

Cloud-Native Data Architectures for Real-Time Marketing Analytics in Database DevOps

Kamini Murugaboopathy

Wonderbow Analytics Private Ltd., India

Conf42 Database DevOps 2026

Agenda

1 From Batch Intelligence to Real-Time Action

3 Data Storage: Lakes, Warehouses & Segmentation

5 Scaling, Observability & Cost Optimization

2 Event-Driven Architecture & Streaming Pipelines

4 DataOps, Schema Evolution & Model Deployment

6 Federated Learning, Implementation Paths & Lessons Learned

From Batch Intelligence to Real-Time Action

The competitive advantage in modern marketing no longer belongs to those who collect the most data, it belongs to those who can act on it as it happens.

Millisecond Behavior

Customer actions unfold faster than any batch window can capture

Seconds, Not Days

Campaign optimization windows have collapsed to seconds

Event-Driven Pipelines

Overnight batch jobs replaced by streaming, real-time workflows

Cloud + DevOps

Their intersection is what makes this transformation operationally viable

Three Structural Breakdowns in Traditional Pipelines

Latency

- 1 Overnight batch processing creates a persistent lag. By the time insights are available, the optimization window has already passed.

Fragmentation

- 2 Dozens of SaaS tools generate data on their own formats and schedules. Batch ETL cannot stitch these into a consistent real-time view.

Scalability

- 3 Static infrastructure cannot absorb spiky marketing workloads a flash sale can drive massive surges within minutes.

Event-Driven at the Foundation

The foundational shift: every customer interaction is treated as a **discrete, immediately actionable event**.

1

Page View

Captured and processed immediately as a signal of user intent. It feeds real-time dashboards, updates session context, and triggers personalization rules based on browsing behavior.

2

Click

Triggers downstream logic in real time. Each click event is routed to the appropriate processing stream, whether for A/B test tracking, recommendation engine updates, or behavioral scoring.

3

Cart Abandonment

Feeds the personalization engine instantly. Within milliseconds, abandonment events activate retargeting workflows, adjust audience segments, and can trigger automated recovery campaigns.

4

Purchase

Updates attribution models, revenue dashboards, and audience segmentation on the spot. It closes the conversion loop and recalibrates predictive models with fresh ground-truth data.

Data Lakes vs. Data Warehouses

Dimension	Data Warehouse	Data Lake
Examples	Snowflake, BigQuery, Redshift	GCS, Amazon S3, Azure ADLS
Best For	Structured analytical queries	Raw, unstructured, semi-structured data
Use Case	Dashboards, attribution, reporting	Model training, clickstream logs, creative assets
Query Style	Interactive, sub-second latency	Batch processing, exploratory analysis

The Lakehouse Approach

Combining both paradigms delivers flexibility across the full spectrum of marketing data while preserving performance for latency-sensitive queries.

Streaming Pipelines: Processing Millions of Events Per Second

1

Apache Kafka as the Central Event Bus

Every customer touchpoint publishes events to Kafka topics. Multiple downstream processors consume simultaneously none blocking the others.

2

Stream Processing Frameworks

Apache Flink, Spark Streaming, Google Dataflow, and AWS Kinesis apply transformations in flight audience scores update continuously, not nightly.

Real-Time Segments

Updates live audience definitions

Personalization

Feeds the recommendation engine

Warehouse Write

Persists events for analytics

Campaign Triggers

Activates optimization workflows

Real-Time Audience Segmentation

The Core Tension

Streaming environments require handling **thousands of profile updates per second** while keeping query latency in single-digit milliseconds. Traditional B-tree indexes, optimized for read-heavy workloads, cannot achieve both simultaneously.



Inverted Indexes

Attribute-based filtering
at scale



Bitmap Indexes

High-cardinality
behavioral flags



Materialized Views

Pre-aggregated
segments for instant
results

DataOps The Discipline Behind the Architecture

The same principles that enable reliable software delivery apply equally to data platforms.

Software DevOps

Version control for code

Automated testing

CI/CD for application code

Observability dashboards

Rollback on failure

DataOps Equivalent

Version control for schemas and pipeline definitions

Data quality validation and regression testing

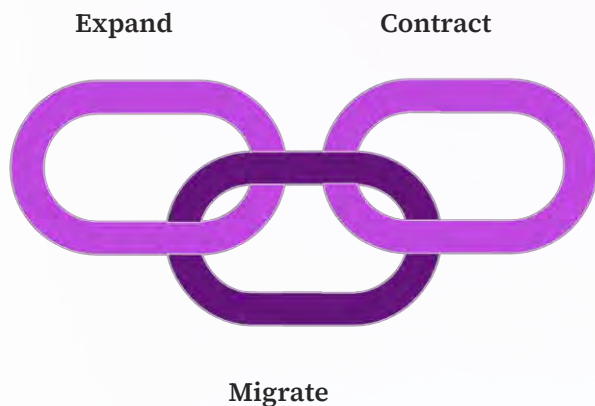
CI/CD for database migrations and pipeline changes

Data pipeline monitoring and alerting

Schema rollback and pipeline recovery

Schema Evolution Without Downtime

In a streaming system, there is no overnight window for schema migrations. Every change must be applied **without breaking downstream consumers**.



The Three Rules of Safe Schema Changes

1

Backward Compatibility

New schema versions must be readable by old consumers. Adding optional fields is safe; removing or renaming fields is not.

2

Forward Compatibility

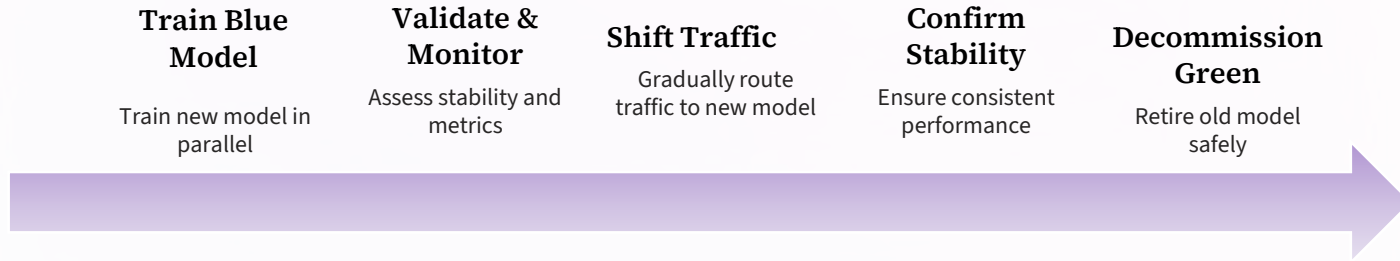
Old schema versions must be readable by new consumers. Consumers must gracefully ignore unknown fields.

3

Full Compatibility

The gold standard: changes are both backward and forward compatible, enabling independent deployment of producers and consumers.

Zero-Downtime Model Retraining



AutoML Integration Goes Further

- Automates feature engineering, model selection, and hyperparameter tuning
- Retraining triggered automatically when data drift is detected
- No manual intervention required for routine retraining cycles
- Drift-to-deployment time reduced from **weeks to hours**

Automated Scaling: Proactive, Not Reactive

Reactive auto-scaling that triggers after load increases is insufficient. Marketing analytics workloads demand **proactive, predictive scaling**.

Leading Indicator

Event ingestion rate signals incoming query load *before* it arrives at the warehouse layer

Pre-Emptive Trigger

Warehouse query layer scales when ingestion exceeds threshold ahead of the spike

Consistent Latency

Query performance remains stable through peak periods, not just normal load

Observability for Streaming Systems

A batch job either succeeds or fails. A streaming pipeline **degrades gradually** consumer lag builds, data quality erodes, latency creeps upward often long before a threshold is breached.

Consumer Lag

Per-topic lag on every Kafka consumer — are processors keeping pace with ingestion?

P99 Query Latency

Tail latency on critical dashboard queries reveals warehouse pressure

Schema Failure Rates

Upstream data quality signals caught before they corrupt downstream models

Model Drift Metrics

Continuous monitoring of prediction distributions signals retraining needs

Data Consistency in Distributed Systems

Why Selective Consistency?

Applying global strong consistency across all workloads imposes unnecessary performance overhead. In distributed streaming systems, not every data point carries the same business risk a missed click event is far less costly than a duplicate purchase record. Deliberately scoping consistency guarantees to where they are actually required delivers both **correctness and throughput**.

Consistency Tiers in Practice

1

Strong Consistency

Required for financial transactions, purchase attribution, and deduplication. Every node sees the same data at the same time. Higher latency cost is justified.

2

Eventual Consistency

Appropriate for behavioral event streams, click tracking, and audience scoring. Temporary divergence is acceptable; the system converges over time.

3

Causal Consistency

Used where ordering matters but strict synchronization does not e.g., session stitching and funnel reconstruction. Preserves cause-and-effect without global coordination overhead.

Cost Optimization

Infrastructure Level

Spot and preemptible instances for batch workloads where interruption is tolerable significant savings at scale

Query Level

Automated query cost controls, result caching, and materialized views prevent expensive ad-hoc queries from driving unexpected spikes

Data Level

Tiered storage policies automatically migrate data from high-performance to low-cost storage based on access frequency

Federated Learning and Privacy-Aware Analytics

How It Works

Model training occurs **locally at the data source** within a browser, on a device, or within a partner's environment. Only model gradients, not raw customer data, are transmitted to a central aggregator.

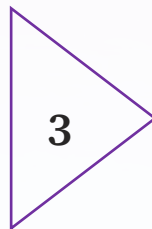
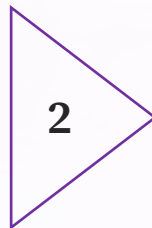
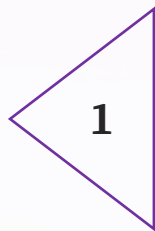
Marketing Use Cases

- Cross-publisher audience modeling without sharing raw user data
- Collaborative filtering across retail partners
- On-device personalization within privacy-compliant boundaries

Three Implementation Paths

Complete Cloud-Native Transformation

Start fresh with event-driven architecture and CI/CD-managed infrastructure. Maximum long-term coherence. Best for new platform initiatives unconstrained by legacy systems.



Hybrid Migration

Maintain batch pipelines for stable workloads while introducing streaming for latency-sensitive use cases. Builds expertise incrementally. **Where most enterprises currently operate.**

Incremental Modernization

Introduce cloud-native components into existing on-premise architecture without full migration. Minimizes disruption but the batch bottleneck often persists at the boundary.

Lessons from Production Deployments

What Consistently Worked

- Infrastructure-as-code from day one version-controlled, not manually provisioned
- Data quality infrastructure built *before* analytical models bad data propagates instantly in real time
- Containerized stream processing with automated rollback on health check failure
- Schema registries with compatibility enforcement preventing silent upstream breakage

What Failed

- **Underestimating operational complexity:** Streaming failure modes appear only under peak load. Stress test at several times expected peak before launch.
- **Ignoring cost observability:** Integrate cost metrics and automated budget alerts into operational dashboards from the start.
- **Schema changes without coordination:** A schema registry with compatibility checks is non-negotiable — not optional.

The Path Forward

Cloud-native data architectures and DevOps together represent a **qualitative change** in marketing analytics capability not an incremental improvement, but a fundamentally different operational model.

01

Data Foundation

Deliberate combination of data lakes and warehouses the Lakehouse approach

03

Mature CI/CD

Schema evolution and zero-downtime model retraining as standard practice

02

Real-Time Indexing

Strategies purpose-built for dynamic segmentation workloads

04

Operational Discipline

Continuous observability, consistency, and cost management across distributed streams

Organizations that build reliable, scalable, high-performance real-time analytics platforms will operate with a **fundamental marketing intelligence advantage** over those still waiting for overnight batch jobs to complete.

Thank You...!

Kamini Murugaboopathy
Wonderbow Analytics Private Ltd., India
Conf42 Database DevOps 2026