



Enterprise Data Lakehouse for Retail & CPG

Executive Summary & Challenges

Modern retail and consumer goods organizations often struggle to scale AI-driven personalization and maintain efficient supply chains due to fragmented legacy systems. Before adopting a unified architecture, companies typically face several critical bottlenecks:

Data Silos:

Customer data lives in CRM systems, transactional data in relational databases, and clickstream data in disparate cloud storage, preventing a unified "Customer 360" view.

Stale Insights:

Traditional data warehouses struggle to process high-velocity e-commerce clickstream and point-of-sale (POS) data, leading to delayed reporting.

Inefficient AI/ML:

Data scientists spend 80% of their time on data engineering and pipeline maintenance rather than building and deploying predictive models.

Supply Chain Volatility:

The inability to accurately forecast demand leads to overstocking of slow-moving items and stockouts of trending products.

Strategic Use Cases

Real-Time Customer 360 & Personalization

Mechanism: Ingests web clickstream data (via Apache Kafka) and merges it with historical purchase data in real-time.

Action: Powers the real-time recommendation engine; if a user abandons a cart, the system instantly triggers personalized emails with dynamic pricing incentives.

Core Platform Capabilities

The solution leverages the Databricks platform to build a governed, scalable architecture:

Delta Lake:

Provides ACID transactions and scalable metadata handling. It ensures reliable streaming and batch data processing for inventory updates and order fulfillment.

Unity Catalog:

Delivers unified governance and security across all data, analytics, and AI assets. This ensures compliance with data privacy regulations (GDPR, CCPA) by managing fine-grained access controls.

Databricks SQL:

Empowers business analysts and merchandisers to run high-performance ad-hoc queries and build dashboards directly on the data lake without moving data.

MLflow:

Manages the complete machine learning lifecycle, from tracking experiments for demand forecasting models to deploying real-time recommendation engines.

Databricks Mosaic AI:

Enables the fine-tuning of large language models (LLMs) to power generative AI customer service chatbots and automated product description generation.

Demand Forecasting & Inventory Optimization

Mechanism: Machine learning models ingest historical sales data, promotional calendars, weather forecasts, and social media trends to predict product demand at the SKU level.

Action: Automatically adjusts reorder points and dynamically allocates inventory across regional fulfillment centers to minimize delivery times.

Dynamic Pricing

Mechanism: Continuously processes competitor pricing data, current inventory levels, and consumer demand elasticity.

Action: Algorithms automatically adjust pricing on the e-commerce storefront to maximize margins or clear out aging inventory.

Generative AI for Customer Support

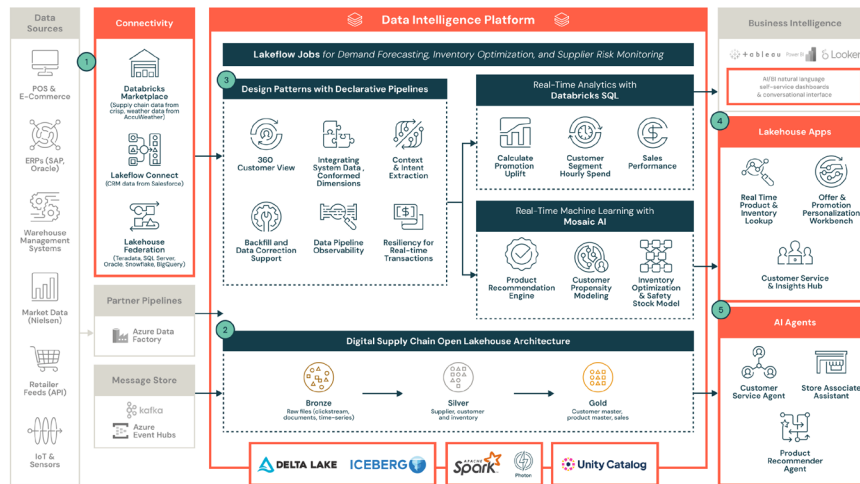
Mechanism: Uses enterprise-specific data (product manuals, return policies, customer history) grounded via Retrieval-Augmented Generation (RAG).

Action: Powers an intelligent virtual assistant capable of handling complex customer queries and assists merchandisers in drafting SEO-optimized product copy.

Reference Architecture

The platform follows a Medallion Data Pipeline to transform raw, fragmented retail data into high-value insights.

Databricks Lakehouse for Retail, CPG & E-Commerce reference Architecture



Architectural Component Breakdown



Ingestion:

Real-time streaming (IoT devices in warehouses, website clickstreams) and batch ingestion (ERP systems, ad network APIs) flow into the platform.



Storage (Delta Lake):

Bronze: Raw, immutable data history.

Silver: Cleaned, filtered, and augmented data (e.g., unified customer IDs).

Gold: Business-level aggregates ready for BI (e.g., daily sales by region, customer churn risk scores).



Serving & Consumption:

BI & Analytics: Power BI / Tableau connected via Databricks SQL endpoints.

Machine Learning: Model serving endpoints hit by the e-commerce backend for real-time recommendations.

Measurable Business Impact

Organizations adopting this unified lakehouse architecture typically achieve significant operational improvements:



30% Reduction

in overall data infrastructure costs by migrating from legacy data warehouses to a unified Lakehouse.



4x Faster

time-to-insight for daily merchandising reports.



20% Decrease

in supply chain waste due to highly accurate ML-driven demand forecasting.



15% Increase

in customer retention through hyper-personalized marketing campaigns.