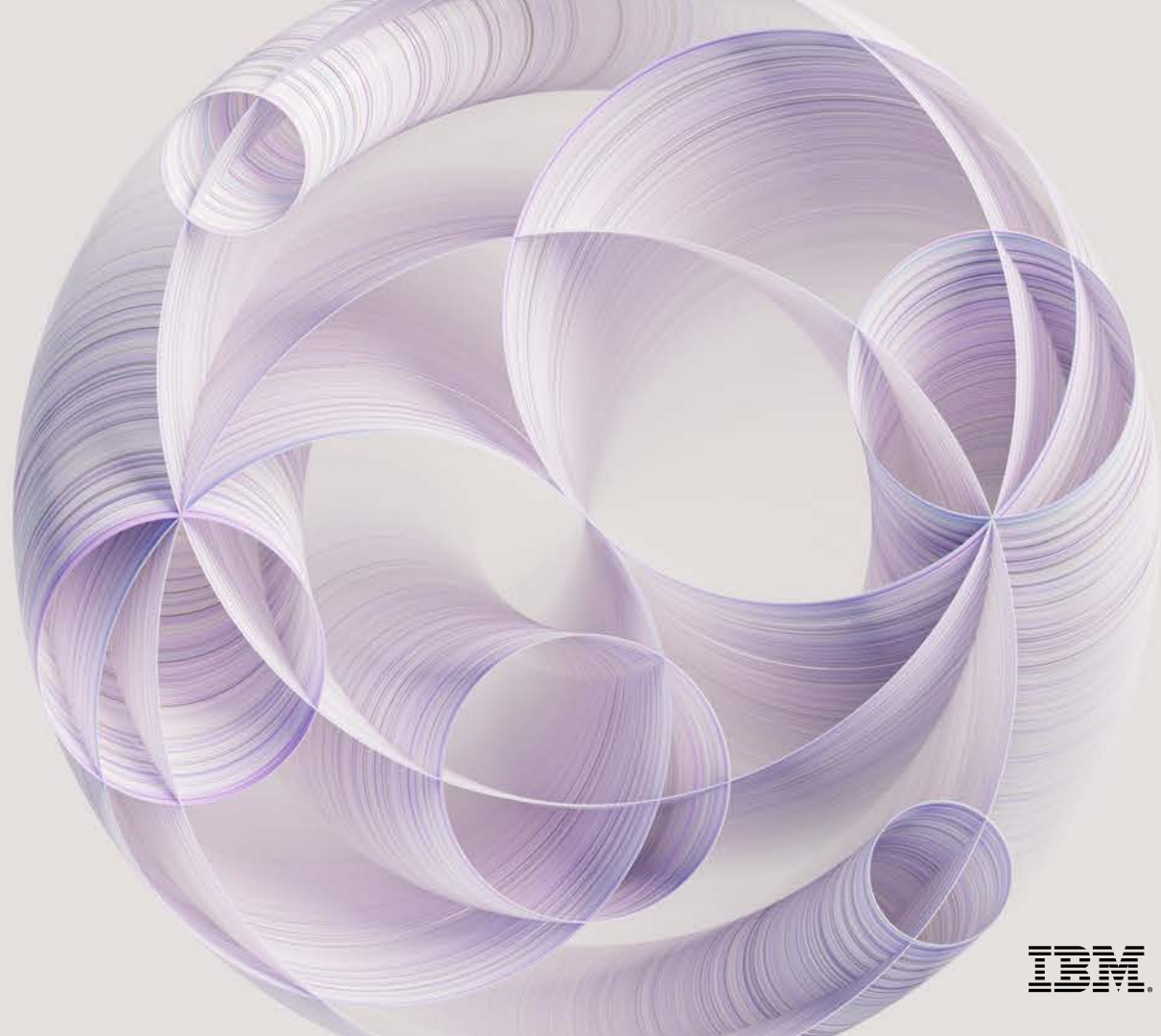# Watsonx.data

Day 1
Fundamentals
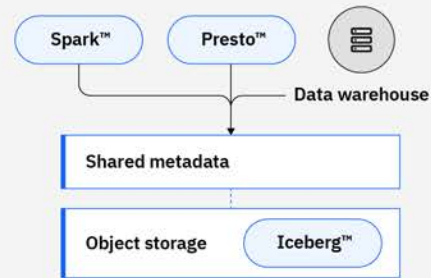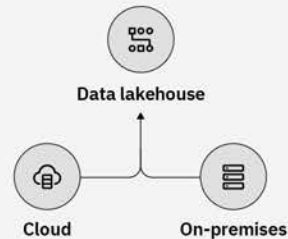
IBM

# Access all your data across hybrid-cloud through a single point of entry

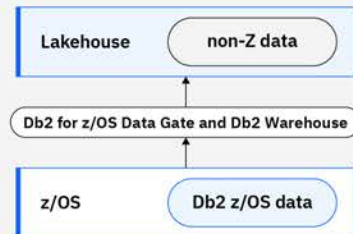An open data store built for hybrid deployment of your analytics and AI workloads

1. Share a single copy of data with tools that can read open data formats to minimize data duplication

Spark™  Presto™  Data warehouse

Shared metadata

Object storage  Iceberg™

2. Connect to and access data remotely across hybrid-cloud with the ability to cache remote sources

Data lakehouse

Cloud  On-premises

3. Synchronize and incorporate Db2 for z/OS data for lakehouse analytics.

Lakehouse  non-Z data

Db2 for z/OS Data Gate and Db2 Warehouse
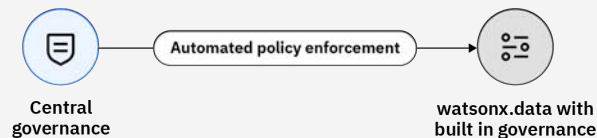
z/OS  Db2 z/OS data

# Get started in minutes with built-in governance, security and automation.

Accelerate time to trusted analytics and AI

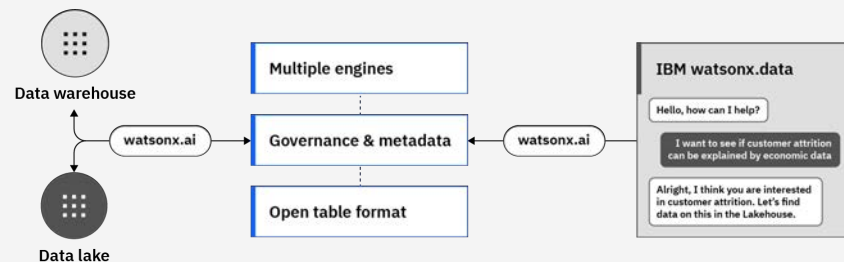Connect to your existing analytics data and deploy fit-for-purpose engines in minutes



**Object storage** → **Metadata and access controls** → **Query engines**

Address enterprise compliance and security using built-in centralized governance across your data ecosystem



**Central governance** — Automated policy enforcement → **watsonx.data with built in governance**

Use foundation models to discover, augment, refine, and visualize watsonx.data data and metadata



**Data warehouse**

watsonx.ai

**Data lake**

Multiple engines

Governance & metadata

Open table format

watsonx.ai

**IBM watsonx.data**

Hello, how can I help?

I want to see if customer attrition can be explained by economic data

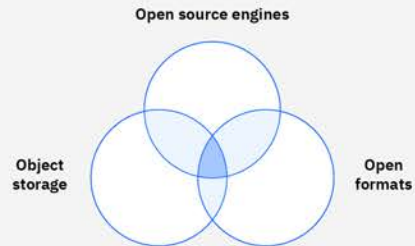Alright, I think you are interested in customer attrition. Let's find data on this in the Lakehouse.

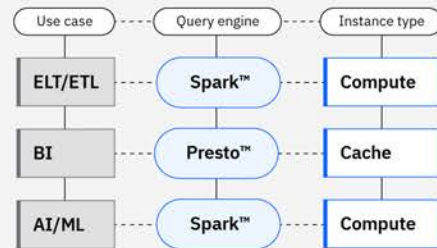# Reduce your data warehouse costs by up to 50%* by optimizing workloads

Optimize workloads from your data warehouse when you take advantage of low-cost object storage and fit-for-purpose query engines

*When comparing published 2023 list prices normalized for VPC hours of IBM watsonx.data to several major cloud data warehouse vendors. Savings may vary depending on configurations, workloads and vendors.
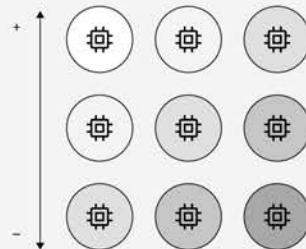


① Share data between multiple analytics engines

Open source engines

Object storage

Open formats

② Use fit-for-purpose compute and cache-optimized instances

| Use case | Query engine | Instance type |
|----------|--------------|---------------|
| ELT/ETL | Spark™ | Compute |
| BI | Presto™ | Cache |
| AI/ML | Spark™ | Compute |

③ Scale up and scale down automatically

# Common data file formats

Computer systems and applications store data in files

Data can be stored in binary or text format

File formats can be open or closed (proprietary/lock-in)

Open formats (Parquet, ORC, and Avro) are commonly used in data lakes and lakehouses

## CSV

- Human-readable text
- Each row corresponds to a single data record
- Each record consists of one or more fields, delimited by commas

## { JSON }

- Human-readable text
- Open file and data interchange format
- Consists of attribute-value pairs and arrays
- JSON = JavaScript Object Notation

## Parquet

- Open-source
- Binary columnar storage
- Designed for efficient data storage and fast retrieval
- Highly compressible
- Self-describing

## Apache ORC

- Open-source
- Binary columnar storage
- Designed and optimized for Hive data
- Self-describing
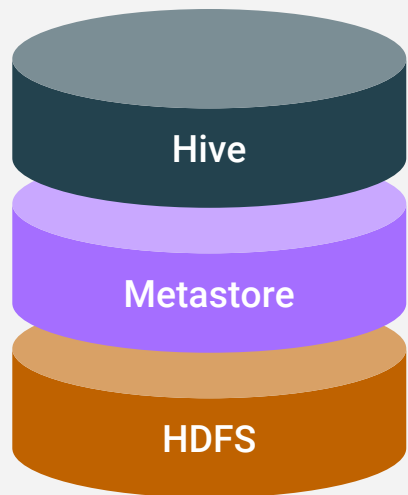- Similar in concept to Parquet

## Avro

- Open-source
- Row-oriented data format and serialization framework
- Robust support for schema evolution
- Mix of text/binary

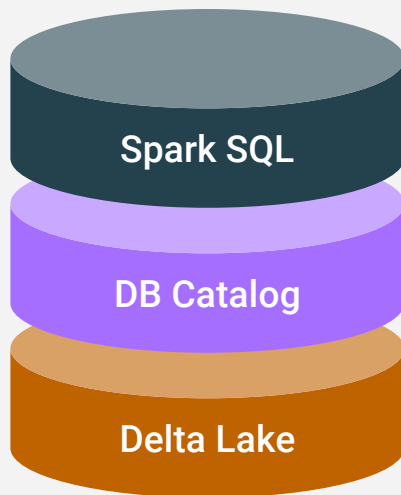# THE OPEN DATA LAKEHOUSE DIFFERENTIATOR

"IBM / Cloudera Shares a Joint Vision"
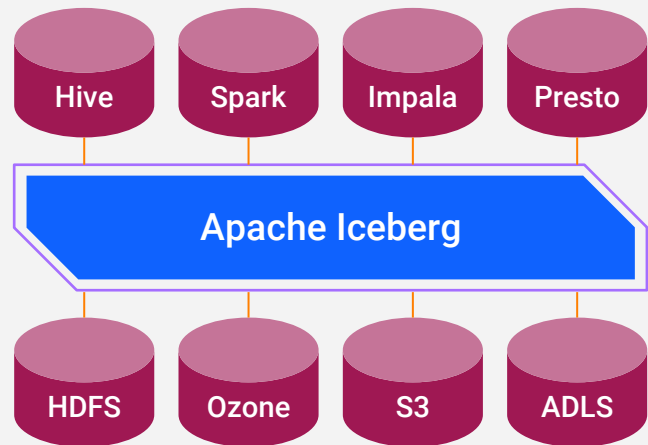


The "Old" Way — Hive SQL over data in HDFS: Hive / Metastore / HDFS

"Some" other way — Spark SQL over Delta tables: Spark SQL / DB Catalog / Delta Lake

The Iceberg way — Multi-function analytics over all your data: Hive, Spark, Impala, Presto → Apache Iceberg → HDFS, Ozone, S3, ADLS

**Powerful Engines**

Open-source Engines for each use case

**Catalog**

IBM use Hive Metastore to catalog Iceberg Table Formats

04

ICEBERG

**Table Format**

Organize the data where it lives into tables using Iceberg

03

02

**Parquet**

**File Format**

Store Data in Apache Parquet Files (Open Columnar Format)

01
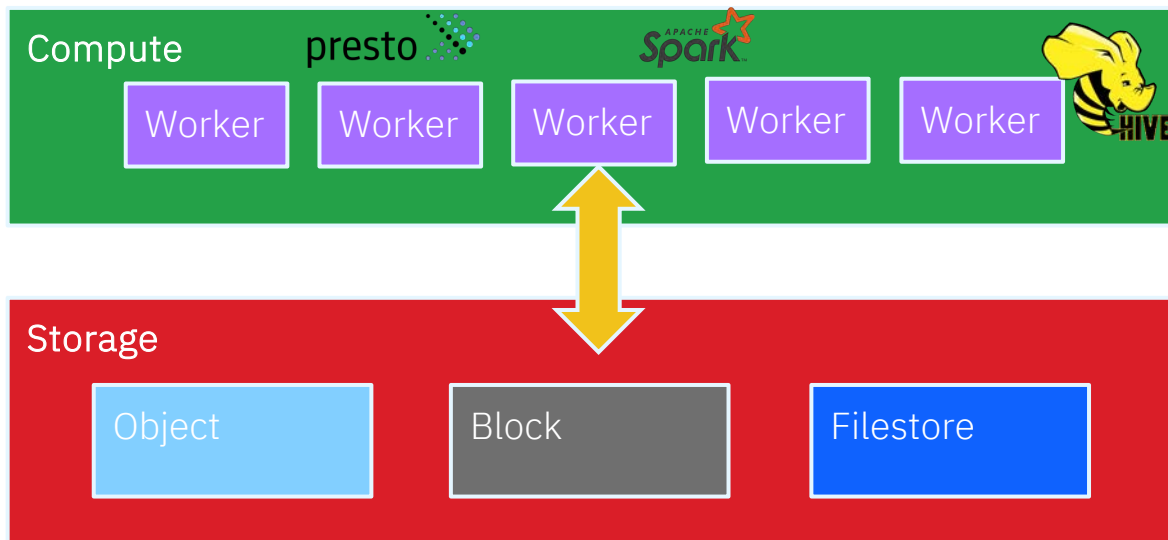
**Storage**

amazon S3

IBM Cloud Object Storage

Hyperscaler Provider (AWS, Azure or IBM)

# Which cluster type?

Depending on your use-case!

There are three models you should consider for your cluster:

- Compute

- Balanced

- Storage

# Which storage?

Recommendations for Storage Decision