

VDAIT



Data Lake vs. Data Warehouse vs. Lakehouse – What's the Right Choice?

Padma Addanki



VDAIT

13653 Leland Road
Centreville, Virginia 20120
571-549-0973

© Veda Infotech Services, LLC. All Rights Reserved. www.vdait.com.

This document is a resource of Veda Infotech Services, LLC (VDAIT.com). The information contained in this document is business sensitive and the intellectual property of VDAIT. This material may not be duplicated, published, disclosed or distributed, in whole or in part, without the express prior written consent of VDAIT. Possession, use, forwarding, or sharing of this document or any part of this document by anyone other than VDAIT RSM Federal constitutes theft and plagiarism. VDAIT Employees and Clients may utilize this resource but remain legally bound to the copyright, trademarks, and confidentiality of this paragraph and may not share this resource with anyone outside their organization. All graphics, tables, content, processes, and methodology are copyrighted and the intellectual property of VDAIT. All rights reserved.

Data Lake vs. Data Warehouse vs. Lakehouse – What's the Right Choice?

In today's data-driven world, understanding the distinctions between data warehouses, data lakes, and the newer concept of data lakehouses is crucial for businesses aiming to manage and analyze their data effectively. But before diving into the differences, it is important to understand why these architectures exist and how they cater to different data needs.

Data Warehouses: The Traditional Storage Solution

Think of a data warehouse as a well-organized library. It's designed to store structured data—information neatly arranged in tables, much like books categorized by genre and author. Before data enters this system, it's cleaned and processed, ensuring consistency and reliability. This setup is ideal for generating reports and conducting business analyses.

In a **Data Warehouse**, data from various operational databases, CRM systems, ERP platforms, and external sources is first extracted using ETL (Extract, Transform, Load) or ELT (Extract, Load, Transform) pipelines. During this process, data is cleaned, transformed, and structured into predefined schemas before being loaded into relational storage systems such as Amazon Redshift, Google BigQuery, Snowflake, or Microsoft Azure Synapse. The stored data is optimized for analytical workloads using indexing, partitioning, and columnar storage formats. Since Data Warehouses follow a **schema-on-write** approach, data is structured upfront to ensure consistency, making it efficient for BI tools, SQL-based queries, and dashboard visualizations. The processed and structured data is then accessed through **SQL engines, OLAP cubes, and business intelligence tools like Tableau, Power BI, or Looker**, enabling organizations to perform advanced analytics, generate reports, and drive data-driven decisions.

Pros:

- **Structured Organization:** Data is stored in a predefined schema, making it easy to retrieve and analyze.
- **High Performance:** Optimized for complex queries, ensuring swift data retrieval.

Cons:

- **Limited Flexibility:** Not well-suited for unstructured data like videos or social media posts.
- **Costly Scaling:** Expanding storage can be expensive due to the need for high-performance hardware.

Data Lakes: The Vast Reservoirs

Now, picture a data lake as a massive reservoir where water from various sources—rivers, rain, streams—collects. Similarly, a data lake stores raw data in its native format, whether it's structured, semi-structured, or unstructured. This means you can dump data from spreadsheets, emails, images, and more into the lake without immediate processing. Which means unlike traditional databases, where data is formatted before storage, a Data Lake follows a schema-on-read approach, meaning that the data is stored in its original format and structured when accessed.

At its core, data from multiple sources such as IoT devices, relational databases, application APIs, or social media streams is ingested in batches or in real time through streaming technologies like Apache Kafka. This raw data is stored in massive datasets in formats like Parquet, ORC, Avro, JSON, or CSV. They are stored in scalable distributed file systems such as Amazon S3, Azure Data Lake Storage (ADLS), Google Cloud Storage (GCS), or Hadoop Distributed File System (HDFS). Since the data is stored in raw format, processing tools such as Apache

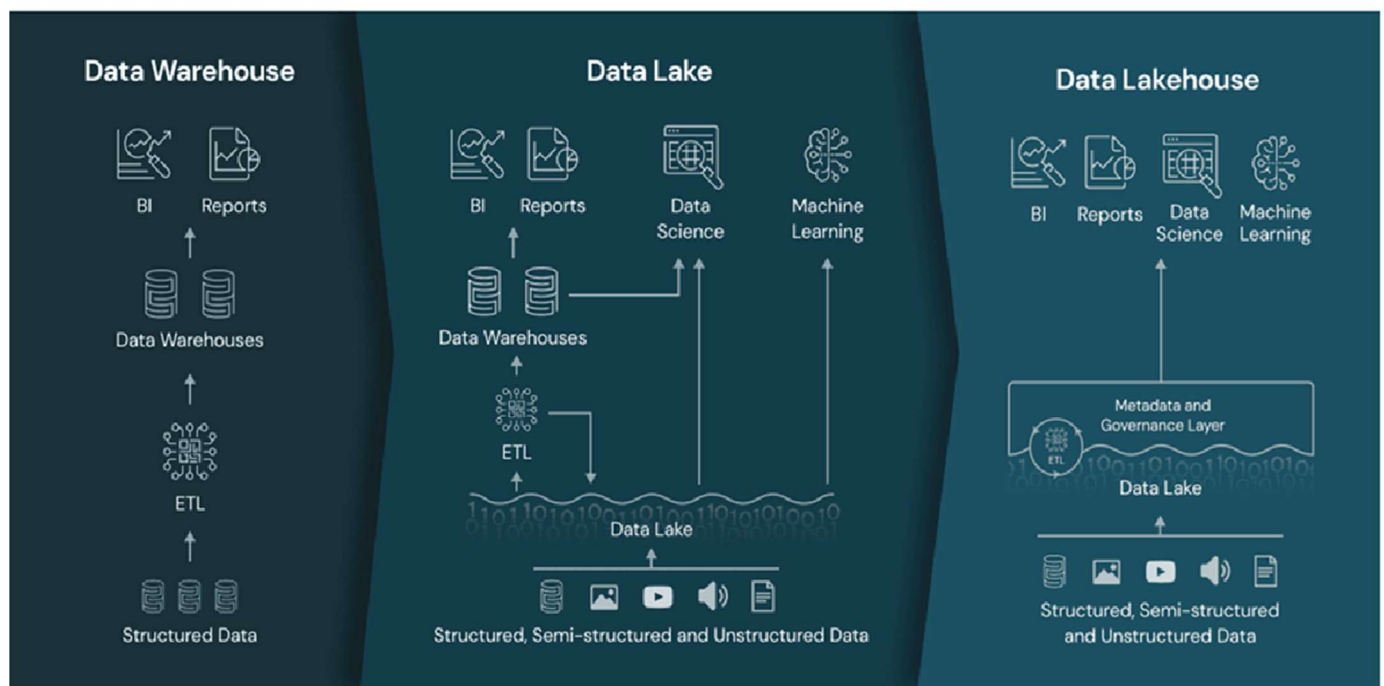
Spark, Hive, or Presto are needed to transform, cleanse and aggregate before analyzing and querying. The analyzed data is then accessed and presented via analytics tools, SQL engines, and machine learning models.

Pros:

- **Flexibility:** Capable of storing diverse data types without a predefined schema.
- **Scalability:** Easily accommodates large volumes of data, making it cost-effective for big data storage.

Cons:

- **Complex Data Management:** Without proper governance, data lakes can become disorganized, leading to so-called "data swamps".
- **Performance Issues:** Retrieving specific data can be slower due to the lack of structure.



Data Storage Architectures: Which One is Right for Your Business?

source: <https://www.databricks.com/glossary/data-lakehouse>

Data Lakehouses: The Best of Both Worlds

Enter the data lakehouse, a hybrid approach that combines the strengths of data warehouses and data lakes. Imagine a modern library that not only has organized shelves for books but also spaces to store raw manuscripts, videos, and digital media. A data lakehouse allows for the storage of raw data while also providing the structure needed for efficient querying and analysis.

At its core, a Data Lakehouse stores data in **low-cost, scalable cloud storage** like **Amazon S3, Azure Data Lake Storage (ADLS), or Google Cloud Storage (GCS)**, similar to a Data Lake. Data is ingested from multiple sources, including **IoT devices, databases, APIs, and real-time streams**, and is stored in formats like **Parquet, ORC, and Avro**. Unlike a traditional Data Lake, which lacks structure, a Data Lakehouse introduces **schema enforcement and transactional capabilities** using open table formats like **Delta Lake, Apache Iceberg, or Apache Hudi**.

To process and analyze data, a Lakehouse supports **SQL-based querying, real-time analytics, and batch processing** using engines like **Apache Spark, Presto, and Hive**. It also integrates seamlessly with **BI tools (Tableau, Power BI) and AI/ML frameworks (PyTorch, Databricks MLflow)**, allowing data scientists and analysts to run complex workloads efficiently. Unlike a Data Warehouse, where data must be **structured and pre-processed before ingestion**, a Lakehouse **preserves raw data while enabling structured querying**, offering both **schema-on-read (like a Data Lake) and schema-on-write (like a Data Warehouse)**. Governance, security, and metadata management are handled using **tools like AWS Glue, Apache Atlas, and Unity Catalog**, ensuring data consistency and access control.

A Data Lakehouse provides the **cost-efficiency and scale of a Data Lake** with the **performance and reliability of a Data Warehouse**, making it a versatile solution for modern **data-driven enterprises** looking to streamline **analytics, AI, and decision-making** from a single, unified platform.

Pros:

- **Unified Storage:** Handles both structured and unstructured data seamlessly.
- **Cost Efficiency:** Offers the scalability of data lakes with the analytical capabilities of data warehouses.

Cons:

- **Emerging Technology:** Being relatively new, best practices are still evolving.
- **Complex Implementation:** Integrating the features of both systems can be challenging.

The following table presents a side-by-side comparison of these architectures, highlighting their key differences in terms of structure, cost, performance, and practical applications. As far as key differences and considerations apply, Data warehouses use a "schema-on-write" approach, defining the structure before storing data. Data lakes employ "schema-on-read," applying structure when data is accessed. Data lakehouses aim to support both methods.

In use cases perspective, Data warehouses are ideal for operational reporting and business intelligence. Data lakes cater to data scientists working on machine learning and advanced analytics. Data lakehouses strive to serve both purposes:

Aspect	Data Warehouse	Data Lake	Data Lakehouse
Data Type	Stores structured data optimized for analytics. Data is cleaned and processed before loading.	Stores raw, unstructured, semi-structured, and structured data in its native format.	Combines structured, semi-structured, and unstructured data in an optimized format.
Data Processing	ETL (Extract, Transform, Load) - data is transformed before being stored.	ELT (Extract, Load, Transform) - <u>data</u> is stored first and transformed later as needed.	Uses a hybrid approach of ETL and ELT, enabling both structured and raw data storage.
Schema	Schema-on-write - <u>Data</u> must conform to a predefined structure before entering the warehouse.	Schema-on-read - <u>Data</u> is stored as-is and structured when queried.	Schema-on-write and schema-on-read - offers flexibility for structured and ad-hoc querying.
Storage Cost	Higher storage costs due to structured and optimized data formats.	Lower storage costs as it stores raw data without heavy preprocessing.	Moderate storage costs - balances performance and affordability.
Query Performance	Fast query performance with optimized indexes and predefined schemas.	Slower queries due to the need for on-the-fly schema creation and transformations.	Optimized query performance with support for both batch and real-time analytics.
Governance & Security	Strong governance and security controls, ensuring compliance with regulations.	Less governance; security measures can be complex due to the varied nature of data.	Improved governance and security, integrating DW controls with Data Lake flexibility.
Flexibility	Less flexible; ideal for structured business intelligence reporting.	Highly flexible, accommodating various data formats and future use cases.	Balances flexibility and performance, supporting diverse workloads.
Use Cases	Business intelligence, operational reporting, and predefined analytics.	AI/ML, real-time analytics, raw data storage, and data discovery.	BI, AI/ML, real-time analytics, and large-scale data applications.
Data Update Frequency	Batch processing - data is updated at scheduled intervals.	Real-time and batch processing - continuously ingests new data.	Supports real-time and batch updates efficiently.
AI & ML Readiness	Not well-suited for AI/ML as it lacks support for unstructured and raw data.	Ideal for AI/ML since it retains all raw data, allowing deep learning and big data analytics.	Well-suited for AI/ML as it provides both raw and processed data in an optimized manner.
Best For	Best for structured, business-oriented data with high query performance needs.	Best for unstructured, large-scale, and AI-driven workloads.	Best for hybrid use cases that need both structured reporting and AI-driven analytics.
Limitations	Expensive, rigid structure, and not ideal for handling large unstructured datasets.	Complexity in managing metadata, query speed issues, and security challenges.	More complex setup, evolving technology, and potential integration challenges.

Conclusion

Choosing the right data architecture—**Data Warehouse**, **Data Lake**, or **Data Lakehouse**—depends on an organization's unique **data needs, use cases, and analytics requirements**. While **Data Warehouses** offer high-performance querying and structured reporting, they lack flexibility for handling diverse data formats. **Data Lakes**, on the other hand, provide unmatched scalability for raw and unstructured data but require additional governance and processing power. The **Data Lakehouse** emerges as a modern solution, blending the best of both worlds by combining the cost efficiency and flexibility of a Data Lake with the governance, structure, and performance of a Data Warehouse.

Organizations looking for **real-time analytics, AI-driven insights, and large-scale data processing** may find the **Lakehouse model** the most effective. However, businesses with well-defined structured reporting needs might still prefer the **traditional Data Warehouse**, while those focusing on **big data exploration, machine learning, and diverse data ingestion** may continue to rely on **Data Lakes**. As **cloud computing, AI, and real-time analytics** evolve, so too will the way we store, manage, and analyze data.

Hashtags: #DataWarehouse #DataLake #DataLakehouse #BigData #CloudComputing #Analytics #AI #MachineLearning #DataEngineering #DataScience #BusinessIntelligence #DataManagement #ETL #DataGovernance #DataArchitecture #CloudData #ModernDataStack #DigitalTransformation

Further Reading & Resources

The below resources offer in-depth insights into the distinctions and convergence between data warehouses, data lakes, and data lakehouses.

Fundamentals of Data Engineering—An O'Reilley book by Joe Reis & Matt Housley

Google Cloud: Data Lake and Data Warehouse Convergence

<https://cloud.google.com/blog/products/data-analytics/data-lake-and-data-warehouse-convergence>

Data Lakehouse <https://www.databricks.com/glossary/data-lakehouse>

Databricks: Data Lakehouse Architecture

<https://www.databricks.com/product/data-lakehouse>

Microsoft Azure: What is a Data Lakehouse?

<https://learn.microsoft.com/en-us/azure/databricks/lakehouse/>

Qlik: Data Lake vs Data Warehouse: 6 Key Differences

<https://www.qlik.com/us/data-lake/data-lake-vs-data-warehouse>

AWS: Data Lake vs Data Warehouse vs Data Mart - Key Differences

<https://aws.amazon.com/compare/the-difference-between-a-data-warehouse-data-lake-and-data-mart/>