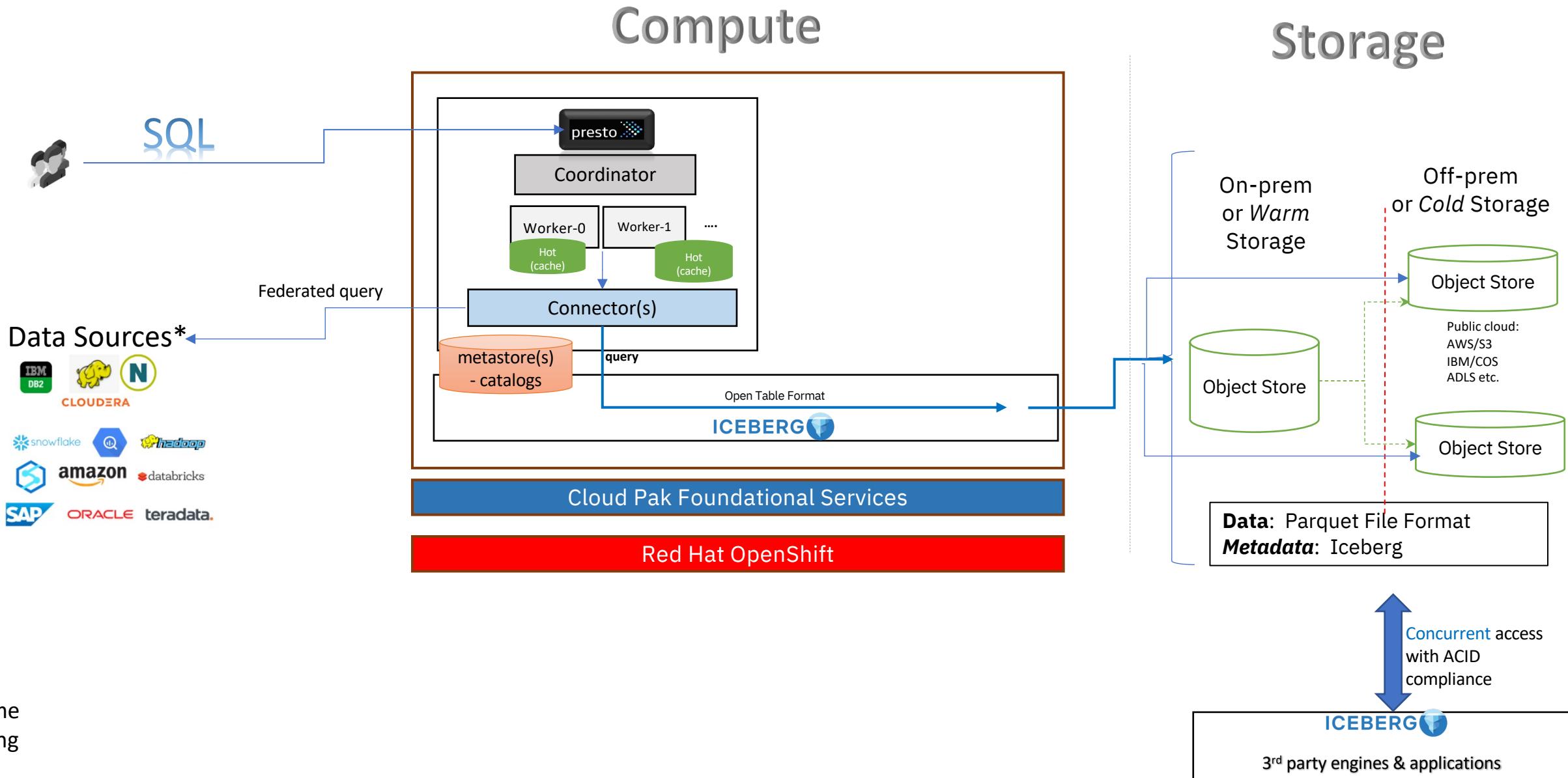
Software
(on-premises)



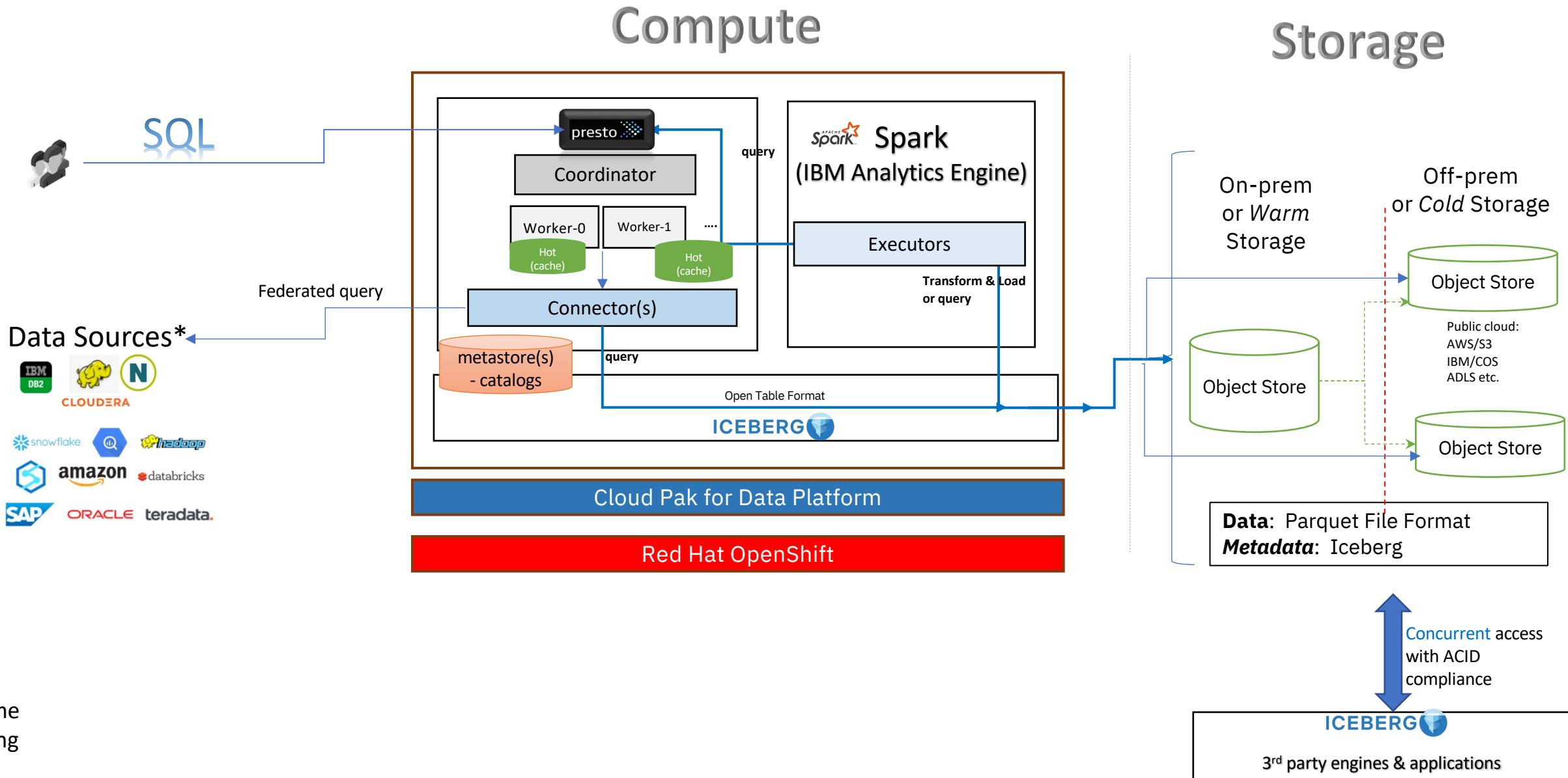
Stand alone or
CP4D Cartridge

The screenshot shows the IBM watsonx.data web application interface. At the top left is the navigation menu with icons for Home, Infrastructure, Catalogs, Buckets, and SQL. The main header says "IBM watsonx.data" and "Welcome, admin." Below this, a message says "You've been logged in for a few seconds." On the left, there's a sidebar with "Architect your lakehouse" (Define and associate infrastructure components to make your data queryable for you and other users) and "Infrastructure manager →". In the center, there's a "Work with your data" section (Build and run queries against and across your data, monitor their progress, and save them for reuse) and a "Query workspace →". At the bottom left, there's a "Welcome to IBM watsonx.data." section with a brief description and "Recommend resources" links to FAQs, Release notes, and API specs. To the right, there are four cards: "Infrastructure components 5" (Engines: 1/1 running, Catalogs: 2/2 queryable; Buckets: 2/2 queryable), "Recent tables 0" (No recent tables. Use the data manager to explore and curate tables across catalogs.), and "Recent ingestion jobs 0" (No recent ingestion jobs. Create an ingestion job to move external data into watsonx.data.). The top right corner of the interface includes the IBM logo, version information ("Version v4.7"), and links to "Version details", "Release notes", and "License agreements". The bottom right corner credits "Copyright IBM Corporation 2018, 2023" and "Powered by".

watsonx.data on Openshift – stand alone



watsonx.data on Cloud Pak for Data



- Single Presto Engine
- Cannot Scale
- Single User
- No Caching
- Used to test/play
- Not for Production

- Install on Linux (Recommended)
- Runs on Docker/Podman

Deployment options

Developer Edition

Runs on Podman or Docker

The screenshot shows the IBM watsonx.data developer edition interface. On the left, there's a sidebar with icons for Home, Catalogs, Databases, SQL, and Infrastructure manager. The main area has a dark background with several sections:

- Welcome, ibmlhadmin.**: A message indicating you've been logged in for a few seconds.
- Architect your lakehouse**: A section for defining and associating infrastructure components.
- Work with your data**: A section for building and running queries against data.
- Welcome to IBM watsonx.data.**: A summary of recommended resources and integration with the API suite.
- Recommend resources**: Links to FAQs, Release notes, and API specs.
- Infrastructure components 7**: A grid showing 1/1 Engines (running), 3/3 Catalogs (queryable), 3/3 Buckets (queryable), and 0/0 Databases (queryable).
- Recent tables 0**: A section stating "No recent tables." with a note about using the data manager to explore and curate tables across catalogs.
- Recent ingestion jobs**: A section stating "Available only via CLI." with a note about data ingestion from the console.
- Cloud resource name (CRN)**: 0000-0000-0000-0000
- Copyright IBM Corp. 2023**
- Console version**: 1002-20230706-024004-sw_dev_ent

File Formats



- Row based
- Indexed
- Less compression than Parquet and ORC but faster write speeds
- Up to 3x faster read times
- Good schema evolution support
- Can be used as a landing area prior to additional data transformations



- Columnar
- Originally also designed for Hadoop
- Works very well with Spark
- Good for traditional OLAP queries because of its columnar format
- Carries schema, is self describing
- Data stored in pages
- Good for complex nested data structures



- Also columnar
- Data stored in stripes
- Indexed rows, compatible with HDFS
- KPIs for compression and runtime are very similar to parquet
- Parquet better for Spark, ORC better for Hive/Hadoop

Open Source Table Formats

- Separation of compute, data, and storage
 1. Leverage low-cost, infinitely scalable object storage
 2. Standardized
 - open file formats (Parquet, ORC, DWRF, JSON, ...)
 - table formats (Apache Iceberg, LF Delta, Apache Hudi)
 3. Accessed by scalable compute engines of choice (Presto, Spark, etc.)



Sample Discovery Questions

- What type of database are they using today
- Is cost an issue or starting to become an issue
- What are the use cases and how do they serve the business
- What are the required SLAs? Who are the users accessing?
How fast do they need answers?
- Are there ETL workloads within the EDW, what are the time windows
- Is there historical data in the EDW that is not frequently accessed but still required to be there?
- Overall data size and instance size

IBM Modernize + Cross Sell

Customer Goal: Customer has an IBM on-premises solution and want to move to next generation of self-managed and SaaS offerings

Customer motivation:

- Looking for next generation features
- Limited compute or data storage resulting in offloading or expanding data infrastructure
- Difficulty scaling up and down for workloads
- New use cases: AI/ML, Real-time analytics, data engineering, data sharing

Challenges to Address:

- System migration compatibility
- Sizing management
- Cost controls on the cloud
- Access to all data across hybrid cloud without duplication/movement

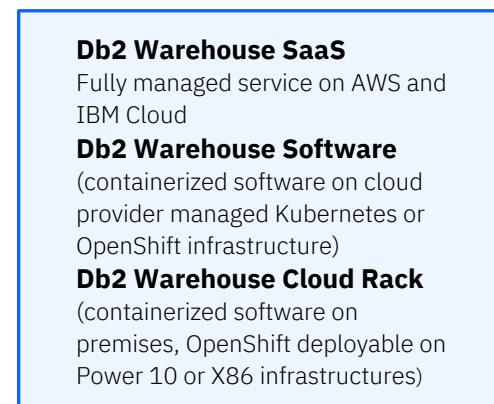
Key Messages:

1. Modernize your database appliance with like for like compatibility and support new use cases
2. Optimize workloads for AI with fit for purpose engines and reduce warehouse costs by 50%, access all governed data across hybrid cloud



Modernize appliance to Db2 Warehouse SaaS, Software or Cloud Rack

Modernize software to Db2 Warehouse SaaS

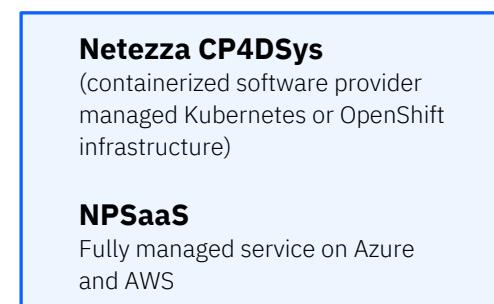


+



Modernize Mako to Netezza CP4D Sys (Hammerhead) or SaaS

Modernize Hammerhead to Netezza SaaS



+

IBM Db2 for z/OS Cross-sell

Customer Goal: Customer has IBM Db2 for z/OS and wants to access transactional data from mainframe for AI use cases

Customer motivation:

- Brand new use cases with mainframe transactional data: AI/ML, Real-time analytics, data engineering, data sharing

Challenges to Address:

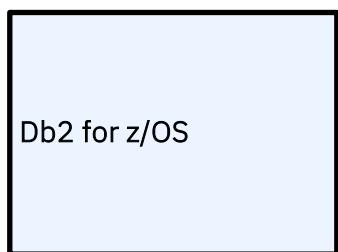
- System migration compatibility

Key Messages z/OS and Db2 Warehouse:

- An integrated, optimized synchronization feature maintains currency between source Db2 for z/OS data on IBM Z and Db2 Warehouse targets
- Db2 Warehouse is a relational databases that delivers advanced data management and analytics capabilities

Key Messages z/OS and .data:

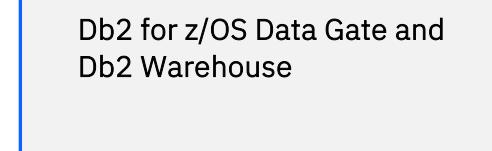
- Access and analyze data from mainframe systems in near-real time
- Consume less processor capacity with built-in synchronization
- Use the most up-to-date data from mainframe for machine learning models



Db2 Data Gate provisions and enables a Db2 Warehouse service



Db2 Data Gate provides synchronized Db2 for z/OS data to IBM Cloud Pak for Data with data stored and managed in Db2 Warehouse.



Watsonx.data can readily access this data via included connectors.

Competitors to IBM data warehouse & watsonx.data

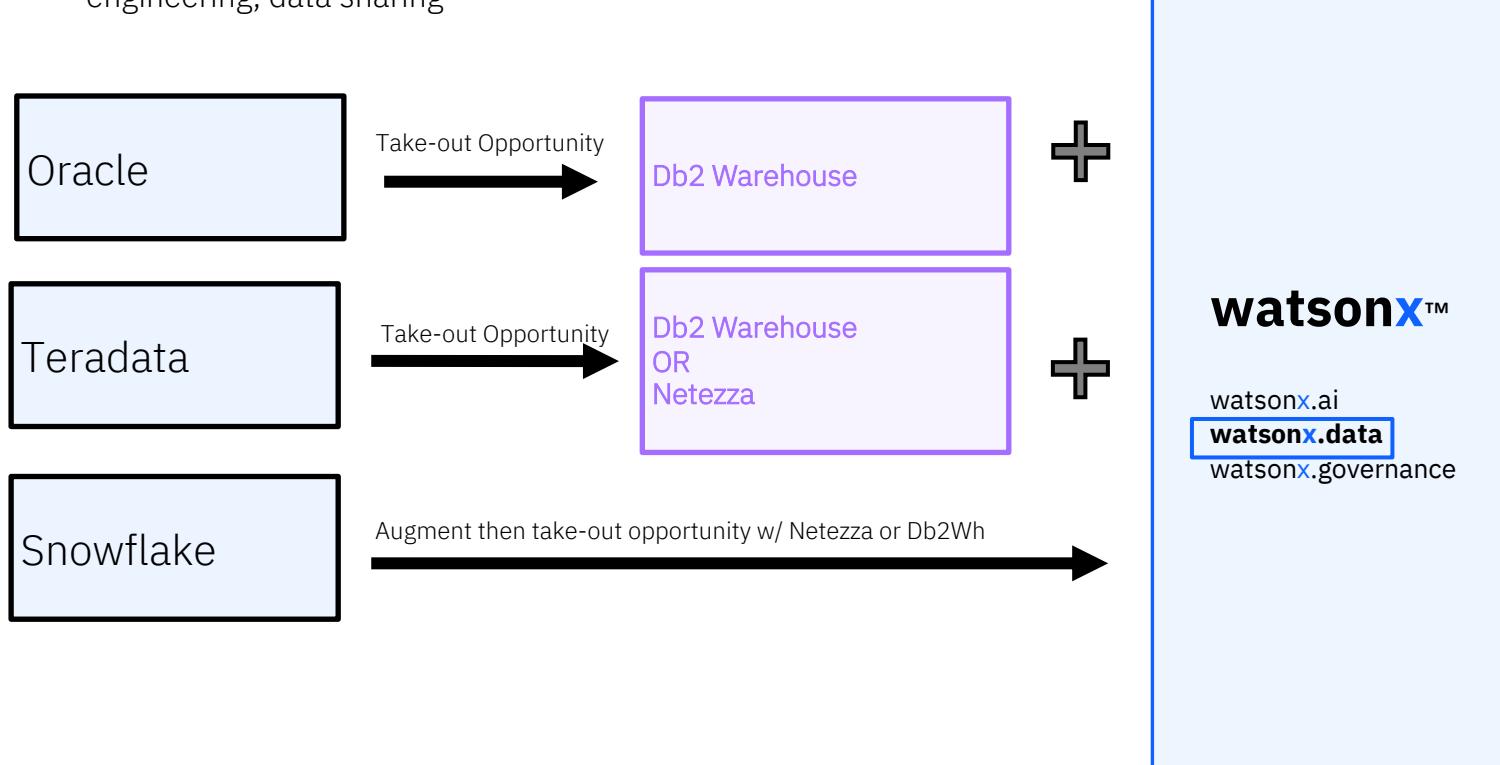
Customer Goal: Customer has competitor solution and is dissatisfied with expense of current solution or desires an IBM solution

Customer motivation:

- Costs too high with existing software or cloud solution
- Difficult relationship with existing software provider
- Need flexibility to scale new use cases and all workloads running on same instance: AI/ML, Real-time analytics, data engineering, data sharing

Challenges to Address:

- System migration compatibility
- Cost controls on the cloud - this is the # one complaint we hear from Snowflake customers!
- Sizing management
- Access to all data across hybrid cloud without duplication/movement



Key Messages:

Why watsonx.data?

- Access all data for AI across hybrid cloud
- Optimize workloads for AI and reduce warehouse costs by 50%
- Built-in governance
- Leverage watsonx platform for trusted AI and generative AI capabilities

Why Db2 Warehouse?

- Simplify migration with Db2 compatibility mode (Oracle, Sybase, PostgreSQL)
- Fully managed SaaS offerings, or containerized deployments for self-managed software
- Seamless integration with .data and IBM ecosystem
- Use object storage to save on storage costs

Why Netezza?

- Leverage similar PostgreSQL construct to reduce migration complexity
- Fully managed SaaS offerings, or containerized deployments for self-managed software
- AI-infused elastic scaling for cost predictability
- Seamless integration with .data and IBM ecosystem

Concepts and market entry points

Three key concepts for IBM watsonx.data

1. Presto is a next-generation open-source SQL engine designed to run efficiently over data lakes.
2. Warehouses and first-generation lakehouses are monolithic, and not optimized to work on all workloads. Only IBM watsonx.data's multi-engine architecture allows for true workload optimization.
3. Iceberg is an open-table format that allows multiple engines to access the same data – this means, Snowflake, Netezza, and IBM watsonx.data can all access data in Iceberg at the same time.

Market entry points + use cases

1. Warehouse Optimization narrative

- Competitive and IBM (Netezza + Db2)
- Talk track – Client is concerned with the spend on traditional warehouse today – looking to optimize for both performance and cost
- Value prop: Cost optimization and openness through the shared meta layer and fit-for-purpose engines
- Use case – Snowflake write-intensive workloads moving to Spark and/or Presto, thus reducing cost of Snowflake virtual data warehouses

2. Modernizing data lake narrative

- Modernizing storage architecture to facilitate shared metadata and fit-for-purpose engines
- Talk track – Converting legacy file storage structures into open-file structures and assigning those into the shared meta layer, thus facilitating fit-for-purpose engines
- Value prop: Modernize with warehouse-like performance for querying data in open formats, built-in governance
- Use case – Move from HDFS to open-source iceberg within a consolidated shared meta layer

Business Drivers



250%

The aggregate volume of data stored is set to grow over 250% in the next 5 years.

82%

82% of enterprises are inhibited by data silos.

80%

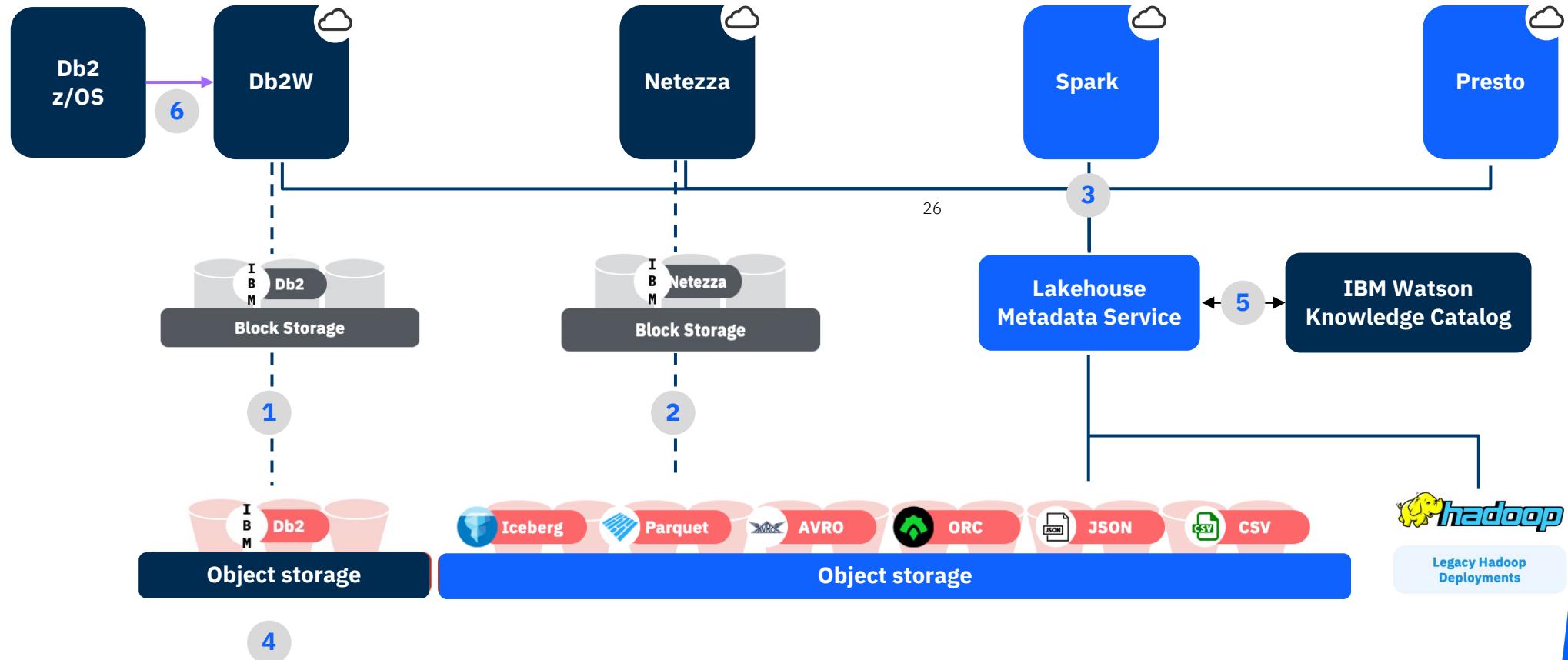
80% of time is spent on data cleaning, integration and preparation.

82%

82% of enterprises say data quality is a barrier on their data integration projects.

The integrated IBM watsonx.data ecosystem for maximum workload coverage and optimal price-performance

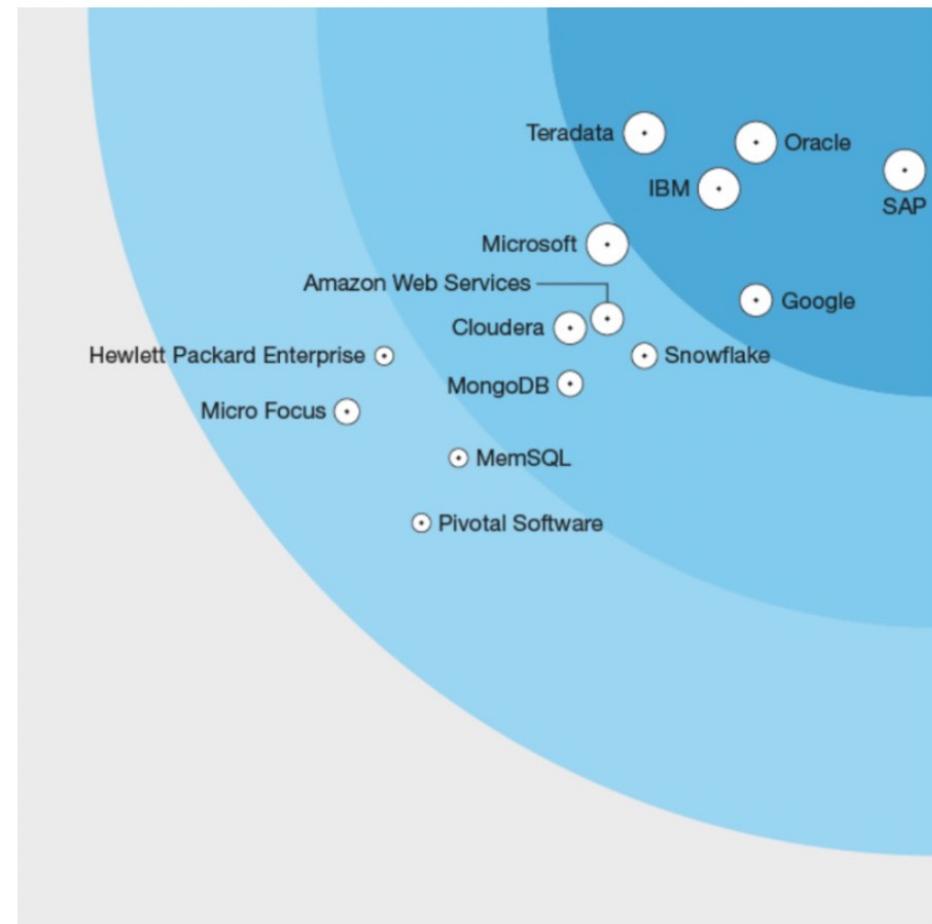
IBM watsonx.data functionality Integrations at GA



- 1 Warehouses can access data in the lakehouse
- 2 Easily Promote data between the warehouse and lakehouse
- 3 Query routing service, multiple engines can access same data lake data
- 4 The lakehouse can access data residing in Db2/Netezza
- 5 WKC policies enforced by the lakehouse via metadata service
- 6 Analyze z data easily and securely with Db2 for z/OS Data Gate



Gartner Magic Quadrant for Data Quality Solutions



Forrester wave: Data Management for Analytics

watsonx.data

Differentiators

- No other data lakehouse offering has **integrated data warehouse engines** in addition to the Apache Spark and open-source query engines
- The cloud hyperscalers (AWS, Microsoft Azure, and GCP) along with Databricks provide no **hybrid cloud deployment capability**
- **Deployment flexibility in other clouds** – no other data lakehouse offering can be deployed as easily across different cloud platforms
- Other data lakehouse competitors do NOT have the level of experience with mission critical applications, and **level of research in query optimization and query processing**, as IBM
- Watsonx.data also can be coupled with **Watson Knowledge Catalog** – a leader in the Gartner Magic Quadrant
- Watsonx.data and its selection of Apache Iceberg and Presto delivers **an open solution versus a single contributor open-source lock-in**