**Event Generator**
Faker + Avro

string, enum
s: world
+: quantum

**Schema Registry**

processValues

**Kafka Cluster**

Topic: shopstream-purchases
3 Partitions, RF = 3

Broker 1
:9094

Broker 2
:9095

Broker 3
:9096

**Stream Processing**

Consumer Group
Regional Batching

parquet

**Data Lake**

MINIO

- Real-world e-commerce streaming pipeline simulation

- Demonstrates complete data journey from generation to storage

- Multi-technology integration (Kafka, Schema Registry, Object Storage)

- Event generation with realistic data
- Schema validation at entry point
- Distributed storage across 3 brokers
- Analytics-optimized storage

---

**Docker Network**

**Kafka Cluster**

Schema Registry
Port: 8082
Schema Management

kafka-1
Broker ID:1
Port: 9094
Controller + Broker

kafka-2
Broker ID:2
Port: 9095
Controller + Broker

kafka-3
Broker ID:3
Port: 9096
Controller + Broker

kafka UI
Port: 8080
Monitoring

MinIO
Port: 9000/9001
Object Storage

---

**Event Generation**

US Events
user_id: 1234, 5678

EU Events
user_id: 2345, 6789

ASIA Events
user_id: 3456, 7890

**Kafka Partitioning by**

Partition 0
hash(user_id) % 3 = 0
Mixed: US, EU, ASIA

Partition 1
hash(user_id) % 3 = 1
Mixed: EU, ASIA, US

Partition 2
hash(user_id) % 3 = 2
Mixed: ASIA, US, EU

**Consumer Batching**

US Batch
All US events
from all partitions

EU Batch
All EU events
from all partitions

ASIA Batch
All ASIA events
from all partitions

**Data Lake Structure**

region=US/
date=2025-10-24/
purchases_*.parquet

region=EU/
date=2025-10-24/
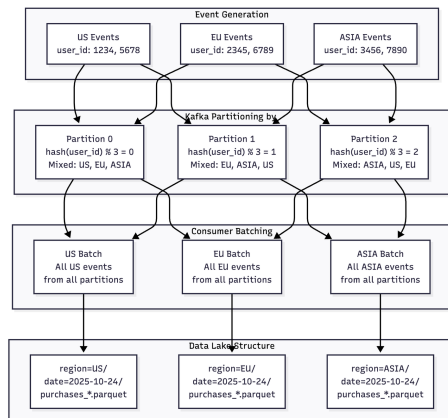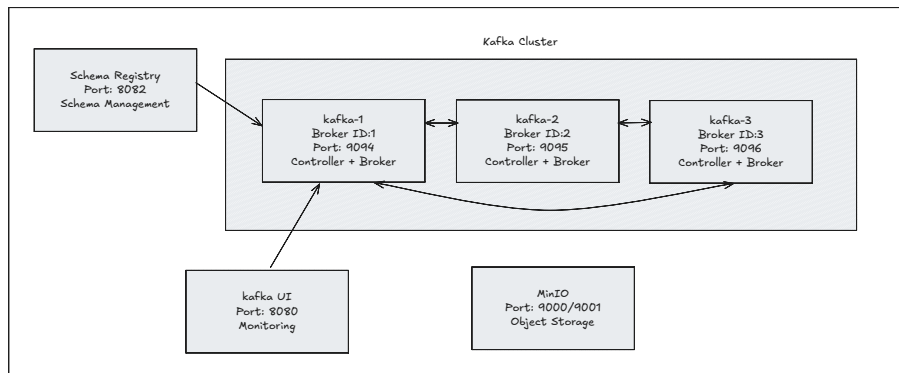purchases_*.parquet

region=ASIA/
date=2025-10-24/
purchases_*.parquet

**Event Generation**

Three regions generate purchase events: US, EU, ASIA
Message key = user_id (not region) for ordering guarantees
Avro serialization with Schema Registry validation
Realistic data using Faker library

**Kafka Partitioning**

hash(user_id) % 3 determines partition assignment
All regions mixed across all 3 partitions
Same user always goes to same partition
Replication factor 3 = every message on every broker

**Consumer Processing**

Single consumer group reads from all partitions
Events reorganized by region field in message
Separate in-memory batches: US, EU, ASIA
Manual offset commits after successful storage

**Batch Triggers**

Size trigger: 50 events per batch
Timeout trigger: 30 seconds maximum wait
Balances efficiency with data freshness
Prevents memory overflow and stale data

**Stream-to-Batch Conversion**

Pandas DataFrame → PyArrow Table → Parquet
Columnar format optimized for analytics
Snappy compression for storage efficiency
In-memory processing without disk I/O

**Data Lake Organization**

Hive partitioning: region=US/date=2025-10-24/
Business logic grouping vs technical distribution
Enables partition pruning for efficient queries
Bronze layer for raw ingested data