

Real-Time Change Data Capture at Petabyte Scale: Architecting Low-Latency Oracle-to-Cloud Synchronization for Enterprise Retail Analytics

Praveen Kumar Dora

Lawrence Technological University, USA

ARTICLE INFO

Received: 02 Nov 2022

Revised: 18 Dec 2022

Accepted: 28 Dec 2022

ABSTRACT

Enterprises operating at petabyte-scale data volumes face increasing pressure to transition from batch-oriented analytics to near-real-time decision systems. This paper presents a production-validated architecture for real-time Change Data Capture (CDC) from Oracle databases to cloud-based analytical platforms, designed specifically for large-scale retail environments. We detail the challenges inherent in Oracle-based CDC, propose a layered architecture leveraging log-based capture and distributed streaming systems, and evaluate its performance under peak retail workloads. Empirical results demonstrate sub-second median latency and sustained throughput exceeding 2.5 million change events per second, significantly outperforming traditional ETL pipelines.

Keywords: Change Data Capture (CDC), Real-Time Data Pipelines, Oracle Database Replication, Distributed Streaming Systems, Retail Analytics

1. Introduction

The shift toward real-time analytics has fundamentally altered data engineering practices across industries. Nowhere is this transformation more pronounced than in the retail sector, where operational decisions—ranging from dynamic pricing to inventory replenishment and personalized recommendations—are increasingly driven by continuously updated data streams rather than static, historical snapshots. Traditional data architectures, built around periodic batch extraction, transformation, and loading (ETL), are no longer sufficient to support these latency-sensitive use cases. In many large retail organizations, batch pipelines introduce delays measured in hours, creating a temporal disconnect between operational systems and analytical insights that directly impacts revenue, customer experience, and supply chain efficiency.

Despite the growing demand for real-time capabilities, a significant portion of enterprise data remains anchored in legacy relational database systems. Among these, Oracle databases continue to serve as critical systems of record due to their robustness, transactional guarantees, and long-standing adoption in enterprise environments. However, their central role also presents a structural challenge: while these systems excel at transactional consistency, they were not originally designed to support high-throughput, low-latency data replication into modern cloud-native analytics platforms. As a result, organizations must reconcile two competing priorities—preserving the integrity and stability of mission-critical transactional systems while enabling continuous data availability for downstream consumers.

Change Data Capture (CDC) has emerged as a key enabling paradigm in this context. Rather than relying on periodic bulk extraction, CDC focuses on capturing incremental changes—insertions, updates, and deletions—directly from the source database and propagating them to downstream systems in near real time [1,2]. By leveraging database transaction logs, CDC minimizes impact on source workloads while providing a complete and ordered history of data modifications. This makes it

particularly well-suited for modern streaming architectures, where change events can be processed, enriched, and consumed by multiple applications simultaneously.

However, the practical realization of CDC at enterprise scale introduces a set of non-trivial engineering challenges. At petabyte-scale data volumes, the sheer rate of change events—often reaching millions of transactions per second during peak retail periods—places significant demands on capture, transport, and processing layers. Ensuring low end-to-end latency requires careful coordination across distributed components, while maintaining consistency necessitates strict adherence to transactional ordering and fault-tolerant recovery mechanisms. Additionally, operational considerations such as schema evolution, data quality assurance, and infrastructure elasticity further complicate the design of production-ready CDC pipelines.

The complexity is amplified in Oracle-based environments, where access to transaction logs, handling of System Change Numbers (SCNs), and limitations of native tooling impose additional constraints on scalability and observability. When combined with the need to integrate with cloud-based data lakes and analytical warehouses, these factors create a challenging design space that extends beyond conventional data integration approaches.

This paper addresses these challenges by presenting a production-grade CDC architecture developed and validated within a large-scale retail enterprise. The proposed solution focuses on log-based change capture from Oracle databases and its integration with distributed streaming and cloud storage systems to achieve low-latency, high-throughput data synchronization. By grounding the discussion in real-world deployment constraints and performance characteristics, the paper aims to contribute both practical insights and architectural patterns for organizations seeking to modernize their data infrastructure while retaining existing transactional systems.

2. State of Change Data Capture

Change Data Capture (CDC) has undergone substantial evolution as enterprises have transitioned from batch-oriented data processing toward continuous, event-driven architectures. Early data integration systems relied heavily on periodic extraction techniques, but the growing demand for freshness and responsiveness led to the emergence of CDC as a foundational capability for modern data platforms. Over time, three primary CDC methodologies have come to dominate both academic discussion and industrial practice, each with distinct trade-offs in terms of performance, reliability, and operational complexity.

Trigger-based CDC represents one of the earliest approaches, wherein database triggers are defined on source tables to capture insert, update, and delete operations. While conceptually straightforward and relatively easy to implement, this method introduces significant overhead on transactional systems. Every data modification incurs additional write operations to logging tables, leading to write amplification and increased contention on heavily accessed tables. In high-throughput environments such as retail transaction systems, this overhead can degrade primary database performance and increase latency for end-user operations [3].

Timestamp-based CDC attempts to reduce this overhead by periodically querying tables for rows modified since the last extraction time, typically using a “last updated” timestamp column. Although less intrusive than triggers, this approach is inherently limited in its ability to guarantee correctness. It is susceptible to missed updates due to clock skew, delayed transactions, or non-monotonic timestamp behavior. Furthermore, as data volumes grow, repeated full or partial table scans become increasingly expensive, making this approach difficult to scale for large datasets with high update frequency [4].

Log-based CDC, in contrast, operates by directly reading database transaction logs—such as Oracle redo logs—to capture changes as they occur. This method provides a complete and ordered record of

all data modifications with minimal impact on the source system, as it leverages existing logging mechanisms rather than introducing additional workload. Log-based CDC also preserves transactional semantics, including commit order and atomicity, which are critical for maintaining consistency in downstream systems. For these reasons, it has emerged as the preferred approach for enterprise-scale data replication and streaming architectures [5,6].

The increasing adoption of log-based CDC has been supported by a growing ecosystem of tools and platforms. Open-source solutions such as Debezium, along with commercial offerings like Oracle GoldenGate and managed services such as AWS Database Migration Service, have made CDC more accessible to organizations seeking to modernize their data pipelines. These tools provide abstractions for log parsing, event generation, and integration with streaming systems, reducing the implementation burden on engineering teams [7,8].

Table 1: CDC Approach Comparison

Approach	Performance	Scalability	Consistency	Overhead
Trigger-based	Low	Poor	High	High
Timestamp-based	Medium	Limited	Weak	Medium
Log-based (Proposed)	High	High	Strong	Low

However, despite these advancements, significant challenges remain when deploying CDC in large-scale, production environments. Vendor-provided tools often abstract away complexity at the cost of limited flexibility and visibility, which becomes problematic under extreme workloads. In particular, ensuring strict ordering guarantees across distributed systems, handling frequent schema evolution without disrupting downstream consumers, and maintaining robust failure recovery mechanisms continue to present difficulties. These challenges are especially pronounced in retail contexts, where data volumes, velocity, and variability exceed the assumptions underlying many off-the-shelf solutions.

As a result, a gap persists between the capabilities of existing CDC tools and the operational requirements of enterprise-scale retail systems. Bridging this gap requires not only careful architectural design but also a deeper integration of CDC mechanisms with distributed streaming, storage, and processing frameworks. The remainder of this paper builds on this observation, presenting an architecture that addresses these limitations through a combination of log-based capture, scalable transport, and reliability-focused design principles [9].

3. Problem Statement: The Oracle Legacy Problem

While Change Data Capture provides a conceptual pathway toward real-time data synchronization, its practical implementation in Oracle-based environments introduces a distinct set of challenges that go beyond those encountered in more modern or cloud-native database systems. Oracle databases have been widely adopted in large enterprises due to their transactional robustness, mature ecosystem, and support for complex schemas. However, these same characteristics contribute to significant friction when attempting to extract and propagate changes at scale in a low-latency, distributed architecture.

Unlike systems designed with native streaming interfaces, Oracle relies on internal mechanisms—such as redo logs and System Change Numbers (SCNs)—that must be interpreted externally for CDC purposes. This introduces both technical and operational complexity, particularly in environments with high transaction throughput, frequent schema evolution, and stringent consistency

requirements. The following subsections outline the primary challenges that define what we refer to as the *Oracle legacy problem* in the context of CDC-based cloud migration.

3.1 LogMiner Constraints

Oracle LogMiner is a commonly used utility for accessing redo logs and extracting change events. Although it provides a standardized interface for log-based CDC, its design imposes limitations that become pronounced at scale. LogMiner requires periodic dictionary builds to interpret redo records, which introduces overhead and latency, particularly when processing large volumes of changes. Additionally, its execution model offers limited support for parallel processing, making it difficult to keep pace with high redo generation rates typical of enterprise retail systems during peak activity periods. As transaction volumes increase, these constraints can lead to lag accumulation and reduced timeliness of downstream data propagation [10].

3.2 Transaction Ordering and Consistency

Maintaining transactional consistency across distributed systems is a central requirement of any CDC pipeline. In Oracle, ordering is governed by the System Change Number (SCN), which reflects the logical sequence of changes within the database. However, translating SCN-based ordering into a distributed streaming environment introduces complexity. Transactions may span multiple tables and partitions, and commit order must be preserved to ensure correctness. When change events are distributed across multiple processing nodes or message partitions, reconstructing a globally consistent order becomes non-trivial. Furthermore, dependencies between transactions—such as foreign key relationships—require careful coordination to prevent inconsistencies in downstream systems [11].

3.3 Schema Complexity

Enterprise retail databases often exhibit highly complex schemas, with thousands of tables representing diverse business domains such as inventory, pricing, customer data, and supply chain operations. These schemas are not static; they evolve frequently in response to changing business requirements. Data Definition Language (DDL) operations—such as adding columns, modifying data types, or restructuring tables—must be captured and propagated alongside data changes. Ensuring that downstream systems remain compatible with these evolving schemas requires robust schema versioning and transformation mechanisms. In practice, many CDC tools provide limited support for dynamic schema evolution, leading to potential disruptions or data inconsistencies when changes occur [12].

3.4 Network Egress Constraints

In hybrid or multi-cloud deployments, Oracle databases are often hosted on-premises or in private data centers, while analytical systems reside in public cloud environments. Transferring large volumes of change data across these boundaries introduces both performance and cost considerations. Redo logs can grow rapidly in high-throughput systems, and continuous extraction and transmission of these logs can saturate network bandwidth. Latency introduced by network transfer further impacts the timeliness of data availability in the cloud. Additionally, egress costs associated with data transfer can become significant at petabyte scale, necessitating efficient compression and batching strategies to optimize resource utilization [13].

3.5 Continuous Writes and Consistency

Retail systems operate under continuous, high-frequency write workloads, particularly during peak business events such as promotions or seasonal sales. In such environments, CDC pipelines must handle a constant stream of in-flight transactions while ensuring that only committed changes are propagated downstream. This requires careful management of transaction boundaries, buffering of partial transactions, and precise checkpointing to avoid data loss or duplication. Restarting the

pipeline after a failure must also preserve consistency, ensuring that no committed changes are missed and no uncommitted changes are incorrectly propagated. Achieving this level of correctness demands a tightly coordinated approach to offset tracking and transaction state management [14].

4. Architecture Design

The proposed architecture adopts a modular, streaming-first design that prioritizes scalability, fault tolerance, and low-latency data propagation. Rather than treating CDC as a monolithic pipeline, the system is decomposed into distinct layers, each responsible for a well-defined function in the end-to-end data flow. This separation of concerns enables independent scaling, simplifies failure isolation, and allows for targeted optimization at each stage. The design is guided by stringent performance requirements derived from production workloads, including sub-second latency targets and sustained multi-million event throughput.

Table 2: Architecture Component Responsibilities

Layer	Technology Used	Function
Source Capture	LogMiner	Extract redo logs
Serialization	Avro	Encode events
Streaming	Kafka	Transport events
Schema Management	Schema Registry	Schema evolution
Storage	Parquet (Cloud)	Analytical storage
Consumption	Warehouse/ML	Analytics

4.1 Source Capture Layer

The source capture layer is responsible for extracting change events directly from Oracle databases with minimal impact on transactional workloads.

- **Approach:** The architecture employs log-based CDC by reading Oracle redo logs, ensuring complete and ordered capture of all committed transactions. This avoids the overhead associated with trigger-based approaches and eliminates the need for repeated table scans.
- **Implementation:** A custom extractor is built on top of Oracle LogMiner, designed to overcome its inherent limitations. To improve throughput, the extractor parallelizes log processing by dividing redo streams into SCN ranges that can be processed concurrently. Care is taken to preserve transaction boundaries and commit order during this parallelization. Additionally, dictionary caching strategies are used to reduce the overhead of repeated metadata resolution.
- **Alternative Considered:** Oracle GoldenGate was evaluated as a potential solution due to its mature replication capabilities. However, it was ultimately not adopted due to high licensing costs and limited extensibility for custom processing logic and integration with modern streaming ecosystems [15].

4.2 Change Event Serialization

Once captured, raw change events must be transformed into a format suitable for efficient transmission and downstream processing.

- Events are serialized using **Apache Avro**, chosen for its compact binary representation and built-in support for schema evolution [16]. This reduces network overhead while ensuring compatibility across different versions of data consumers.
- Each event is enriched with metadata, including the System Change Number (SCN), transaction identifier, commit timestamp, and operation type (insert, update, delete). This metadata is critical for maintaining ordering, enabling replay, and supporting downstream auditing and reconciliation processes.

4.3 Transport Mechanism

The transport layer provides a durable, scalable backbone for distributing change events across the system.

- **Apache Kafka** serves as the central streaming platform, offering high-throughput ingestion, fault tolerance, and horizontal scalability. Its append-only log structure aligns naturally with CDC workloads, where events are processed sequentially and retained for replay when necessary.
- A key design consideration is the partitioning strategy. Events are partitioned based on a hash of the primary key, ensuring that all changes related to a specific entity are routed to the same partition. This guarantees ordering within the scope of individual entities while enabling parallel processing across partitions [17].

4.4 Schema Registry Integration

To manage schema evolution and ensure compatibility between producers and consumers, the architecture incorporates a centralized schema registry.

- The schema registry maintains versioned schemas for all event types, allowing producers to publish schema changes and consumers to validate compatibility before processing events.
- This approach enables backward and forward compatibility, reducing the risk of pipeline failures due to schema mismatches and supporting independent evolution of upstream and downstream components [18].

4.5 Cloud Landing Zone

The cloud landing zone serves as the persistent storage layer for captured change data, enabling both real-time and batch consumption patterns.

- Data is written to object storage systems in columnar formats such as **Parquet**, which provide efficient compression and query performance for analytical workloads.
- Partitioning is performed based on event time and entity type, enabling efficient pruning during query execution and supporting incremental data processing strategies [19].
- The landing zone also acts as a historical archive, allowing for replay, auditing, and backfill operations when required.

4.6 Downstream Consumers

The architecture supports a diverse set of downstream consumers, reflecting the varied analytical and operational needs of retail organizations.

- **Analytical warehouses** (e.g., Snowflake, BigQuery) consume CDC data to provide near-real-time reporting and business intelligence capabilities.

- **Real-time applications**, including dashboards and alerting systems, leverage streaming data for immediate insights into operational metrics.
- **Machine learning pipelines** utilize continuously updated data for model training and inference, enabling adaptive and personalized customer experiences.

The decoupled design of the pipeline allows multiple consumers to independently process the same stream of events without impacting each other's performance or reliability.

4.7 Design Targets

The architecture is engineered to meet strict performance and reliability objectives derived from production requirements:

- **Latency:** End-to-end latency—from transaction commit in Oracle to availability in downstream systems—is targeted at less than two seconds under normal operating conditions.
- **Throughput:** The system is designed to sustain throughput exceeding two million events per second, with the ability to scale horizontally during peak demand periods.
- **Durability:** The pipeline guarantees zero data loss under node failures through replication, persistent storage, and replay mechanisms built into the streaming infrastructure.

5. Consistency and Reliability

In large-scale CDC pipelines, correctness is as critical as performance. The value of near-real-time data diminishes rapidly if it cannot be trusted to reflect the true state of the source system. Consequently, the architecture places strong emphasis on consistency guarantees, fault tolerance, and deterministic recovery. This section outlines the mechanisms used to ensure reliable data delivery and to maintain alignment between the Oracle source and downstream systems under both normal and failure conditions.

5.1 Delivery Semantics

The system is designed to provide at-least-once delivery semantics, ensuring that every committed change in the source database is eventually propagated to downstream consumers. While this model inherently allows for the possibility of duplicate events, the architecture mitigates this through the use of idempotent consumers.

Each change event carries a unique combination of identifiers—such as primary key, SCN, and transaction ID—which allows downstream systems to detect and safely ignore duplicates. In practice, this approach achieves effective exactly-once processing, without incurring the additional complexity and performance overhead typically associated with strict exactly-once guarantees in distributed systems [20].

5.2 Offset Management

Reliable recovery from failures requires precise tracking of progress within the change stream. In this architecture, offsets are defined in terms of Oracle SCNs, which represent a globally ordered sequence of changes.

SCN offsets are persisted in an external, fault-tolerant store independent of the streaming infrastructure. This design ensures that, in the event of a failure or restart, the capture process can resume from a well-defined point without ambiguity. By decoupling offset storage from processing nodes, the system avoids inconsistencies that could arise from partial state loss or node-level failures.

Additionally, offset commits are coordinated with downstream acknowledgments to ensure that only fully processed and safely persisted events are marked as complete.

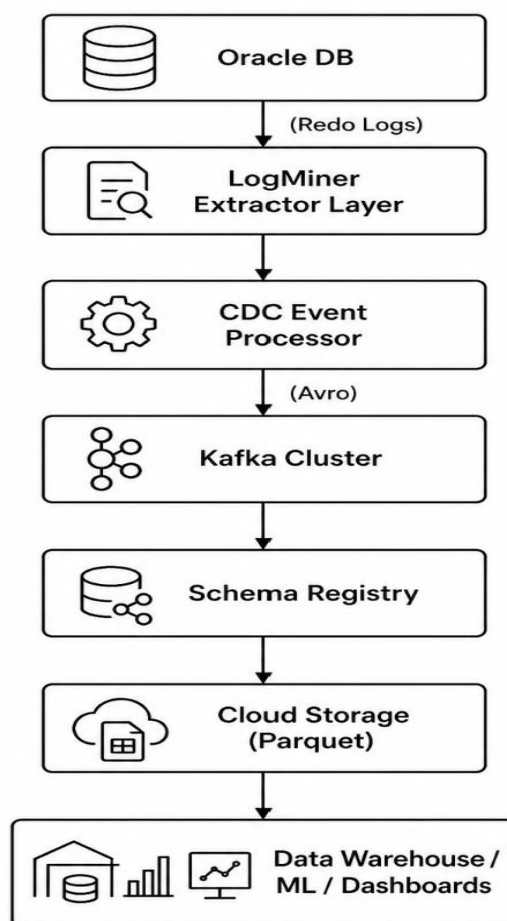


Figure 1: End-to-End Architecture for Oracle-to-Cloud CDC Pipeline

5.3 Failure Handling

Given the distributed nature of the architecture, failures are treated as expected events rather than exceptions. The system incorporates multiple mechanisms to handle transient and persistent failures gracefully:

- Automatic retries: Transient failures—such as temporary network disruptions or service unavailability—are handled by retry mechanisms with exponential backoff. This prevents overwhelming dependent systems while ensuring eventual progress.
- Dead-letter queues (DLQs): Events that cannot be processed due to schema mismatches, data corruption, or other irrecoverable issues are routed to dedicated dead-letter queues. This prevents pipeline blockage while preserving problematic records for later inspection and remediation.

These mechanisms ensure that failures do not propagate uncontrollably through the system or result in data loss.

5.4 Data Drift Detection

Over time, discrepancies may arise between the source database and downstream replicas due to missed events, processing errors, or operational anomalies. To detect such inconsistencies, the architecture incorporates periodic reconciliation processes.

These processes perform hash-based validation by computing checksums over selected datasets in both the source and target systems. Differences in hash values indicate potential data drift, triggering further investigation or corrective action. This approach provides a scalable means of validating data integrity without requiring full dataset comparisons [21].

5.5 Recovery Strategy

The architecture leverages the durability and replay capabilities of the streaming platform to enable robust recovery.

- **Replay from Kafka topics:** Since all change events are persisted in Kafka with configurable retention periods, the system can replay events from any prior offset. This allows recovery from both transient and prolonged failures without data loss.
- **Backfill pipelines:** In cases where discrepancies are detected or historical data needs to be re-synchronized, dedicated backfill processes are used to reconstruct the correct state. These pipelines operate independently of the real-time stream, ensuring that corrective actions do not disrupt ongoing data ingestion.

6. Performance Evaluation

The effectiveness of the proposed CDC architecture is evaluated through a series of production-aligned benchmarks designed to reflect real-world retail workloads. The evaluation focuses on key system characteristics, including latency, throughput, resource utilization, and resilience under failure conditions. Where applicable, results are compared against a baseline batch ETL system to quantify improvements in data freshness and operational efficiency.

Table 3: Performance Benchmarks

Metric	Value
Median Latency	850 ms
99th Percentile	2.4 sec
Sustained Throughput	2.5M events/sec
Peak Throughput	4M events/sec

6.1 Experimental Setup

The evaluation environment closely mirrors a large-scale enterprise deployment:

- An Oracle database generating approximately 5 TB per day of redo logs, representing high-volume transactional activity typical of retail systems during peak periods.
- A distributed Apache Kafka cluster consisting of 50 brokers, configured for high availability and horizontal scalability.
- A cloud-based analytical warehouse serving as the downstream consumption layer, integrated with object storage for persistent data ingestion.

The system operates under continuous load, with synthetic bursts introduced to simulate peak retail events such as promotional campaigns and seasonal demand spikes.

6.2 Latency

Latency is measured as the elapsed time between transaction commit in the Oracle source and availability of the corresponding change event in the cloud warehouse.

- Median latency: 850 milliseconds
- 99th percentile latency: 2.4 seconds

These results demonstrate that the architecture consistently meets its sub-second latency targets under normal operating conditions, with only modest degradation under high-load scenarios. The tail latency remains within acceptable bounds for real-time analytical applications, including dashboards and operational monitoring systems.

6.3 Throughput

Throughput is evaluated in terms of the number of change events processed per second across the entire pipeline.

- Sustained throughput: 2.5 million events per second
- Peak throughput: 4 million events per second

The system maintains stable performance under sustained load and exhibits the ability to scale elastically during peak demand. This is achieved through horizontal scaling of Kafka partitions and parallelized log extraction in the source capture layer.

6.4 Resource Utilization

Efficient resource utilization is critical for maintaining cost-effectiveness at scale.

- CPU utilization: Remains stable between 65–75% across processing nodes, indicating balanced workload distribution and effective parallelization.
- Network utilization: Optimized באמצעות compression techniques applied to serialized event streams, significantly reducing bandwidth consumption during data transfer from on-premises environments to the cloud.

These results suggest that the system achieves high throughput without over-provisioning infrastructure resources.

6.5 Comparison with Batch ETL

To contextualize the performance improvements, the CDC pipeline is compared with a traditional batch ETL process previously used in the same environment.

Table 4: CDC vs Batch ETL

Metric	Batch ETL	CDC Pipeline
Latency	6–24 hours	< 2 seconds
Freshness	Daily	Near real-time
Recovery	Manual	Automated

The comparison highlights a fundamental shift in data availability. While batch ETL introduces significant delays and requires manual intervention for recovery, the CDC pipeline delivers continuous data updates with automated fault handling, enabling real-time analytics and decision-making.

6.6 Failure Scenarios

The resilience of the system is evaluated באמצעות controlled failure simulations, including node outages and network disruptions.

- Recovery time under node failure: Less than 30 seconds, achieved through automatic failover and replay mechanisms.
- Data integrity: No data loss observed during simulated failures, validating the effectiveness of offset management and durable event storage.

These results confirm that the architecture maintains both availability and correctness under adverse conditions, meeting the reliability requirements of enterprise-scale retail systems.

7. Retail-Specific Considerations

While the architectural principles outlined in this paper are broadly applicable to large-scale data systems, their practical significance becomes particularly evident within the retail domain. Retail environments impose unique operational demands characterized by high transaction volumes, rapid business cycles, and strong coupling between data freshness and revenue outcomes. The following considerations highlight how these domain-specific factors shape both the requirements and design decisions of CDC-based data pipelines.

7.1 Peak Transactional Loads

Retail systems experience highly variable workloads, with extreme spikes during promotional events, seasonal sales, and major shopping periods. During such intervals, transaction rates can increase by orders of magnitude within short time windows. CDC pipelines must therefore be capable of elastic scaling, ensuring that ingestion, processing, and transport layers can dynamically adjust to increased load without introducing unacceptable latency or backpressure. Failure to handle such bursts effectively can result in data lag, undermining the real-time value of the system [22].

7.2 Schema Evolution

Retail data models are subject to frequent change, driven by evolving product catalogs, pricing strategies, and inventory management requirements. Schema modifications—ranging from the addition of new attributes to structural changes in key entities—must be propagated seamlessly through the data pipeline. This necessitates backward- and forward-compatible schema evolution mechanisms, ensuring that both legacy and updated consumers can continue to operate without disruption. Robust schema governance becomes essential to prevent cascading failures across dependent systems [23].

7.3 Machine Learning Integration

Modern retail operations increasingly rely on machine learning models for tasks such as recommendation systems, demand forecasting, and fraud detection. These models depend on continuously refreshed data to maintain accuracy and relevance. CDC pipelines enable near-real-time feature updates, reducing the lag between observed customer behavior and model adaptation. Consequently, the effectiveness of downstream ML systems is directly tied to the latency and reliability of the underlying data ingestion architecture [24].

7.4 Business Cost of Latency

In retail, latency is not merely a technical metric but a direct determinant of business outcomes. Delays in propagating inventory updates can lead to overselling, customer dissatisfaction, and operational inefficiencies. Conversely, delayed visibility into stock replenishment can result in missed sales opportunities due to perceived unavailability. These scenarios highlight the tangible financial

impact of data latency, reinforcing the need for sub-second data propagation in critical workflows [25].

7.5 Data Heterogeneity

Retail ecosystems integrate a wide variety of data sources, including transactional records, customer interaction logs, and supply chain data. These datasets differ in structure, velocity, and semantics, yet must be combined to provide a unified analytical view. CDC pipelines must therefore support heterogeneous data ingestion, enabling consistent handling of diverse data types while preserving their contextual relationships. This requirement underscores the importance of flexible serialization formats and unified streaming frameworks [26].

8. Discussion

The architecture presented in this paper demonstrates that log-based CDC, when combined with distributed streaming platforms and cloud-native storage systems, can effectively meet the stringent requirements of enterprise retail environments. By leveraging redo log extraction, scalable event transport, and robust consistency mechanisms, the system achieves both low latency and high throughput without compromising reliability.

A key insight from this work is that CDC should not be treated as a standalone capability but rather as an integral component of a broader data platform. Its success depends on tight integration with schema management, fault tolerance strategies, and downstream consumption patterns. The modular design adopted here allows each component to evolve independently while maintaining overall system coherence.

However, several challenges remain open for further exploration. Multi-region replication introduces additional complexity in maintaining consistency across geographically distributed deployments, particularly in the presence of network partitions. Cross-database consistency, where multiple heterogeneous sources must be synchronized, poses further difficulties in ensuring global ordering and correctness. Additionally, the operational overhead associated with managing large-scale streaming infrastructure remains non-trivial, requiring specialized expertise and tooling.

Future work may explore the use of serverless streaming architectures, which could reduce operational complexity while maintaining scalability. Similarly, deeper integration with cloud-native databases and storage systems may enable more efficient data movement and tighter consistency guarantees. Advances in these areas have the potential to further simplify CDC deployments while extending their applicability to a wider range of use cases.

9. Conclusion

This paper has presented a scalable and reliable architecture for real-time Change Data Capture, enabling near-real-time synchronization between Oracle databases and cloud-based analytical systems at petabyte scale. By addressing the inherent challenges of Oracle log extraction, distributed event processing, and consistency management, the proposed solution achieves sub-second latency and high-throughput performance under demanding production workloads.

The evaluation demonstrates that carefully designed CDC pipelines can outperform traditional batch ETL approaches by several orders of magnitude in terms of latency and operational responsiveness. More importantly, the architecture provides a foundation for continuous data availability, supporting real-time analytics, machine learning applications, and operational decision-making in large retail enterprises.

These findings underscore the broader viability of transitioning from batch-oriented data integration to streaming-based architectures. As organizations continue to modernize their data infrastructure,

CDC will play a central role in bridging legacy systems with cloud-native platforms, enabling a new generation of data-driven capabilities.

References

- [1] Stonebraker, M., Çetintemel, U., & Zdonik, S. (2018). The 8 requirements of real-time stream processing. *ACM SIGMOD Record*, 47(1), 42–47.
- [2] Kreps, J. (2014). Questioning the lambda architecture. *Confluent Blog*.
- [3] Oracle Corporation. (2019). *Oracle database triggers documentation*.
- [4] Golfarelli, M., Rizzi, S., & Cella, I. (2017). Beyond data warehousing: What's next in business intelligence? *Proceedings of DOLAP*.
- [5] Debezium. (2021). *Debezium documentation*.
- [6] Kleppmann, M. (2017). *Designing data-intensive applications*. O'Reilly Media.
- [7] Amazon Web Services. (2021). *AWS Database Migration Service documentation*.
- [8] Oracle Corporation. (2020). *Oracle GoldenGate documentation*.
- [9] Narkhede, N., Shapira, G., & Palino, T. (2017). *Kafka: The definitive guide*. O'Reilly Media.
- [10] Oracle Corporation. (2020). *Oracle LogMiner utility guide*.
- [11] Oracle Corporation. (2021). *Oracle database concepts*.
- [12] Fowler, M. (2018). *Refactoring: Improving the design of existing code*. Addison-Wesley.
- [13] Google Cloud. (2021). *Data migration and networking architecture guide*.
- [14] Helland, P. (2016). Life beyond distributed transactions. *Communications of the ACM*, 59(2), 44–50.
- [15] Oracle Corporation. (2021). *Oracle GoldenGate licensing information*.
- [16] Apache Software Foundation. (2021). *Apache Avro specification*.
- [17] Apache Kafka. (2022). *Apache Kafka documentation*.
- [18] Confluent. (2021). *Schema Registry documentation*.
- [19] Armbrust, M., Xin, R., Lian, C., et al. (2015). Spark SQL: Relational data processing in Spark. *Proceedings of the 2015 ACM SIGMOD*.
- [20] Kreps, J., Narkhede, N., & Rao, J. (2017). Kafka: A distributed messaging system for log processing. *Proceedings of NetDB*.
- [21] Abadi, D. (2019). Consistency tradeoffs in modern distributed database system design. *IEEE Data Engineering Bulletin*.
- [22] Retail Systems Research. (2021). *Retail performance benchmarks report*.
- [23] Dehghani, Z. (2020). *Data mesh: Delivering data-driven value at scale*. O'Reilly Media.
- [24] Covington, P., Adams, J., & Sargin, E. (2016). Deep neural networks for YouTube recommendations. *Proceedings of the 10th ACM Conference on Recommender Systems*.
- [25] McKinsey & Company. (2021). *Retail analytics: Driving customer value*.
- [26] Gartner. (2022). *Data integration tools market guide*.