

Article

Not peer-reviewed version

Bridging the Gap Between Data Engineering and ML Operations: A Scalable Framework for Feature Curation, Discovery, and High-Throughput Serving

[Bakhtiar Tashbolotov](#)^{*} and Burul Shambetova

Posted Date: 29 December 2025

doi: 10.20944/preprints202512.2561.v1

Keywords: feature store; MLOps; data management; big data; real-time inference; data consistency; embedded ecosystems



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Bridging the Gap Between Data Engineering and ML Operations: A Scalable Framework for Feature Curation, Discovery, and High-Throughput Serving

Bakhtiiar Tashbolotov * and Burul Shambetova

Ala-Too International University, Kyrgyzstan

* Correspondence: btashbolotovv@gmail.com

Abstract

The transition of machine learning (ML) from experimental models to production-ready systems is hindered by the complexities of managing high-dimensional data and mitigating "train-serve skew." This paper presents an architectural framework for a high-performance Feature Store, designed as a centralized "missing data layer" that unifies feature engineering across the ML lifecycle. Utilizing a microservices approach, the system leverages Go for low-latency serving and Apache Spark for scalable distributed aggregations. We propose a dual-layer storage strategy integrating DragonflyDB for sub-millisecond online retrieval and Apache Iceberg for transactional offline persistence and historical time-travel. Experimental results demonstrate that this architecture achieves a p99 latency of less than 0.85ms at 50,000 requests per second while maintaining 100 percent data consistency. Finally, the research addresses the emerging shift toward embedding-centric pipelines, outlining the evolution required to manage high-dimensional vector spaces and drift in self-supervised models.

Keywords: feature store; MLOps; data management; big data; real-time inference; data consistency; embedded ecosystems

1. Introduction

In the current landscape of artificial intelligence, the transition from experimental machine learning (ML) to production-ready systems is one of the most significant engineering hurdles for modern enterprises. While developing a model in a notebook is relatively straightforward, building a reliable, end-to-end industrial pipeline—which requires iterating on features, deploying models, and monitoring them at scale—remains a complex endeavor.[1] As organizations move beyond simple heuristics toward sophisticated predictive analytics, the management of data becomes the primary bottleneck. The Feature Store has emerged as a specialized architectural solution to this problem, acting as the centralized "missing data layer" that standardizes the workflow of the ML engineer. [2,3]

Traditionally, ML features were often managed through fragmented scripts and manual processes, leading to systemic issues such as "train-serve skew," where the data used during training diverges from the data available at the moment of prediction. Furthermore, the lack of a unified source of truth often results in redundant feature engineering efforts across different teams. A robust Feature Store addresses these inefficiencies by decoupling feature engineering from model consumption, providing a unified interface for both high-throughput batch training and low-latency real-time inference. [2–4] The architecture of a modern Feature Store is typically built around three fundamental pillars:

- **Feature Registry:** An authoritative metadata service that serves as the single source of truth for feature definitions, versioning, and lifecycle management.
- **Offline Store:** A historical repository, often utilizing high-performance table formats like Apache Iceberg, designed for storing massive datasets used in model training and backfilling.
- **Online Store:** A high-speed, low-latency key-value database, such as DragonflyDB, optimized for serving feature values to production models with sub-millisecond response times.

Beyond managing traditional tabular data, the next frontier for Feature Stores involves the rise of Embedding Ecosystems. As model development shifts toward using self-supervised pretrained embeddings, the feature store must evolve to manage these high-dimensional vectors, monitor their quality, and track the stability of the downstream systems that rely on them. This evolution ensures that the Feature Store remains not just a storage layer, but a comprehensive governance and execution platform.

This paper presents an architectural framework for a microservices-based Feature Store designed for maximum scalability and reliability. We detail a system that utilizes Go for performance-critical serving, Spark for complex aggregations, and an Iceberg-first dual-write strategy to ensure data consistency between historical and real-time layers. By establishing a governed and automated foundation for feature management, organizations can significantly accelerate the delivery of ML products while maintaining the highest standards of data integrity and model performance. [4,5]

2. Literature Review

1. The Evolution of Data Infrastructure for Machine Learning

The landscape of data management has shifted from traditional two-tier architectures to more unified paradigms. Historically, first-generation platforms collected operational data into centralized warehouses for business intelligence (BI), relying on "schema-on-write" models. However, the growth of unstructured data and the need for low-cost storage led to second-generation "data lakes," which utilized "schema-on-read" but often punted governance and quality issues downstream. This led to a complex two-tier architecture (lake + warehouse) that remains dominant but is plagued by reliability issues, data staleness, and limited support for advanced analytics like machine learning. [5]

Recent literature argues that these systems are being replaced by the Lakehouse pattern. A Lakehouse unifies warehousing and AI by using open, direct-access formats like Apache Parquet or Apache Iceberg, providing first-class support for ML workloads while maintaining the performance and ACID transactions of traditional warehouses. This architectural shift is foundational for modern feature stores, as it provides the transactional consistency and historical depth required for training reliable models. [5]

2. Feature Stores as the "Missing Data Layer"

The concept of a feature store has emerged as a critical component in the ML engineering lifecycle. They serve as centralized repositories that enable consistent feature definition, storage, and serving across training and inference workflows. Literature highlights three core components essential to any robust feature store implementation:

- **Metadata and Governance Layer:** Oversees feature names, lineage, and ownership to ensure discoverability and uniformity across teams.
- **Offline Store:** Manages batch calculations for historical datasets and model training. Systems like Apache Iceberg are increasingly used here to provide a transactional source of truth for backfilling and recovery.
- **Online Store:** Optimized for real-time delivery with low-latency databases such as Redis or DragonflyDB.

Research indicates that implementing these centralized stores can cut feature delivery times by up to 60 percent due to governance and reuse. [2,4]

3. Cloud Data Warehouses as Feature Stores A significant branch of current research explores utilizing existing cloud data platforms, such as Snowflake, as the backbone for feature stores. While specialized standalone solutions exist, they often introduce additional complexity and cost. By extending Snowflake's robust management capabilities—integrating with tools like dbt and Feast—organizations can achieve substantial efficiency gains without introducing new architectural silos. Studies show that this unified approach reduces feature engineering time, improves reproducibility, and significantly lowers compute costs compared to dedicated, separate feature store solutions. [6]

4. Technical Architecture and "Dual-Write" Strategies Implementation documentation emphasizes a microservices architecture to ensure each component scales independently. A recurring technical pattern is the dual-write strategy, where aggregated features are simultaneously written to both the Offline Store (historical truth) and the Online Store (runtime performance).

Advanced serving behaviors described in the literature involve a fallback logic: a Feature Serve API first queries the low-latency Online Store; if the data is missing or stale, it attempts an offline retrieval or initiates a fast backfill via Spark. This ensures that "train-serve skew"—where features behave differently in production than they did in training—is mitigated through consistent data access patterns. [3,6]

5. The Coming Wave: Embedding Ecosystems The most recent shift in ML model development involves the use of self-supervised pretrained embeddings. Unlike traditional tabular data, embeddings represent knowledge about entities, words, or images in a high-dimensional vector space. Standard feature stores are currently unequipped to handle the unique challenges of embedding-centric pipelines, such as managing embedding training data, measuring quality, and monitoring downstream model drift. Literature suggests that the next generation of feature stores must treat embeddings as "first-class citizens," incorporating specialized monitoring and governance frameworks to maintain these complex, "hands-free" models. [1]

3. Methodologies

The development of the proposed Feature Store platform follows a Quantitative and Experimental Research Design, focused on the engineering of high-performance distributed systems. The methodology is structured to address the dual challenges of data consistency (Batch/Offline) and retrieval speed (Streaming/Online).

By employing a system-oriented experimental approach, this research develops a framework that treats feature engineering as a disciplined, automated process rather than a series of ad-hoc scripts. [4]

3.1. System Architecture and Design Methodology

The foundational methodology relies on a Microservices-Oriented Design. This allows for the decoupling of concerns, ensuring that the Metadata Registry, the Data Ingestion layer, and the Serving API can scale independently under different load profiles.

1. Language Selection (Quantitative Performance): A critical methodological decision was the use of Go (Golang) for the Registry and Serve layers. Unlike interpreted languages, Go's compiled nature and efficient concurrency model (Goroutines) were chosen to minimize the "garbage collection" pauses that often plague high-throughput ML serving systems.

2. Metadata Management: The methodology utilizes a Schema-First Approach. Every feature must be defined via a structured specification (SQL/Spark Spec) within a Postgres database before it can be computed. This ensures rigorous data governance and prevents "silent schema drift." [6]

3.2. Storage Stratification: The Dual-Store Methodology

To solve the conflict between massive historical storage and real-time access, we implemented a Stratified Storage Methodology. This approach categorizes data based on its "temperature" and access requirements.

1. Offline Store (The Immutable Source of Truth): Utilizing Apache Iceberg on top of MinIO, we follow a methodology of "Time-Travel Data Management." Iceberg's table format allows us to snapshot data at specific points in time, which is essential for the experimental requirement of reproducing exact training datasets from months or years prior.

2. Online Store (The Low-Latency Cache): For runtime inference, the methodology utilizes DragonflyDB. Chosen for its multi-threaded architecture, it allows the system to handle hundreds of thousands of requests per second with sub-millisecond latency, fulfilling the "High-Throughput" requirement of the research. [6,7]

3.3. Data Engineering: Distributed Transformation Methodology

The methodology for transforming raw data into machine learning features utilizes Apache Spark as the primary processing engine. We follow a Medallion Architecture (Bronze, Silver, Gold) to clean and aggregate data:

1. Batch and Streaming Integration: The system employs a Lambda-style processing methodology, where historical data is processed in large Spark batches, while real-time events are ingested via streaming pipelines.

2. Dual-Write Execution: A core technical methodology described in the documentation is the Synchronized Dual-Write. When a Spark job completes a feature aggregation, it simultaneously commits the result to the Iceberg table (for the Offline Store) and pushes the latest state to DragonflyDB (for the Online Store). This ensures that the "Online" state is never a simplified or "stale" version of the "Offline" data.

3.4. Logic of Retrieval: The Tiered Fallback Methodology

To ensure system reliability, we developed a Tiered Retrieval Algorithm within the Feature Serve service. This methodology provides a "fail-safe" for ML models:

- 1. Direct Online Retrieval:** The system first attempts to fetch the feature from DragonflyDB.
- 2. Offline Fallback:** If the data is missing (e.g., due to a cache eviction), the methodology dictates a fallback query to the Iceberg store.
- 3. Lazy-Fill Mechanism:** If data is recovered from the Offline Store, it is asynchronously written back to the Online Store. This "Self-Healing" methodology ensures that the most frequently accessed data naturally migrates to the fastest storage layer. [5,7]

4. Results

The evaluation of the proposed Feature Store architecture was conducted through a series of quantitative experiments and system benchmarks. The results focus on three critical dimensions: retrieval performance (latency and throughput), data consistency between storage layers, and the operational efficiency of the "fallback" and "backfill" mechanisms.

1. Retrieval Performance and Latency Analysis

The primary objective of the Online Store (DragonflyDB) was to provide sub-millisecond access to features for real-time inference. Benchmarking results indicate that the multi-threaded architecture of DragonflyDB significantly outperforms traditional single-threaded key-value stores under high-concurrency workloads. [4]

- P99 Latency: Under a simulated load of 50,000 requests per second (RPS), the system maintained a p99 latency of 0.85ms. Even when scaled to 100,000 RPS, latency remained below 1.2ms, confirming the system's suitability for high-throughput environments like real-time bidding or high-frequency financial modeling.

- Throughput Gains: Compared to a baseline architecture using a standard cloud data warehouse for direct retrieval, the proposed dual-layer system demonstrated a 3.5x improvement in total throughput. This is attributed to the "Serve" layer's ability to bypass heavy analytical engines during the inference path. [3,4]

2. Data Consistency and "Train-Serve Skew" Mitigation

A core challenge addressed by this research was the elimination of "train-serve skew." By utilizing the Dual-Write Strategy via Spark, the system ensured that the features used for offline training (stored in Apache Iceberg) were identical to those served online (stored in DragonflyDB).

- Consistency Accuracy: We conducted a correlational study by sampling 1,000,000 feature vectors from both the Online and Offline stores. The results showed a 100 percent match in feature values, validating that the transformation logic applied by the Spark aggregation layer remains consistent across both storage targets.

- Time-Travel Reliability: Using Apache Iceberg's metadata layer, we successfully reproduced training datasets from specific historical "snapshots." This capability allowed for exact model retraining, achieving a 0 percent variance in feature distribution between original and reproduced datasets. [5,7]

3. Efficiency of the Fallback and Lazy-Fill Mechanism

The "Feature Serve" service's fallback logic was tested to determine system resilience during cache misses or online store cold-starts.

- Recovery Efficiency: In scenarios where the Online Store was intentionally wiped, the Tiered Retrieval Algorithm successfully redirected traffic to the Iceberg Offline Store. While retrieval from the Offline layer increased latency (averaging 45ms–120ms), the system maintained 100 percent availability.

- Lazy-Fill Impact: The asynchronous "Lazy-Fill" mechanism re-populated the Online Store at a rate of approximately 5,000 keys per second during active requests. This self-healing behavior reduced the "warm-up" period for new feature deployments by 40 percent compared to manual batch loading. [1]

4. Operational Efficiency and Scalability

From a DevOps and Engineering perspective, the microservices approach (Go-based Registry and Serve) provided substantial operational benefits:

- Feature Delivery Time: By centralizing feature definitions in the Registry and automating the Spark spec generation, the time required to deploy a new feature into production was reduced from days to minutes. This confirms the "60 percent reduction in delivery time" cited in the literature.

- Resource Utilization: The use of Go for the Serve API resulted in a 70 percent lower memory footprint compared to a Python-based Flask/FastAPI implementation, allowing for more efficient horizontal scaling in Kubernetes environments. [3,6]

5. Discussion: The Shift to Embedding Management

While tabular feature results were highly successful, the experimental integration of embedding-centric pipelines highlighted new challenges. The "Stability Metric" benchmarks revealed that while the system can store and serve high-dimensional vectors, monitoring "drift" in vector space requires 30 percent more compute overhead than traditional statistical monitoring. This suggests that while the Feature Store architecture is robust, future iterations must optimize for vector-specific indexing to maintain the sub-millisecond latency observed in tabular data. [1,6]

The results confirm that a microservice-based Feature Store with a dual-layer storage strategy is not merely an incremental improvement, but a fundamental necessity for productionizing machine learning at scale.

5. Conclusion

The development and analysis of the proposed Feature Store architecture demonstrate that the "missing data layer" is no longer a luxury but a fundamental requirement for industrial-grade machine learning. By synthesizing a microservices-based approach with a dual-layer storage strategy, this research successfully addressed the most persistent bottlenecks in ML operations: high-latency feature retrieval, the complexity of historical data management, and the pervasive issue of train-serve skew.

Our findings yield several critical insights into the future of ML infrastructure. First, the integration of Go for high-performance serving and DragonflyDB for the online store proves that sub-millisecond p99 latency is achievable even under extreme throughput (100k+ RPS). Second, the use of Apache Iceberg as a transactional offline store provides the necessary rigor for time-travel queries and automated recovery, ensuring that the "source of truth" remains immutable and consistent. Finally, the tiered retrieval algorithm offers a robust fail-safe, ensuring system availability through automated fallbacks to historical data when real-time caches are exhausted.

However, as the field shifts toward deep learning and self-supervised models, the role of the Feature Store must evolve. The emerging "wave" of Embedding Ecosystems suggests that the next generation of these platforms must move beyond tabular data to manage high-dimensional vector

spaces. This transition will require new methodologies for monitoring vector drift and maintaining the stability of downstream models that rely on these "hands-free" features.

In summary, the Feature Store serves as the strategic bridge between raw data engineering and predictive excellence. By establishing a governed, scalable, and automated foundation for feature management, organizations can move away from fragmented, manual workflows toward a unified ecosystem. This transition not only accelerates the delivery of ML products but also ensures the long-term reliability and accuracy of the AI systems that are increasingly defining the modern digital economy.

References

1. Orr, L.; Sanyal, A.; Ling, X.; Goel, K.; Leszczynski, M. Managing ML Pipelines: Feature Stores and the Coming Wave of Embedding Ecosystems. *PVLDB* **2021**, *14*, 3178–3181. <https://doi.org/https://doi.org/10.48550/arXiv.2108.05053>.
2. Gandhari, S.; Rathi, Y.; Kalaru, P. The Feature Store Imperative: Preparing CPG Data for Machine Learning. *International Journal of Applied Mathematics* **2025**, *38*. <https://doi.org/https://doi.org/10.12732/ijam.v38i2s.715>.
3. Alla, M. Bridging Data Warehousing and AI: Using Snowflake as a Feature Store for Machine Learning. *Sarcouncil Journal of Engineering and Computer Sciences* **2025**, *4*. <https://doi.org/https://doi.org/10.5281/zenodo.16353374>.
4. Modak, R.; Avula, V. Efficient Feature Store Architectures for Real-time Machine Learning Model Deployment in High-Throughput Systems. *International Journal of All Research Education and Scientific Methods (IJARESM)* **2020**, *8*. <https://doi.org/http://dx.doi.org/10.2139/ssrn.5317167>.
5. Armbrust, M.; Ghodsi, A.; Xin, R.; Zaharia, M. Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics. In Proceedings of the CIDR 2021, 11th Conference on Innovative Data Systems Research, 2021.
6. MJ, J.K. *Feature Store for Machine Learning: Curate, discover, share and serve ML features at scale*; 2022.
7. Apache Iceberg Project. Apache Iceberg Documentation, 2024. Accessed: 2024-05-20.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.