

E-ISSN: 0976-4844 • Website: www.ijaidr.com • Email: editor@ijaidr.com

Designing Scalable Streaming Data Pipelines with Apache Kafka Schema Enforcement, Real-Time Cleansing, and Event-Driven RAG Patterns

Saurabh Atri

satri@ieee.org, srbwin@gmail.com

Abstract:

Modern data products depend on low-latency, trustworthy streams that can evolve without breaking downstream applications. This article presents a practical blueprint for building scalable streaming data pipelines on Apache Kafka [1]. We focus on three pillars: (1) schema enforcement using a central registry and compatibility policies [2-4]; (2) real-time cleansing and enrichment with stateless and stateful operators on Kafka Streams or Apache Flink [5,6]; and (3) event-driven Retrieval-Augmented Generation (RAG) patterns where model inference is triggered by events and grounded in fresh, streamed context [11]. We provide reference architecture, configuration examples, correctness and cost metrics, and operational playbooks to reach predictable performance.

Keywords:

Apache Kafka, Schema Registry, Avro, Protobuf, Kafka Streams, Apache Flink, Data Quality, Streaming ETL, RAG, Vector Index, Event-Driven Architectures.

1. Introduction

Data producers and consumers change at different speeds. Without strong contracts and online hygiene, pipelines accumulate technical and data debt: schema drifts, null floods, and out-of-order records. Kafka offers durable, scalable logs, exactly-once semantics, and backpressure-friendly consumers [1,7], but teams still need disciplined schema enforcement and cleansing to keep streams reliable. At the same time, applications increasingly require retrieval-augmented LLMs that react to events in near real time [11]. This article lays out a cohesive, production-friendly approach that addresses these needs end to end.

2. Reference Architecture

The reference stack separates the control plane (schemas, policies) from the data plane (topics, processors). Producers publish Avro/Protobuf messages that are registered in a Schema Registry with compatibility guarantees [2-4]. Kafka Streams or Flink jobs cleanse, deduplicate, and enrich events [5,6], writing trustworthy facts to curated topics and serving feature views. Common sources include application SDKs, IoT, and change data capture (CDC) streams such as Debezium [9]. An event-driven RAG service consumes curated facts, refreshes embeddings in a vector store (e.g., FAISS/HNSW/ScaNN) [12-14], and executes retrieval and generation when trigger events arrive. Observability and governance run across all layers.



E-ISSN: 0976-4844 • Website: www.ijaidr.com • Email: editor@ijaidr.com

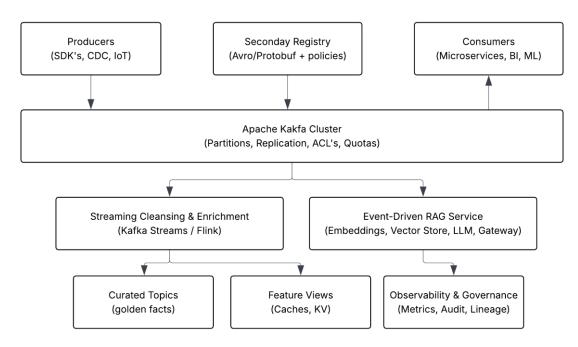


Figure 1. Reference architecture for Kafka-based streaming with schema enforcement, cleansing, and event-driven RAG.

3. Schema Enforcement

A registry assigns immutable IDs to schema versions and enforces compatibility so producers and consumers can evolve independently [2].

Compatibility modes include Backward, Forward, and Full; choose based on reader/writer evolution needs [2].

Treat schemas as code in version control; require defaults for new optional fields; avoid destructive changes on hot topics [2,3].

Validate on write via serializers that check payloads against registered schemas; apply broker quotas and ACLs [1-4].

```
Example: Avro schema with a forward-compatible addition [3].

{
    "type": "record",
    "name": "Payment",
    "namespace": "io.example.payments",
    "fields": [
        {"name": "payment_id", "type": "string"},
        {"name": "amount", "type": "double"},
        {"name": "currency", "type": "string"},
        {"name": "ts", "type": {"type":"long", "logicalType":"timestamp-millis"}},
        {"name": "merchant_category", "type": ["null", "string"], "default": null}
    ]
}
```

4. Real-Time Cleansing and Enrichment

Stateless operators: normalize casing/encodings, standardize units, and route bad records to a dead-letter topic [5,6].



E-ISSN: 0976-4844 • Website: www.ijaidr.com • Email: editor@ijaidr.com

Stateful operators: deduplicate with windows, reorder using grace periods, and compute aggregates; use RocksDB state stores with Streams or Flink backends [5,6,15].

Enrichment: joins and KV caches for low-latency lookups [5,6].

PII: prefer tokenization; use format-preserving encryption only when the format must be preserved.

Kafka Streams example: deduplication with a time window [5].

KStream<String, Payment> payments = builder.stream("raw.payments");

KStream<String, Payment> deduped = payments

 $.selectKey((k,v) \rightarrow v.getPaymentId())$

.transformValues(() -> new DedupTransformer(Duration.ofMinutes(10), "seen-store")); deduped.to("curated.payments");

5. Event-Driven RAG Patterns

Index facts as they change; trigger retrieval only on relevant events to reduce staleness and compute waste [11].

Indexing: chunk and embed curated facts; upsert to a vector store (FAISS/HNSW/ScaNN) [12-14]. Serving: on triggers, retrieve top-k context and generate grounded answers [11].

Design events with routing keys, TTLs, and quality tags; separate raw from curated facts.

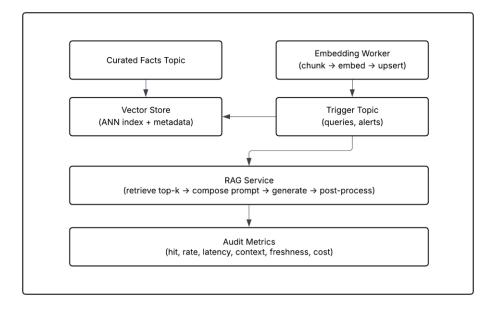


Figure 2. Event-driven RAG: indexing curated facts, then retrieving context on demand when trigger events arrive.

6. Operations: Scaling, Latency, and Reliability

Partitions: scale by partitions and consumer instances; keep keys stable; use sticky partitioner when keys are sparse [1,8].

Latency budgets: monitor end-to-end p50/p95; split budgets across embed/retrieve/generate for RAG [11-14]; size consumer pools for bursts [1].

Exactly-once semantics: idempotent producers and transactions (EOS v2) or deterministic deduplication [1,5,7].

Failure domains: rack-aware replication, multi-AZ, quotas to contain noisy neighbors [1]. Example latency budget for a curated stream and RAG call.



E-ISSN: 0976-4844 • Website: www.ijaidr.com • Email: editor@ijaidr.com

Stage	Budget (p95)	Notes
Producer → Broker	20 ms	Local AZ, batching enabled
Broker → Cleansing	40 ms	1-2 partitions per core
Cleansing → Curated	60 ms	State store hits warmed
Curated \rightarrow Embed Upsert	150 ms	Batch 16-64 chunks
Trigger → Retrieve	80 ms	Top-k ANN within scope
RAG Generate	300-800 ms	Token limits, cache hits

7. Security, Governance, and Cost Control

Enforce authN/authZ, topic ACLs, and schema approvals; version runtime policy and keep audit trails [1-4].

Encrypt in transit and at rest; store secrets in a vault [1].

Track lineage and decisions; for RAG, log context and citations for trust/debugging [11].

Control costs with right-sized retention, tiered storage, and compacted topics for dimensions [1].

8. Evaluation and Benchmarks

Correctness: schema-validation pass rate, duplicate rate after dedup, enrichment hit rate [5,6].

Performance: throughput/partition, end-to-end latency, consumer lag, failover recovery [1].

RAG: context freshness, retrieval hit rate, grounding quality, cost per answer [11-14].

Load test with replay traffic; sweep partitions, batches, and consumer counts; test restarts and hot-producer spikes [1].

9. Implementation Checklist

- Register schemas with compatibility policies [2-4].
- Validate on write; reject incompatible payloads [2-4].
- Package cleansing operators; reuse across jobs [5,6].
- Separate raw, refined, curated topics; RAG reads curated only [11-14].
- Size partitions and apply sticky partitioning if needed [1,8].
- Instrument latency, lag, dedup hit rate, DLQ volume [1].
- Enforce ACLs/TLS/secret management [1-4].
- Game day: broker failover, backlog catch-up, schema evolution [1,2].

10. Conclusion

Kafka moves events at scale; contracts and hygiene make them dependable. Schema-first design, real-time cleansing, and event-driven RAG let one backbone serve analytics and AI with predictable cost and reliability.

REFERENCES:

- [1] Apache Kafka Documentation. https://kafka.apache.org/documentation/
- [2] Confluent Schema Registry Documentation. https://docs.confluent.io/platform/current/schema-registry/index.html
- [3] Apache Avro Specification. https://avro.apache.org/docs/
- [4] Protocol Buffers (Protobuf) Documentation. https://protobuf.dev/
- [5] Kafka Streams Documentation. https://kafka.apache.org/documentation/streams/
- [6] Apache Flink Documentation. https://nightlies.apache.org/flink/
- [7] Apache Kafka KIP-98: Exactly Once Semantics. https://cwiki.apache.org/confluence/display/KAFKA/KIP-98%3A+Exactly+Once+Semantics



E-ISSN: 0976-4844 • Website: www.ijaidr.com • Email: editor@ijaidr.com

- [8] Apache Kafka KIP-480: Sticky Partitioner. https://cwiki.apache.org/confluence/display/KAFKA/KIP-480%3A+Sticky+Partitioner
- [9] Debezium Documentation (Change Data Capture). https://debezium.io/documentation/
- [10] RocksDB: A Persistent Key-Value Store for Fast Storage. https://rocksdb.org/
- [11] Patrick Lewis et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP. NeurIPS 2020.
- [12] Jeff Johnson, Matthijs Douze, Hervé Jégou. Billion-scale similarity search with GPUs. IEEE T-BD, 2019. (FAISS)
- [13] Yu A. Malkov, Dmitry A. Yashunin. Efficient and robust approximate nearest neighbor search using HNSW. TPAMI, 2020.
- [14] ScaNN: Efficient Vector Similarity Search at Scale. Google Research, 2020.
- [15] Kafka Streams State Stores (RocksDB). https://kafka.apache.org/28/documentation/streams/developer-guide/state.html