
| RESEARCH ARTICLE

Lakehouse Architecture: Unifying Data Warehousing and Advanced Analytics in the Cloud Era

Poshan Kumar Reddy Ponnareddy

S V University, India

Corresponding Author: Poshan Kumar Reddy Ponnareddy, **E-mail:** poshan@gmail.com

| ABSTRACT

The issue of the existence of the chronic separation between data lakes and data warehouses in enterprise data management remains a thorn in its flesh, creating operational headaches and performance issues. Lakehouse architecture offers something different: a unified approach that brings warehouse capabilities to open file formats sitting in cloud object storage. The magic happens through clever metadata management, smart concurrency controls, and thoughtful data organization. What emerges? ACID transactions and query speeds that rival expensive proprietary systems, yet the flexibility and affordability of data lakes remain intact. This matters because modern analytics faces real problems—syncing data between systems takes too long, machine learning pipelines suffer from training-serving mismatches, and specialized tools multiply like rabbits. Lakehouse cuts through this mess. Standard format accessibility means data scientists stop fighting with the infrastructure. Meanwhile, coordinated caching and intelligent partitioning keep business intelligence queries humming along. Feature stores and experiment tracking? Those can live right inside the Lakehouse, trimming the bloated ML technology stack. The shift toward unified platforms isn't just trendy—it actually reduces complexity, speeds up analytics, and lets organizations make data-driven decisions without all the friction.

| KEYWORDS

Lakehouse Architecture, Data Warehousing, Cloud Object Storage, Machine Learning Infrastructure, Acid Transactions.

| ARTICLE INFORMATION

ACCEPTED: 01 November 2025

PUBLISHED: 26 November 2025

DOI: 10.32996/jcsts.2025.7.12.21

1. Introduction

Data management in enterprises has always walked a tightrope between getting things done efficiently and staying flexible enough for different analyses. The most visible symptom? Lakes and warehouses live in separate worlds. This split forces organizations into uncomfortable positions when trying to extract insights from their data. Modern analysis demands constant juggling—moving data around, transforming formats, reconciling differences across systems. All this burns through human hours and computing power at an alarming rate. Managing data quality and keeping governance consistent across these disconnected platforms requires substantial frameworks and oversight [1]. Coordination between the lakes and warehouses turns into a nightmare that is repeated every time, and every time, there is a new chance that something may go wrong or data may get off course. The existence of these parallel systems prevents organizations from adapting quickly because analysts and data scientists have to alternate between interfaces and discover version conflicts, and construct complex pipelines to transport data between storage optimized for completely different end goals.

The costs of running dual architectures go far beyond just operational headaches—they hit performance hard, especially for newer computational tasks. Looking at how AI applications actually perform end-to-end in distributed environments reveals something troubling: data management overhead dominates the picture. Studies tracking real AI workloads show that loading, preprocessing, and moving data often takes longer than the actual computation—the training and inference that supposedly matter most. Numbers tell the story: data-related tasks can eat up seventy percent of total runtime, with much of that waste

Copyright: © 2025 the Author(s). This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) 4.0 license (<https://creativecommons.org/licenses/by/4.0/>). Published by Al-Kindi Centre for Research and Development, London, United Kingdom.

coming from proprietary warehouse formats that don't play nice with extraction, plus the cost of moving data between storage layers [2]. This becomes particularly painful for machine learning workflows that loop repeatedly through large training sets, and for production systems where every millisecond of latency counts. Traditional two-tier architectures don't just cause inconvenience—they create genuine bottlenecks that limit what advanced analytics and machine learning can accomplish.

Lakehouse architecture challenges this status quo by rethinking how storage and compute should relate in enterprise platforms. Rather than maintaining separate systems, it runs full warehouse functionality directly on top of open file formats stored in cloud object storage. The technical approach combines sophisticated metadata tracking, optimistic concurrency controls for transactions, and advanced query optimization to match the performance of purpose-built warehouses without surrendering the cost benefits and flexibility that make data lakes attractive. This architectural shift tackles both the operational tangles documented in modern data analysis [1] and the performance walls encountered when shuffling data through multiple tiers [2]. The promise? Streamlined workflows and faster execution across everything from traditional business intelligence to cutting-edge machine learning applications.

Aspect	Characteristics	Impact
Data Quality Management	Requires robust frameworks across heterogeneous platforms	Constrains organizational agility
Synchronization Overhead	Multiple interfaces and conflicting versions	Introduces failure points
AI Application Performance	Data operations dominate system latency	Constitutes the primary bottleneck
ML Workflow Efficiency	Data loading exceeds computational phases	Affects iterative training cycles

Table 1: Data Analysis Challenges and AI Performance Characteristics [1,2]

2. The Architectural Limitations of Traditional Data Platforms

The majority of enterprise data architectures are divided into tasks: raw and semi-structured data is handled on lakes, and polished and structured data is stored in warehouses designed to handle business intelligence queries. This division stems from fundamentally different design priorities, rooted in how distributed data processing evolved. Data lakes arose when traditional relational databases couldn't handle massive amounts of unstructured and semi-structured content. Apache Hadoop provided the breakthrough—enabling organizations to store petabyte-scale datasets affordably across clusters of commodity hardware. The Hadoop ecosystem bundled the Hadoop Distributed File System for storage with MapReduce for distributed processing, transforming enterprise data management by delivering fault-tolerant storage and parallel processing that scaled horizontally across hundreds or thousands of machines [3]. Such systems focused more on flexibility and low cost so they could accommodate time-series data, text, images, audio, and semi-structured data that defines high volumes of enterprise streams. Warehouses used the converse approach, maximising query performance and data consistency using proprietary storage formats with compute engines closely coupled. Columnar designs, intensive compression, and advanced query optimization provided response times of less than a second to complex analytical queries.

Running this dual architecture exacts a heavy toll across multiple dimensions. Engineering teams constantly extract, transform, and load data from lakes into warehouses, with each transfer representing another potential failure or inconsistency. Multiple data copies multiply storage bills while creating governance nightmares—ensuring security policies and access controls stay consistent across platforms becomes increasingly difficult. Machine learning applications face particular friction: streaming large datasets from proprietary warehouse formats proves so inefficient that data scientists resort to maintaining yet another copy in formats suited for model training. Research on deploying ML systems in production highlights data validation as a critical yet frequently underestimated infrastructure component. Evidence shows that training-serving skew—where features computed during training differ from what's available during inference—ranks as a primary cause of model degradation. Studies pin over eighty percent of production ML failures on data quality problems rather than algorithmic issues, with trouble stemming from schema changes, missing feature values, distribution shifts, and mismatches between training and serving data [4]. Separate lakes and warehouses amplify these challenges because ETL transformations inject additional opportunities for subtle inconsistencies that surface as training-serving skew.

Cloud-native data warehouses achieved commercial success by decoupling storage from compute, yet these systems remain secondary stores in most large organizations—the bulk of enterprise data still lives in lakes. Cloud warehouses added external

table querying for lake formats, but can't extend management features like ACID transaction support to data outside their proprietary storage. This limitation perpetuates exactly the architectural complexity organizations want to escape.

Component	Technology Foundation	Challenge
Data Lake Infrastructure	Hadoop Distributed File System	Horizontal scalability across nodes
Processing Framework	MapReduce distributed computing	Fault-tolerant parallel processing
ML Data Validation	Training-serving consistency	Schema changes and distribution shifts
Production Failures	Data quality deficiencies	Feature value inconsistencies

Table 2: Distributed Data Processing and ML Production Issues [3,4]

3. Technical Foundations of the Lakehouse Paradigm

Lakehouse architecture builds on several key innovations that bring warehouse-like capabilities to open file formats in cloud object storage, fundamentally reimagining transactional guarantees and query optimization without proprietary engines. The foundation rests on sophisticated metadata management implemented through transaction logs that provide ACID guarantees over collections of immutable Parquet files in cloud stores. Systems like Delta Lake achieve these guarantees through an append-only transaction log recording every change to table state—file additions, removals, metadata modifications—with each transaction receiving a monotonically increasing version number establishing total ordering. The transaction log leverages optimistic concurrency control, where writes proceed independently and conflict resolution happens at commit time through atomic rename operations that cloud object stores support. Performance evaluations show Delta Lake hitting transaction throughput exceeding ten thousand commits per second for small transactions while maintaining snapshot isolation semantics, with reads benefiting from time-travel capabilities enabling queries against any historical table version by referencing the corresponding transaction log entry [5]. The metadata layer tracks comprehensive statistics for each data file—row counts, column-level minimums and maximums, null counts, distinct value estimates—enabling sophisticated query optimization that exploits data organization without auxiliary index structures that would compromise the open format requirement.

Query optimization in Lakehouse systems exploits data layout strategies that dramatically boost analytical query performance through intelligent physical record organization within files. Skipping-oriented partitioning employs multi-dimensional clustering algorithms that co-locate records with similar attribute values, letting query engines eliminate irrelevant files based on predicate evaluation against file-level statistics. Research on optimal data layouts for analytical workloads demonstrates that carefully designed partitioning schemes can slash data scanning requirements by factors of ten to one hundred for queries with selective predicates, though effectiveness varies based on distribution characteristics and query patterns. Advanced partitioning approaches harness space-filling curves, particularly Z-order curves derived from Morton encoding, to map multi-dimensional data points onto one-dimensional sequences while preserving locality properties across multiple dimensions simultaneously. Empirical evaluations across diverse datasets reveal Z-order clustering achieves query performance improvements ranging from two to eight times compared to single-dimensional partitioning strategies, with especially pronounced benefits for queries involving conjunctive predicates across multiple columns where traditional partitioning fails to provide effective pruning [6]. The effectiveness stems from creating data files with highly concentrated value ranges for frequently queried attributes, maximizing the probability that query predicates will completely exclude entire files during query planning.

Integrating these foundational techniques—transactional metadata management and intelligent data layout optimization—enables Lakehouse systems to deliver query performance approaching proprietary data warehouses while maintaining complete openness and accessibility of underlying data files stored in standard Parquet format.

Technical Feature	Implementation Mechanism	Capability
Transaction Log	Append-only monotonic versioning	ACID guarantees over object stores
Concurrency Control	Optimistic conflict resolution	Atomic commit operations
Metadata Statistics	File-level aggregate tracking	Query optimization without indexes
Partitioning Strategy	Space-filling curve clustering	Multi-dimensional locality preservation
Data Layout	Z-order Morton encoding	Predicate-based file elimination

Table 3: Lakehouse Transaction Management and Partitioning Strategies [5,6]

4. Performance Characteristics and Optimization Strategies

The balance between the competing demands of the different workloads in Lakehouse systems is a diligent process that is achieved through advanced caching hierarchies and smart data placement techniques that take advantage of the nature of the workload. For analytical queries characteristic of business intelligence workloads, Lakehouse architectures achieve performance parity with traditional warehouses through aggressive caching of frequently accessed data in high-performance storage tiers and sophisticated data layout optimization for cold data residing in object storage. Coordinated memory caching strategies represent a critical optimization for distributed data processing frameworks executing parallel jobs across clusters of machines, where traditional independent caching approaches suffer from redundant data replication and poor cache utilization. Advanced caching coordination mechanisms track file access patterns across distributed worker nodes and implement unified cache management policies that eliminate redundant copies while maximizing aggregate cache hit rates. Research examining coordinated caching in production MapReduce clusters processing multi-terabyte datasets demonstrates that coordination-based approaches achieve hit rates exceeding ninety percent for iterative analytical workloads, compared to sixty to seventy percent for uncoordinated caching schemes, with the performance differential becoming increasingly pronounced as cluster sizes scale beyond one hundred nodes [7]. The cache coordination system tracks access patterns as they occur, and dynamically changes the replication factors depending on file popularity, making sure that popular data blocks in the files have an adequate number of replicas in the entire cluster and that infrequently used blocks do not use much cache capacity. Benchmark measurements show that optimized query engines on Lakehouse storage are able to run complex analytical queries in the same latency as workloads on cloud-native warehouses, especially where the workload has predictable access patterns that can be utilized to effectively use the cache.

Machine learning and data science workloads benefit significantly from the direct accessibility of data in open formats, with columnar storage architectures providing substantial advantages for feature-intensive analytical processing. The ability to stream data directly from object storage eliminates the inefficiencies associated with extracting large datasets from proprietary warehouse formats. Columnar storage formats, particularly Apache Parquet, organize data by column rather than by row, enabling highly efficient compression ratios and selective column reading patterns that dramatically reduce I/O requirements for analytical queries. Parquet's design achieves typical compression ratios ranging from three to one up to ten to one, depending on data characteristics, with categorical columns and sorted numeric sequences exhibiting particularly favorable compression behavior. The columnar layout enables query engines to read only the specific columns required for computation, reducing data transfer volumes proportionally to column selectivity. For machine learning feature extraction workflows that typically access ten to thirty percent of available columns, this architectural characteristic translates to data transfer reductions of seventy to ninety percent compared to row-oriented formats [8]. The elimination of data copying between systems not only improves training throughput but also simplifies the machine learning lifecycle by reducing the surface area for data consistency issues.

Lakehouse systems have scalability properties due to their cloud-native design, which allows organizations to dynamically add processing capacity to meet workload requirements at commodity cost and virtually unlimited storage capacity.

Optimization Technique	Operational Approach	Performance Benefit
Coordinated Caching	Unified cache management policies	Eliminates redundant replication
Access Pattern Tracking	Real-time monitoring across nodes	Dynamic replication adjustment
Columnar Organization	Column-oriented data layout	Selective reading capability
Compression Strategy	Format-specific algorithms	Reduced storage and I/O
Feature Extraction	Partial column access	Minimized data transfer

Table 4: Caching Optimization and Columnar Storage Benefits [7,8]

5. Implications for Machine Learning and Data Science Workflows

The consolidation of storage frameworks has extreme significance to machine learning functions and data science procedures, fundamentally altering how companies build frameworks to teach and launch production machine learning designs, as well as sustain them. Modern companies have established elaborate ecosystems of ML-specific software, such as versioning systems of data, feature stores, and model registries, which tend to recreate feature availability in database management systems. The machine learning lifecycle encompasses a complex set of interdependent stages, including data preparation, feature engineering, model training, hyperparameter optimization, model evaluation, deployment, and monitoring, with each stage generating artifacts that must be tracked and versioned for reproducibility. Organizations typically execute thousands of experimental training runs annually, with each experiment potentially varying in dataset composition, feature engineering logic, algorithm selection, and hyperparameter configurations. The challenge of managing this experimentation complexity has driven the

development of specialized ML lifecycle management platforms that provide systematic tracking of experiments, parameters, metrics, and model artifacts. Empirical evidence from production ML systems reveals that data scientists commonly explore dozens to hundreds of model variations during development, with typical projects involving fifty to two hundred distinct experimental runs before identifying production-ready candidates [9]. The Lakehouse paradigm suggests that many of these specialized systems could be simplified or eliminated by leveraging built-in data management capabilities, as the combination of ACID transactions, time-travel queries, and rich metadata tracking provides foundational primitives necessary for reproducible ML workflows without requiring separate infrastructure layers.

Feature store functionality, which manages the storage and serving of ML features across training and inference contexts, represents a prime candidate for consolidation within Lakehouse architectures. Feature stores emerged as a response to the pervasive challenge of training-serving skew, wherein discrepancies between feature computation logic in training and production environments lead to model performance degradation and operational failures. The architectural complexity of maintaining consistency stems from the fundamental difference between offline batch processing used for training and online real-time processing required for inference. Organizations without feature store infrastructure typically reimplement feature engineering logic separately for training pipelines and serving systems, with these parallel implementations inevitably diverging over time as engineers make incremental modifications. Research examining feature store adoption patterns documents that production ML systems typically maintain between fifty and five hundred distinct features, with feature engineering pipelines comprising ten to thirty transformation stages that aggregate data from multiple source tables, apply complex business logic, and compute temporal aggregations over historical windows [10]. By implementing feature computations as declarative queries over versioned data, organizations can leverage native snapshot isolation and time-travel capabilities to ensure consistency between training and production environments. The ability to query historical snapshots of data directly supports model retraining and experimentation workflows without requiring separate versioning infrastructure, as data scientists can reference specific transaction log versions to reconstruct exact training datasets used for historical model builds.

The fact that the declarative machine learning systems could be effective against Lakehouse storage creates additional opportunities in the form of optimization, and the immediate access to the data allows forming the new patterns of distributed collaboration, including the model of data mesh organization.

6. Conclusion

Lakehouse architecture is an important breakthrough in the data platform design of enterprises, resolving the old conflicts between analytical flexibility and operational efficiency with architectural cohesion. Conventional two-tier architectures consisting of distinct data lakes and warehouses come with heavy operational costs, including constant maintenance of ETL pipelines, the cost of data duplication, the complexity of governance, and performance bottlenecks to machine learning workloads. The Lakehouse paradigm removes such challenges by providing extensive warehouse functionality on top of open file formats in cloud object storage, taking advantage of intelligent metadata management, optimistic concurrency control, and intelligent data layout optimization to provide competitive performance without compromising openness or flexibility. Transaction log support removes ACID guarantees and time-travelling support on regular Parquet files, and coordinated caching and skipping-oriented partitioning provide query performance akin to proprietary warehouses. The immediate availability of data in open formats is especially useful in machine learning processes, where architectural jostling on data extraction of proprietary systems is removed and streaming access patterns to support iterative model training are made possible. Lakehouse infrastructure can store feature stores, as well as ML lifecycle management functions, to minimize the spread of specialized tooling and ease the technology stack. It seems as though the march to integrated data platforms is destined due to the sheer volumes of already existing data in data lakes and the operational benefits of simplification in architecture. Moving to Lakehouse architectures allows organizations to fulfill a wide range of analytical workloads, such as conventional business intelligence, complex machine learning, on a unified storage tier, allowing easier, more responsive decision-making on data and cutting down on operational overhead and complexity. Further development of cloud-native query engines, declarative ML frameworks, and distributed collaboration patterns will continue to expand the functions and use of Lakehouse systems in the enterprise setting.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors, and the reviewers

References

- [1] Dani P, (2024) Apache Parquet for Data Engineers: Optimized Data Storage, Estuary, 2024. [Online]. Available: <https://estuary.dev/blog/apache-parquet-for-data-engineers/>

-
- [2] Daniel R, et al., (2020) Missing the forest for the trees: End-to-end AI application performance in edge data centers, 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/9065599>
- [3] Databricks, (n.d) What is Hadoop? [Online]. Available: <https://www.databricks.com/glossary/hadoop>
- [4] Eric B, et al., (2019) DATA VALIDATION FOR MACHINE LEARNING, MLSys, 2019. [Online]. Available: <https://mlsys.org/Conferences/2019/doc/2019/167.pdf>
- [5] Ganesh A, et al., (2012) PACMan: coordinated memory caching for parallel jobs, 2012. [Online]. Available: <https://dl.acm.org/doi/10.5555/2228298.2228326>
- [6] Kim H and Jim D, (2023) Feature Store: The missing data layer for Machine Learning pipelines? Hopsworks, 2023. [Online]. Available: <https://www.hopsworks.ai/post/feature-store-the-missing-data-layer-in-ml-pipelines>
- [7] Liwen S, et al., (2016) Skipping-oriented partitioning for columnar layouts, ACM Digital Library, 2016. [Online]. Available: <https://dl.acm.org/doi/10.14778/3025111.3025123>
- [8] Matei Z, et al., (2018) Accelerating the machine learning lifecycle with MLflow, People@ EECS. 2018. [Online]. Available: https://people.eecs.berkeley.edu/~matei/papers/2018/ieee_mlflow.pdf
- [9] Michael A, et al., (n.d) Delta Lake: High-performance ACID table storage over cloud object stores, VLDB,[Online]. Available: <https://www.vldb.org/pvldb/vol13/p3411-armbrust.pdf>
- [10] Qualtrics, (n.d) How to analyse survey data: Best practices, tips and tools, [Online]. Available: <https://www.qualtrics.com/en-au/experience-management/research/analysis-reporting/>