# A Real-Time Analytics Architecture for Enterprise Order Lifecycle Visibility and Backlog Management

**Ranveer Potel**

Potel Projects LLC, USA

Email: potelprojects@gmail.com

**ABSTRACT:** Enterprise order management systems traditionally rely on batch-oriented data pipelines that introduce significant latency between transactional events and analytical insights. In large-scale global supply chains, delays in reporting order status, fulfillment availability, and backlog trends can hinder operational responsiveness and reduce customer satisfaction. This paper presents the design and implementation of a real-time analytics architecture that enables near-instant visibility into the order lifecycle across distributed enterprise systems. The proposed architecture integrates change data capture (CDC), distributed streaming pipelines, scalable NoSQL storage, and search-based indexing to deliver sub-second analytics for order tracking, backlog management, and fulfillment decision support. The system replaces legacy batch reporting processes and enables continuous synchronization between transactional enterprise resource planning (ERP) systems and analytical dashboards. Experimental results demonstrate a reduction in data latency from approximately 24 hours to sub-second availability, a reduction in report load time from several minutes to seconds, and improved scalability supporting large-scale operational analytics across global business units. The architecture has been deployed across multiple enterprise operational units and serves as a reference implementation for event-driven analytics in supply chain environments.

**KEYWORDS:** Real-time analytics, enterprise order management, streaming data pipelines, change data capture, distributed systems, supply chain analytics, backlog management, ERP integration, NoSQL databases, event-driven architecture.

## I. INTRODUCTION

Modern supply chains require continuous visibility into order lifecycles to ensure efficient fulfillment and accurate customer commitments. Enterprises managing global order volumes must coordinate multiple operational systems including inventory, procurement, production, logistics, and customer service applications. Traditional enterprise analytics architectures rely on periodic batch processing, which introduces delays between transactional events and analytical reporting [1]. In industries where supply and demand conditions can shift within hours, such architectural limitations carry direct operational and financial consequences.

In large-scale order management environments, reporting delays can exceed twenty-four hours, limiting the ability of operational teams to respond to demand fluctuations, supply disruptions, or order fulfillment anomalies. Furthermore, backlog management and order promising decisions depend on timely insight into supply availability and demand priorities. When analytics platforms cannot reflect the current state of orders in motion, fulfillment teams may over-commit or under-allocate inventory, resulting in customer dissatisfaction and revenue leakage.

Event-driven analytics architectures have emerged as a promising solution to these challenges. By processing data streams as they are generated, rather than queuing them for periodic extraction, such systems enable near real-time alignment between operational transactions and analytical dashboards [2]. However, deploying these architectures within large enterprise environments presents substantial engineering challenges: transactional systems are heterogeneous, data volumes are high, schemas evolve over time, and reliability requirements are stringent.

This paper proposes a real-time analytics framework designed to enhance order lifecycle visibility, backlog management, and fulfillment planning across enterprise systems. The architecture leverages change data capture (CDC), distributed streaming platforms, scalable NoSQL storage, and real-time indexing to enable sub-second analytics

for operational reporting. The remainder of this paper is organized as follows. Section II reviews background and related work. Sections III through VII describe the proposed architecture and its constituent components. Section VIII presents performance evaluation results. Section IX discusses enterprise deployment and adoption outcomes. Section X addresses open challenges and directions for future research, followed by conclusions in Section XI.

## II. BACKGROUND AND RELATED WORK

Enterprise order management systems integrate multiple components that support demand capture, inventory allocation, fulfillment scheduling, and customer delivery commitments. A key capability within these systems is order promising, which determines when and from where a customer order can be fulfilled. This process typically draws on data from multiple heterogeneous systems and depends heavily on the freshness of supply and demand information.

### A. Order Promising Mechanisms
Order promising mechanisms operate across two primary paradigms: Available-to-Promise (ATP) and Capable-to-Promise (CTP). ATP evaluates on-hand inventory and scheduled receipts to determine fulfillment feasibility within a defined time horizon. CTP extends this by incorporating production and procurement capacity, enabling commitment decisions for make-to-order or configure-to-order scenarios. Both mechanisms typically consider the following inputs:
- Inventory availability across multiple supply sources and warehouses
- Production schedules from manufacturing execution systems
- Procurement lead times and supplier commitments
- Logistics constraints including carrier availability and transit times
- Customer priority tiers and channel-level allocation policies

In practice, order promising accuracy degrades significantly when these inputs are stale. Batch-refreshed analytics systems frequently present an outdated view of available supply, leading to fulfillment commitments that cannot be honored or unnecessary demand deferrals.

### B. Traditional Batch Analytics Architectures
Traditional enterprise analytics for order management rely on data warehouses refreshed through batch extraction, transformation, and loading (ETL) pipelines [1]. These pipelines extract data from transactional systems on a scheduled cadence—often nightly—and load transformed datasets into centralized repositories for reporting. While batch architectures perform well for historical analysis and financial reconciliation, they are poorly suited to operational monitoring tasks where sub-hourly data freshness is required.

Several limitations characterize batch-oriented analytics in order management contexts: (i) latency between events and visibility can exceed 24 hours; (ii) ETL jobs introduce tight coupling between source systems and reporting schemas, impeding agility; (iii) warehouse-scale aggregation queries can be slow to execute under concurrent user load; and (iv) operational teams must often maintain informal side channels—spreadsheets, email threads, manual queries—to supplement delayed reporting.

### C. Streaming Analytics Approaches
Recent research and industry practice have explored streaming architectures as an alternative to batch ETL pipelines [2][3]. Streaming systems process data as events rather than in bulk, enabling continuous computation and low-latency propagation of state changes. Key technologies enabling this shift include:
- Change Data Capture (CDC): Log-based mechanisms that capture row-level changes from database transaction logs with minimal impact on source systems
- Distributed message brokers: Fault-tolerant event streaming platforms that decouple producers from consumers and provide durable, ordered event logs
- Stream processing frameworks: Distributed computation engines capable of performing stateful transformations, windowed aggregations, and join operations over high-throughput event streams
- Real-time indexing systems: Search engines capable of ingesting and indexing high-velocity data streams, enabling low-latency query execution over large, constantly updated datasets

Integrating these technologies within large enterprise order management environments remains an active area of engineering and research. Specific challenges include handling schema evolution across heterogeneous source systems, maintaining consistency under high write concurrency, ensuring pipeline observability, and providing graceful degradation under failure conditions [2]. The architecture described in this paper addresses these challenges in the context of a large-scale enterprise deployment.

## III. SYSTEM ARCHITECTURE

The proposed architecture transforms a traditional batch analytics pipeline into a streaming analytics platform capable of processing transactional updates in near real time. The design follows a layered model in which each component performs a well-defined function and communicates with adjacent layers through standardized interfaces. This modularity supports independent scaling, technology substitution, and fault isolation.

The architecture consists of six primary layers:

• Transactional Systems Layer: Source ERP and order management applications that generate operational events

• Change Data Capture Layer: Log-based replication mechanism that converts database mutations into structured change events

• Streaming Ingestion Layer: Distributed event broker that buffers, routes, and durably stores change events for downstream consumption

• Distributed Processing and Storage Layer: Real-time computation engine and scalable NoSQL store that transform and persist enriched order data

• Indexing and Search Layer: Near real-time indexing platform that enables fast query execution over processed order data

• API and Analytics Delivery Layer: REST services and user-facing dashboards that expose processed analytics to operational users and enterprise applications

The overall data flow proceeds as follows: transactions in source ERP systems produce database log entries that are captured by the CDC layer and published as structured events into the streaming ingestion layer. A distributed processing engine consumes these events, applies transformations and enrichment logic, and writes results to persistent NoSQL storage. The indexing layer continuously synchronizes with the storage layer to maintain a queryable index of current order state. Downstream APIs and dashboards read from this index to deliver low-latency analytics to end users. Figure 1 provides a conceptual representation of this layered architecture.
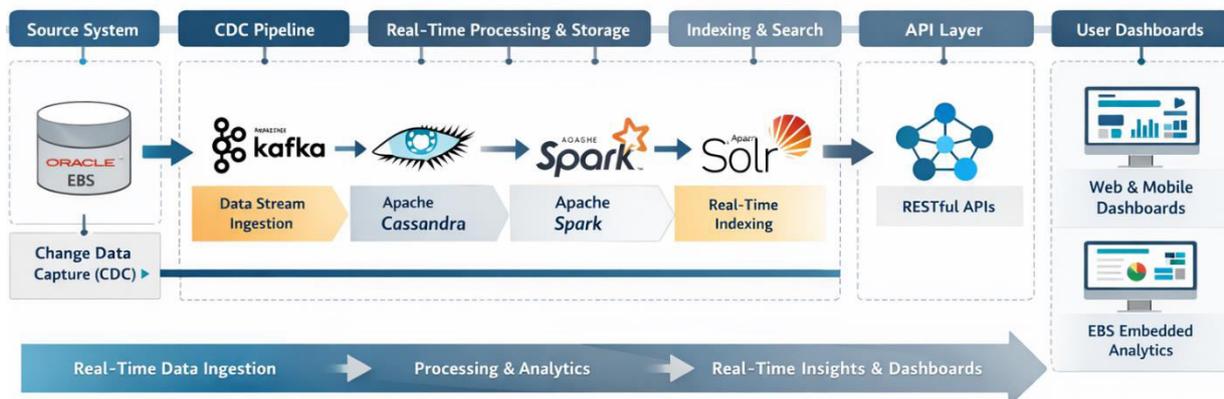


Fig. 1. High-level architecture of the proposed real-time analytics platform.

## IV. REAL-TIME DATA PIPELINE DESIGN

### A. Transactional Data Sources

Enterprise order management systems serve as the primary sources of operational data. These systems are typically built on relational database management systems (RDBMS) and process high volumes of concurrent transactions across geographically distributed business units. Each order lifecycle event generates one or more records within the transactional schema, including:

• Order creation and modification records capturing customer, channel, and product attributes
• Inventory allocation events representing commitments against available supply
• Production scheduling records for make-to-order demand fulfillment
• Shipment confirmation and carrier scan events for logistics visibility
• Backlog adjustment records reflecting demand prioritization and deferral decisions

Because source systems span multiple technology stacks and business domains, the architecture must handle heterogeneous schemas and varying transaction rates. Source systems in the evaluated deployment collectively generated on the order of tens of thousands of transactional events per minute during peak operational periods.

### B. Change Data Capture (CDC)

Rather than relying on periodic database extracts, the architecture captures transactional updates using a log-based change data capture mechanism. CDC reads directly from the database transaction log—the same durable write-ahead log that ensures ACID transactional guarantees—and produces a structured stream of row-level change events.

Each CDC event encodes the following metadata alongside the changed data payload: the operation type (insert, update, or delete), the affected table and schema, a monotonic log sequence number for ordering, and a high-resolution timestamp. This enriched event format allows downstream consumers to reconstruct the full history of changes to any order entity, enabling both real-time aggregation and point-in-time recovery scenarios.

Key operational benefits of log-based CDC include:

• Minimal performance impact on source transactional systems, as log reading is decoupled from query processing
• Continuous data replication with sub-second propagation latency from source commit to event publication
• Complete change history, including intermediate states that batch snapshots would not capture
• Support for schema change detection, enabling downstream consumers to adapt to source schema evolution

### C. Streaming Data Ingestion

Captured change events are published into a distributed event streaming platform that serves as the durable, fault-tolerant backbone of the analytics pipeline. The streaming layer is organized as a set of partitioned topic queues, with partitioning aligned to order identifiers to preserve per-order event ordering while enabling parallel consumption across multiple downstream consumers.

The streaming ingestion layer provides several critical operational capabilities:

• Durable event storage with configurable retention periods, enabling replay and recovery without dependency on source systems
• Consumer group isolation, allowing multiple downstream applications to consume the same event stream independently at their own pace
• Horizontal scalability through partition rebalancing, enabling throughput to be increased by adding consumers or brokers without architectural changes
• Delivery guarantees at the exactly-once or at-least-once level, depending on consumer implementation, ensuring correctness under failure conditions

Topic partitioning strategy is a critical design decision in this layer. Partitioning by order ID ensures that all events related to a given order are processed in sequence by a single consumer instance, preserving causal consistency. However, uneven key distributions—such as unusually high transaction rates for specific large customers—can create partition hotspots. The architecture addresses this through a secondary hash-based sub-partitioning scheme applied to high-cardinality customers.

## V. DISTRIBUTED PROCESSING AND STORAGE

### A. Real-Time Processing Engine

A distributed stream processing engine consumes events from the ingestion layer and performs stateful, real-time transformations that derive analytical metrics from raw transactional data. The processing engine is designed for fault

tolerance through state checkpointing and supports exactly-once processing semantics to prevent double-counting of transactional events under failure and recovery scenarios.

Core processing operations include:

• Data normalization: Resolving encoding differences, unit conversions, and reference data lookups across heterogeneous source schemas

• Order lifecycle aggregation: Maintaining current order state by applying a stream of incremental updates against a keyed state store

• Supply-demand correlation: Joining order demand records with available inventory and scheduled supply receipts to compute real-time fulfillment feasibility

• Backlog computation: Continuously updating backlog metrics by aggregating unfulfilled demand across time horizons, product families, and business units

• Fulfillment status updates: Propagating shipment and delivery events to order records, enabling real-time order completion tracking

Windowed aggregation operations are used to compute time-series metrics such as order intake rate, fulfillment cycle time distributions, and rolling backlog trends. These windows are evaluated continuously as new events arrive, rather than at fixed batch intervals, ensuring that dashboard consumers always receive the most current available metrics.

Stateful operations are backed by a distributed, replicated state store embedded within the processing engine. State is checkpointed to durable storage at configurable intervals, enabling recovery from processor failures without loss of accumulated analytical state. In the evaluated deployment, a checkpoint interval of 30 seconds was selected as a balance between recovery time objective and checkpoint overhead.

### B. Scalable Data Storage

Processed events and derived analytical state are persisted to a distributed NoSQL database optimized for high write throughput and horizontal scalability. The storage layer is designed to accommodate continuous, high-velocity writes from the processing engine without introducing write-path bottlenecks.

The data model is designed around the primary query patterns of the analytics platform. Order entity records are keyed by order identifier and maintained as document-oriented structures that aggregate all relevant attributes—line items, status history, allocation records, shipment data—into a single logical entity. This denormalized design avoids multi-entity joins at query time, ensuring consistent low-latency reads regardless of dataset scale.

Key storage design properties include:

• Partition-aware data distribution aligned with the streaming layer's partitioning scheme, minimizing cross-partition operations

• Tunable consistency models allowing time-sensitive dashboard queries to trade strict consistency for lower read latency

• Multi-region replication to support globally distributed operational teams with low read-path latency

• TTL-based data expiration policies to manage storage costs for historical order records beyond the active analytics window

## VI. REAL-TIME INDEXING AND SEARCH

To enable interactive analytics and user-driven dashboard queries, the architecture incorporates a real-time search and aggregation indexing layer. This layer is populated continuously from the NoSQL storage layer through a dedicated synchronization process that propagates writes to the index within seconds of their storage-layer commit.

The indexing layer supports the following analytical query capabilities:

• Full-text and structured attribute search across order fields including customer name, product description, order reference, and geographic region

• Pre-computed aggregation queries for backlog summaries, fulfillment performance metrics, and demand volume trend analysis

• Faceted filtering allowing users to slice operational data by order status, business unit, product family, date range, and priority tier

• Time-series histogram analysis enabling visualization of order intake, fulfillment, and cancellation trends across configurable time windows

The indexing platform is scaled horizontally to handle concurrent query load from operational dashboards and API consumers. Index sharding is aligned to business unit boundaries, ensuring that frequently executed regional queries are served by dedicated shard resources without contention from global aggregation workloads.

A critical operational requirement for the indexing layer is near real-time index refresh latency. In the evaluated deployment, end-to-end latency from transactional source commit to index-queryable state was measured at under 3 seconds at the 95th percentile during normal operating conditions, and under 8 seconds at the 99th percentile during peak transaction periods. These figures represent a reduction of several orders of magnitude compared to the overnight batch refresh cadence of the legacy system.

## VII. ANALYTICS DELIVERY AND INTEGRATION

### A. API Services

The analytics platform exposes a set of REST-based API services that provide structured access to real-time order data for a range of consuming applications. APIs are versioned and documented through an OpenAPI specification to support integration across enterprise teams with varying technical capabilities.

API design follows resource-oriented principles, with distinct endpoint families for order entity queries, backlog summary metrics, fulfillment performance reporting, and supply availability lookups. Response payloads are structured to minimize client-side computation, with pre-computed aggregations and inline enrichment data returned directly from the indexing layer.

The API tier enables:

- Web-based operational dashboards consuming real-time order and backlog data
- Mobile monitoring applications for field and logistics operations teams
- Embedded analytics within ERP and customer service platforms through iframe and widget integration patterns
- Automated monitoring and alerting systems that poll critical backlog and fulfillment metrics on configurable schedules

API rate limiting and caching strategies are implemented to protect the indexing layer from excessive query load during peak usage periods. Frequently requested aggregation results are cached with short time-to-live values, ensuring that high-concurrency dashboard refreshes do not generate proportional indexing layer load.

### B. Operational Dashboards

User-facing dashboards provide real-time visualization of key order management metrics for multiple operational audiences. Dashboard design follows a role-based information architecture: supply chain planners, fulfillment operations teams, finance controllers, and customer service representatives each access views tailored to their specific operational responsibilities.

Metrics surfaced through dashboards include:

- Global order volume trends and intake rate time series
- Current backlog status by product family, business unit, and fulfillment horizon
- Fulfillment performance including on-time shipment rate, cycle time distribution, and exception counts
- Order lifecycle progress with per-order status drill-down and timeline visualization
- Supply availability heatmaps and anticipated fulfillment date confidence intervals

Dashboards are refreshed on a sub-minute cadence using server-sent event streams, ensuring that operational teams observe near-current data during active business hours. Alert thresholds are configurable at the dashboard level, enabling teams to receive proactive notifications when backlog depth, order cycle time, or fulfillment exception rates exceed defined operational bounds.

## VIII. PERFORMANCE EVALUATION

The proposed architecture was evaluated in a production enterprise environment supporting high transaction volumes across multiple geographic regions. Performance benchmarks were conducted under representative operational load conditions, with transaction rates reflective of peak demand periods.

### A. Data Latency

End-to-end data latency—measured from transactional source commit to dashboard-queryable state—was reduced from an overnight batch refresh cycle of approximately 24 hours to a median of under 1 second under normal operating conditions. At the 99th percentile, end-to-end latency measured under 8 seconds during peak load periods. This represents a reduction in latency of approximately four orders of magnitude compared to the legacy system.

### B. Report Performance

Dashboard query execution time, as measured at the indexing layer, was reduced from several minutes under the legacy batch warehouse to consistently sub-second median response times under the streaming architecture. Complex aggregation queries involving multi-dimensional filtering across millions of order records executed within 2–3 seconds at the 95th percentile, meeting interactive dashboard responsiveness requirements.

### C. Comparative Summary

Table I summarizes the performance comparison between the legacy batch system and the proposed real-time platform across key operational metrics.

### TABLE I
### Performance Comparison: Legacy Batch System vs. Real-Time Analytics Platform

| Metric | Legacy Batch System | Real-Time Platform |
|---|---|---|
| Data Latency | ~24 hours | < 1 second |
| Report Load Time | Several minutes | Seconds |
| Order Visibility | Delayed (end of day) | Immediate |
| Operational Response | Reactive | Proactive |
| Concurrent Users | Limited (~50) | Scalable (500+) |
| Pipeline Failure Recovery | Next batch cycle | Automatic, milliseconds |

### D. Scalability

The architecture demonstrated linear scalability with respect to transaction throughput during load testing. The streaming ingestion layer sustained throughput of over 50,000 events per second without message loss at a three-node broker configuration. The processing engine maintained sub-second processing latency up to 30,000 events per second per worker, with horizontal scaling through worker addition required beyond this threshold. The indexing layer sustained concurrent query loads of over 500 simultaneous active dashboard sessions without exceeding target query response time thresholds.

### E. Fault Tolerance

Fault tolerance was validated through controlled failure injection scenarios including broker node failure, processing engine worker failure, and storage node failure. In all scenarios, the architecture recovered within the operational recovery time objective of under 60 seconds, with no data loss attributable to the failure event. CDC-based replication ensured that events generated during recovery periods were buffered in the ingestion layer and processed upon service restoration.

## IX. ENTERPRISE DEPLOYMENT AND ADOPTION

The platform was deployed incrementally across enterprise operational units responsible for supply chain management, manufacturing operations, finance, and customer service. Deployment followed a phased migration approach: legacy batch reporting pipelines remained operational in parallel with the streaming platform during a transition period, allowing teams to validate real-time data against established reporting baselines before decommissioning batch dependencies.

Adoption outcomes documented across participating business units included:
• Elimination of manual data reconciliation workflows that previously required analysts to merge data from multiple delayed reporting sources
• Improved backlog prioritization enabled by real-time visibility into demand depth and supply availability, reducing instances of missed fulfillment commitments
• Faster identification of fulfillment disruptions, with median time-to-awareness of supply exceptions reduced from next-day to within minutes of event occurrence

• Enhanced decision support for order promising, with fulfillment planners reporting increased confidence in availability commitments due to data freshness improvements

• Cross-departmental operational alignment enabled by a shared, continuously updated view of order lifecycle state accessible to supply chain, finance, and customer operations teams simultaneously

Stakeholder feedback collected through structured interviews with operational users highlighted data freshness and dashboard reliability as the most operationally significant improvements. Several teams reported changes to their operating procedures as a direct consequence of real-time visibility, shifting from end-of-day review cycles to continuous intraday monitoring practices.

The architecture has since been adopted as a reference implementation within the enterprise IT organization for deploying additional real-time analytics solutions beyond the order management domain, including procurement visibility, logistics tracking, and financial close monitoring.

## X. DISCUSSION

The implementation demonstrates that real-time analytics architectures can effectively replace traditional batch reporting systems in large enterprise environments. By integrating CDC pipelines, distributed streaming platforms, scalable storage, and real-time indexing, organizations can achieve continuous visibility into operational data with latency characteristics that enable proactive rather than reactive decision making.

Such architectures are particularly beneficial for supply chain and order management systems where timely information directly impacts customer commitments and revenue capture. The feedback loop between operational events and analytical insight is shortened to the point where corrections can be applied within the same operational shift in which anomalies occur, rather than being deferred to the following business day.

### A. Operational Monitoring
Streaming pipelines introduce new operational monitoring requirements that have no direct analog in batch ETL environments. Consumer lag—the accumulation of unprocessed events in the ingestion layer—is a critical health indicator that must be tracked and alerted against to ensure latency SLAs are maintained. Unlike batch jobs, which either succeed or fail at discrete points in time, streaming pipelines can degrade gradually, requiring continuous monitoring of throughput, latency, error rates, and state store utilization. Operational teams must be trained to interpret these metrics and respond appropriately, representing a capability investment beyond the technical infrastructure itself.

### B. Schema Evolution Management
Schema changes in source transactional systems present a recurring challenge for CDC-based architectures. When source schemas evolve—through column additions, type changes, or table restructurings—CDC events change structure, potentially breaking downstream consumers that parse fixed schemas. The architecture addresses this through a schema registry that maintains versioned schema definitions and enforces compatibility rules before schema changes are deployed to production. Consumers are designed to handle schema versions gracefully through forward compatibility policies, but coordination between source system development teams and the analytics platform remains an ongoing operational requirement.

### C. Cross-System Data Consistency
Distributed streaming architectures introduce eventual consistency characteristics that differ from the strong consistency models of traditional data warehouse ETL pipelines. In the proposed architecture, order entities assembled from multiple source system streams may exhibit temporary inconsistencies during the brief window between related events arriving from different CDC sources. These transient inconsistencies are typically resolved within seconds but must be accounted for in dashboard design and analytical interpretation. For use cases requiring strict consistency— such as financial reconciliation or audit reporting—supplementary batch processes over the NoSQL storage layer continue to provide consistent point-in-time snapshots.

### D. Future Research Directions
Several directions for future research and development are identified based on operational experience with the deployed architecture. First, the integration of machine learning models for demand forecasting and anomaly detection within the streaming pipeline presents an opportunity to extend the platform from descriptive to predictive analytics, enabling proactive alerting on emerging backlog risks or fulfillment anomalies before they materialize. Second, automated

backlog prioritization using reinforcement learning techniques could replace rule-based prioritization policies, dynamically balancing fulfillment commitments against supply constraints in real time. Third, the application of formal stream processing semantics to order promising logic—replacing the current batch-oriented available-to-promise computation with a continuously evaluated streaming model—represents a substantial architectural opportunity for further latency reduction in customer commitment workflows.

## XI. CONCLUSION

This paper presented a real-time analytics architecture designed to provide continuous visibility into enterprise order management operations. By replacing batch-based reporting pipelines with a streaming analytics framework integrating change data capture, distributed event streaming, scalable NoSQL storage, and real-time search indexing, the proposed solution achieves a reduction in data latency from approximately 24 hours to sub-second availability, dramatically improves query performance, and enables operational teams to respond to fulfillment events as they occur rather than after the fact.

The architecture was evaluated in a production enterprise deployment and demonstrated linear scalability, fault-tolerant operation, and measurable operational impact across supply chain, finance, and customer operations functions. Deployment experience confirmed that real-time data freshness fundamentally changes how operational teams interact with order management analytics, shifting established workflows from retrospective review to continuous intraday monitoring.

The implementation confirms that event-driven analytics architectures are operationally viable in large enterprise environments and deliver significant value in domains where information latency directly affects revenue, customer experience, and operational efficiency. As enterprise systems continue to generate increasing volumes of transactional data, real-time analytics platforms of the type described here will play an increasingly critical role in supporting data-driven decision making across global supply chains.

Future work will explore the integration of predictive analytics and machine learning capabilities within the streaming pipeline, the application of real-time order promising logic as a replacement for batch-oriented available-to-promise computation, and formal evaluation of consistency models appropriate for financial-grade analytics over eventually consistent streaming architectures.

## X. ACKNOWLEDGMENT

## REFERENCES

[1] M. Stonebraker et al., "The End of an Architectural Era: It's Time for a Complete Rewrite," Proceedings of the VLDB Endowment, vol. 7, no. 13, pp. 1150–1160, 2014.
[2] T. Akidau et al., Streaming Systems: The What, Where, When, and How of Large-Scale Data Processing. Sebastopol, CA: O'Reilly Media, 2018.
[3] J. Kreps, N. Narkhede, and J. Rao, "Kafka: A Distributed Messaging System for Log Processing," in Proc. 6th International Workshop on Networking Meets Databases (NetDB), Athens, Greece, Jun. 2011.
[4] P. Hunt, M. Konar, F. P. Junqueira, and B. Reed, "ZooKeeper: Wait-free Coordination for Internet-scale Systems," in Proc. USENIX Annual Technical Conference (ATC), Boston, MA, Jun. 2010.
[5] A. Lamb et al., "The Vertica Analytic Database: C-Store 7 Years Later," Proc. VLDB Endowment, vol. 5, no. 12, pp. 1790–1801, 2012.
[6] F. Hueske and V. Kalavri, Stream Processing with Apache Flink. Sebastopol, CA: O'Reilly Media, 2019.
[7] P. Carbone et al., "Apache Flink: Stream and Batch Processing in a Single Engine," IEEE Data Engineering Bulletin, vol. 38, no. 4, pp. 28–38, 2015.

[8] G. DeCandia et al., "Dynamo: Amazon's Highly Available Key-value Store," in Proc. 21st ACM Symposium on Operating Systems Principles (SOSP), Stevenson, WA, Oct. 2007, pp. 205–220.

[9] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google File System," in Proc. 19th ACM Symposium on Operating Systems Principles (SOSP), Bolton Landing, NY, Oct. 2003, pp. 29–43.

[10] R. Cattell, "Scalable SQL and NoSQL Data Stores," ACM SIGMOD Record, vol. 39, no. 4, pp. 12–27, 2011.

[11] D. Abadi, "Consistency Tradeoffs in Modern Distributed Database System Design," IEEE Computer, vol. 45, no. 2, pp. 37–42, 2012.

[12] M. Zaharia et al., "Apache Spark: A Unified Engine for Big Data Processing," Communications of the ACM, vol. 59, no. 11, pp. 56–65, 2016.