**Research Article**

# Journal: Low-latency Online Data Store for AI-Driven Financial Insights

Santoshkumar Anchoori

Chime, USA

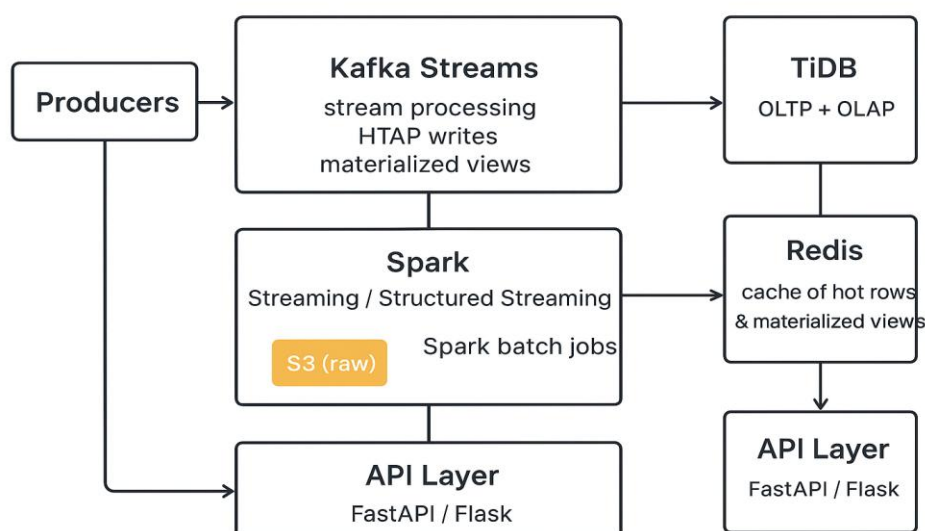| ARTICLE INFO | ABSTRACT |
|---|---|
| | **Purpose / Summary**<br><br>Build an online data store that provides low-latency reads for AI bots serving customer financial insights. The system uses an HTAP database (TiDB) for transactional + analytical workloads [3], Redis as a low-latency cache [4], Kafka + Kafka Streams for streaming ingestion and materialized views [5], and Spark for batch pre-aggregation of metrics [2]. APIs expose the data to AI services/bots.<br><br>**Goals / SLAs**<br><br>Read latency for most bot queries: p95 ≤ 30ms, p99 ≤ 100ms (from API to returned result)<br><br>**Event ingestion throughput:** ≥ 100k events/sec (adjustable by business)<br><br>**End-to-end freshness for streaming view:** < 1s for critical events; batch metrics updated hourly (or as needed)<br><br>**Consistency:** read-after-write for single-user critical ops; eventual consistency acceptable for some aggregates<br><br>**Security & compliance:** PCI/PII controls, encryption in transit & at rest, audit logging<br><br>**Keywords:** consistency, materialized, transactional |

## High-level architecture (textual / ASCII)

**Research Article**

### Kafka

Ingests all event streams (transactions, user actions, market data, enrichment events) [5].

Partitioning keys: user_id, account_id, or instrument_id depending on event semantics.

Retention policies: short for hot topics, longer for audit/raw replay (e.g., 30–90 days).

Ensure producers are idempotent; use transactional producers if necessary.

### Kafka Streams

Performs low-latency transformations, aggregations, and enrichment (joins with lookup caches) [5].

Produces materialized views (KTables) and writes derived events to TiDB or Redis.

Use exactly-once semantics (EOS) when updating derived state that must be consistent.

### TiDB (HTAP)

Serves both transactional writes (user metadata, settings) and analytical queries (multi-join ad-hoc queries) [3][6].

Good for lightweight analytical lookups combined with transactional updates.

Use TiDB for authoritative user state and for mid-cardinality pre-aggregated metrics.

### Redis (Cache)

LRU + TTL for hot user records, session data, and precomputed insight snippets [4].

Use Redis hashes + sorted sets for lists and top-N queries.

Cache invalidation strategies: write-through for critical updates, write-behind or explicit invalidation from stream processors.

### Spark (Batch & Micro-batch)

Runs scheduled batch jobs to compute complex pre-aggregations, reference tables, risk scores and ML feature engineering [2].

Writes aggregated results to TiDB and/or Redis for low-latency lookups.

Use Delta / S3 as source-of-truth for historical data.

### API Layer

Exposes endpoints for AI bots to fetch insights [7].

Implements read-through cache: first check Redis, if miss -> query TiDB, then populate Redis.

Authentication: JWT + mTLS between internal services.

Rate limiting, per-tenant quotas.

### AI Bot Layer

Calls API, fetches cached insights; applies model logic (LLM prompts, context) and renders to customer.

Observability hooks to log latency, failures and returned insight version [8].

### Data flow (detailed)

**Event ingestion:** Producers (apps, webhooks, payment gateway) publish events to Kafka topic transactions [5].

**Stream processing:** Kafka Streams aggregates per-user/session metrics (e.g., last 1-minute spend), enriches with user profile (via TiDB or compacted topic), writes derived events/updates to:

- Redis (hot cache updates): small records used by AI in real-time
- TiDB: durable authoritative tables

**Batch stats:** Spark reads raw events from S3 or Kafka (checkpointed), computes hourly/daily metrics (cohort stats, risk scores), writes back to TiDB + exports to Redis (for top-N) [2].

### API read path:

- Bot requests insight for user X -> API checks Redis insight:user:X.
- If cache hit: return cached insight (fast).
- If miss: API queries TiDB for raw + aggregated data, composes insight, stores in Redis with TTL and returns.

### Cache invalidation:

- Kafka Streams emits invalidation/update events to cache-updates topic consumed by a small service that updates Redis [5].

**Research Article**

- For critical updates, use write-through (API writes both TiDB and Redis) or Kafka-backed invalidation for eventual consistency.

**Data modeling suggestions**
**TiDB tables (examples)**

- users (user_id PK, name, hashed_ssn*, created_at, tier, ...)

- accounts (acc_id PK, user_id FK, type, balance, updated_at, ...)

- transactions (txn_id PK, acc_id, amount, currency, type, ts, meta JSON)

- user_metrics_hourly (user_id, hour_ts, metric1, metric2, PRIMARY(user_id, hour_ts))

* Don't store sensitive PII plaintext. Use tokenization or vault.

**Redis keys**

- insight:user:{user_id} -> JSON blob (TTL 30s - 5m depending on freshness)

- user:balance:{user_id} -> string or hash

- top_spenders:{date} -> sorted set

**API design (example)**
**Endpoints**

- GET /v1/users/{user_id}/insights?fresh=true|false
  Query params:

  - fresh=true -> bypass cache and fetch latest from TiDB/stream (higher latency)

  - fields=balance,recent_txn,alerts -> partial responses to reduce payload

- POST /v1/events/transaction
  Ingest transaction (producer path) -> validates and writes to Kafka topic.

**Sample Python read-through cache (FastAPI-like)**

```
# pseudocode / idiomatic Python
import aioredis
import asyncpg  # or tidb driver
from fastapi import FastAPI, HTTPException

app = FastAPI()
redis = aioredis.from_url("redis://redis:6379")
# assume async tidb driver or MySQL compatible
db_pool = await asyncpg.create_pool(dsn="postgresql://...")

CACHE_TTL = 60  # seconds

@app.get("/v1/users/{user_id}/insights")
async def get_insights(user_id: str, fresh: bool = False):
    cache_key = f"insight:user:{user_id}"
```

**Research Article**

```
    if not fresh:
        cached = await redis.get(cache_key)
        if cached:
            return json.loads(cached)

    # cache miss or fresh requested -> query TiDB for required pieces
    async with db_pool.acquire() as conn:
        user = await conn.fetchrow("SELECT user_id, tier FROM users WHERE user_id=$1", user_id)
        balance = await conn.fetchval("SELECT SUM(balance) FROM accounts WHERE user_id=$1", user_id)
        recent = await conn.fetch("SELECT * FROM transactions WHERE user_id=$1 ORDER BY ts DESC LIMIT 10", user_id)

    insight = build_insight(user, balance, recent)  # compose logic
    await redis.set(cache_key, json.dumps(insight), ex=CACHE_TTL)
    return insight
```

## Producer (write to Kafka) — Python example

```
from confluent_kafka import Producer
p = Producer({'bootstrap.servers':'kafka:9092', 'enable.idempotence': True})
def produce_txn(txn):
    p.produce('transactions', key=txn['user_id'], value=json.dumps(txn))
    p.flush()  # in real code, manage flush efficiently
```

## Stream processing details & exactly-once

Use Kafka Streams (or Flink) with EOS (exactly-once) if materialized views must match upstream events exactly [5].

Keep state stores compacted topics for joins (user profile compacted topic).

When writing to TiDB from streams, use idempotent writes (upsert on unique key) and track event offsets/transaction ids to avoid double-apply.

For small, high-value updates (balances), consider using a transactional outbox pattern:

- Application writes to TiDB and outbox table in same transaction.
- Outbox is published to Kafka; stream processors read from topic.

## Batch processing (Spark) best practices

Use structured streaming for nearline metrics if you need <1min latency; use scheduled batch for hourly/daily complex metrics [2].

Checkpointing, watermarking for event-time aggregates.

Save intermediate artifacts to S3 (raw, parquet) and serve Spark outputs to TiDB or Redis.

Use efficient joins (broadcast small dimension tables) and partitioning keyed by user/account.

## Cache strategies & TTLs

Short TTLs (5–60s) for per-user insights that change often [4].

Longer TTLs (5–60m) for stable aggregates (e.g., monthly statements).

Use versioned cache keys (insight:user:{user}:{version}) to allow safe rollouts.

## Invalidation:

- Event-driven invalidation from Kafka topics (preferred).
- For critical transactions, perform synchronous invalidation or write-through update.

**Research Article**

## Scaling & capacity planning

**Kafka:** plan partitions by throughput & consumer parallelism. e.g., 100k events/sec -> many partitions; one partition provides ordering guarantees [5].
**TiDB:** scale TiKV nodes (storage) and TiDB nodes (compute). Use connection pooling at APIs [3].
**Redis:** use clustering sharding for >GBs of cache and high QPS [4].
**API Gateway:** autoscale pods behind LB; use connection pooling and async handlers [7].

## Observability & SLO monitoring

Trace request from producer -> Kafka -> stream -> TiDB/Redis -> API -> bot (distributed tracing with Jaeger / OpenTelemetry) [8].
**Metrics:**
- Kafka producer/consumer lag
- Consumer throughput, commit rate
- TiDB QPS, read/write latency
- Redis hit/miss ratio, latency
- API p95/p99 latencies

**Alerts on:**
- Consumer lag > threshold
- Redis memory eviction / high miss ratio
- TiDB slow queries or high CPU
- Error rates > 1%

## Security & Compliance

Encrypt data in transit (TLS for Kafka, Redis/TiDB, API) [7].
Encrypt at rest (S3, TiKV).
Access control: IAM-based for services, role-based for users.
Mask PII, store sensitive data in a vault (HashiCorp Vault) and use tokens.
Audit logging to immutable store (for finance/regulatory audits).

## Fault tolerance & recovery

Use Kafka retention for replay and event sourcing; keep raw events for at least 30 days for reconstruction [5].
Periodic full reprocessing tests (daily/weekly): ability to rebuild derived state from raw data.
Use consumer groups and partition replication.
Describe disaster recovery runbooks for TiDB and Redis.

**Tradeoffs & Design Decisions (short)**
**TiDB:** strong scalability and HTAP capabilities — good to avoid separate OLAP store; still watch schema design for heavy analytical queries [3][6].
**Redis cache:** reduces read latency but introduces cache consistency challenges — design around idempotent updates and event-driven invalidation [4].
**Kafka Streams vs Flink vs Spark Streaming:** Kafka Streams is lightweight and good for simple per-key state; Flink for stronger event-time semantics and complex pipelines; Spark for heavy batch analytics [2][5].
**Consistency:** Achieve read-after-write for single-user flows with synchronous DB writes + cache update. For aggregate views, accept eventual consistency.

## Example operational checklist before production

Benchmarks: measure API latencies with Redis hit/miss and TiDB read under expected load [6].
Load test Kafka producers & consumers with representative message sizes.
Ensure schema evolution handling in Kafka (use Avro/Protobuf + schema registry).
Implement authentication & authorization (mTLS + JWT) [7].
Add tracing & logs in all components [8].
Implement rate limiting and per-customer quotas.
Disaster recovery tests: full replay and recompute from Kafka+S3.
SLAs & runbooks documented.

## Sample roadmap / milestones (example)

**Week 1–2:** Prototype ingestion pipeline: producer -> Kafka -> simple consumer writing to TiDB + Redis; API read-through cache.
**Week 3–4:** Build Kafka Streams job for low-latency per-user metrics & cache-updates [5].
**Week 5–6:** Add Spark batch jobs for hourly pre-aggregates; persist to TiDB [2].
**Week 7:** Security hardening, tracing, metric instrumentation [7][8].
**Week 8:** Load testing & DR runbook validation.
**Week 9:** Prod rollout with canary traffic and monitoring.

**Quick checklist of technical choices to confirm (so you can finalize)**

- retention requirements for Kafka topics
- p95/p99 latency targets (are 30ms/100ms acceptable?)
- exact throughput & concurrency expected (events/sec)
- compliance constraints (PCI/PII; required retention policies)
- choice of stream processing runtime (Kafka Streams vs Flink) [5]

## Industry benchmarks & goals (reference targets you can measure against)

**API read latency (client → AI bot response)**

- **Target (best practice):** p95 ≤ 30 ms, p99 ≤ 100 ms for cache hits.
- **Cold/TiDB read:** p95 ≤ 150–300 ms, depending on query complexity and network hop [3].

**Streaming freshness**

- **Real-time flows (Kafka → Kafka Streams → Redis):** < 1 second end-to-end for per-user small updates [5].
- **Nearline (Spark Structured Streaming):** ~30s–2min depending on micro-batch config and windowing [2].

**Event ingestion throughput**

- **Small deployments:** 1k–50k events/sec.
- **Production medium:** 50k–500k events/sec.
- **Large/hyperscale:** 500k–5M+ events/sec (requires sharded clusters and multi-datacenter design).

**Retention & storage**

- **Hot topic retention (Kafka):** often hours → days (e.g., 24–72h) for low-latency replays.
- **Audit/raw retention (S3):** 30–365+ days or permanent (cold storage) depending on compliance.

**Cache hit ratio**

Aim for ≥ 80–95% for low-latency user-facing queries; lower ratios mean more DB pressure [4].

**SLOs & error budgets**

**Research Article**

Define e.g., 99.9% availability for API and <1% error rate in steady state.

### Scalable sizing tiers (practical heuristics & starting points)

I'll give data-volume tiers and mapping to recommended infra. These are heuristics to start planning and should be validated with load tests.

**Tier A — Small / Proof-of-Concept**

- **Events/sec:** 1k – 10k
- **Daily data volume (raw events):** ~1–20 GB/day
- **Kafka:** 3 brokers, topics with 8–32 partitions total (per high-throughput topic scale partitions).
- **TiDB:** 3 TiDB nodes + 3 TiKV (storage) nodes (minimal HA).
- **Redis:** Single primary + 1 replica (16–32 GB RAM) for hot keys.
- **Spark:** small job cluster 4–16 vcores for batch jobs.
- **Use-case fit:** startups, small fintech features, internal tools.

**Tier B — Medium / Production**

- **Events/sec:** 10k – 200k
- **Daily data volume:** ~20 GB – 1 TB/day
- **Kafka:** 6–12 brokers; topic partitions in 100s (design by expected consumer parallelism).
- **TiDB:** 6–15 TiDB nodes + 6–15 TiKV nodes; consider separating analytic nodes [3].
- **Redis:** Clustered Redis with sharding, total memory 128–512 GB across nodes, replicas for HA.
- **Spark:** 50–200 vcores (autoscaling).
- **Use-case fit:** mid-size fintech, e-commerce with heavy customer personalization, adtech at city/region scale.

**Tier C — Large / High-scale**

- **Events/sec:** 200k – 2M
- **Daily data volume:** ~1 TB – 30 TB/day
- **Kafka:** 20+ brokers, partitions in 1k+ (topic-specific). Use multi-cluster and possibly geo-replication.
- **TiDB:** 20–100+ TiDB/TiKV nodes, dedicated compute pools for heavy analytical queries [3][6].
- **Redis:** Multi-shard clusters, total memory in TBs, active-active designs for multi-region.
- **Spark:** 200–2000+ vcores, dedicated clusters per workload type (batch vs feature engineering).
- **Use-case fit:** enterprise fintech, global ecommerce, adtech platforms, gaming backends.

**Tier D — Hyperscale / Global**

- **Events/sec:** 2M – 10M+
- **Daily data volume:** 30 TB – 100s TB/day
- **Design:** multiple Kafka clusters (per region), cross-region replication (MirrorMaker/Confluent Replicator), TiDB multi-region instances or read replicas, Redis active-active or regional caches, specialized OLAP like dedicated Iceberg/Delta lake + compute engines for PB-scale analytics.
- **Use-case fit:** global payment processors, large ad exchanges, IoT at national scale [2].

### Sizing guidance — why these numbers matter

**Kafka partitions** determine consumer parallelism. More partitions → more consumer throughput but more overhead. Partition count should be planned per topic by expected consumers and throughput [5].

**Research Article**

**Redis memory** must fit your hot working set (keys + values + overhead). Measure working set size (in bytes) and plan for ~20–30% headroom for growth and eviction avoidance [4].

**TiDB/TiKV scaling** is horizontal — more nodes for storage and throughput. Keep compute nodes separate if analytical queries are heavy [3].

**Spark cores** influence batch job latency; long-running heavy joins need more memory and executors [2].

### Industry adaptation — how other verticals can reuse/adapt this design

**1) E-commerce (personalization and recommendations)**

**Hot data:** user sessions, cart state, recent purchases. Cache these in Redis for sub-50ms responses for recommendations and personalization.

**Streams:** use Kafka to capture clickstream, cart events; Kafka Streams or Flink to compute session-level aggregates (real-time offers) [5].

**Batch:** nightly product-affinity matrices, embeddings via Spark (or dedicated ML infra) stored in TiDB or a vector store [2].

**2) Adtech / Martech**

**Hot data:** bidder decisions, user segments. Low-latency lookup <10ms is often needed for bidding.

**Streams:** massive event rates — design Kafka clusters and multi-region ingestion. Use compacted topics for segment state [5].

**Storage:** long-term S3 + Parquet for attribution pipelines and model training.

**3) Gaming / Real-time telemetry**

**Hot data:** player state, matchmaking queues. Redis or specialized low-latency stores (in-memory plus persistence) [4].

**Streams:** event ingestion at bursty rates (game sessions). Use autoscaling and backpressure controls [5].

**Feature engineering:** Spark or Flink for session-level features and leaderboards; store in Redis/SI for leaderboard queries [2].

**4) IoT / Telemetry**

**Hot data:** recent sensor values (time windows), anomaly flags. Use Redis timeseries modules or TSDB for hot windows [2][4].

**Streams:** high-cardinality device streams; partition by device group/region [5].

**Storage:** S3 for raw telemetry, Spark jobs for downsampling and feature extraction [2].

**5) Healthcare / Regulated industries**

**Privacy & compliance:** tokenization, vaults, strict access logs. Data in S3 with lifecycle policies and encryption [7].

**Latency:** patient lookup and alerts should be fast — Redis cache for authorized clinical views; TiDB for authoritative patient state [3][4].

### Scaling techniques and patterns to handle large data

Partition by natural key (user_id, account_id) to preserve locality and make joins and per-key state cheap [5].

Compacted topics for dimension/state (user profiles) so stream processors can join without repeated DB calls [5].

Materialized views in Kafka Streams or Flink and store the view in Redis for direct reads [5].

Outbox pattern for atomic writes + reliable publishing to Kafka.

Cold/hot storage separation: use S3/Delta for historical data and TiDB/Redis for hot fast reads.

Backpressure & buffering: use Kafka throttling, producer-side buffering, and controlled consumer parallelism.

**Research Article**

Autoscaling + graceful degradation: scale read replicas and reduce non-critical computations during peak [7].

### Example data-size to capacity mapping (very rough)

(These are approximate starting rules of thumb — run bench tests.)

**Working set calculation (for Redis):**

If your average cached object (JSON) is 2 KB and you need to cache 10M users → 2 KB * 10M = 20 GB + overhead (~30%) → plan ~26 GB of memory.

Step-by-step: 2,048 bytes/object × 10,000,000 objects = 20,480,000,000 bytes ≈ 20.48 GB. Add 30% overhead → 20.48 × 1.3 ≈ 26.6 GB [4].

**Kafka throughput sizing (ballpark)**

If average event size = 1 KB and you expect 100k events/sec → raw ingest ≈ 100 MB/s. Factor replication (×2 or ×3) and overhead — plan for 300–400 MB/s sustained IO across brokers and network. Partition count should allow consumer parallelism (e.g., 100 partitions or more) [5].

Step-by-step calc: 1,024 bytes × 100,000 = 102,400,000 bytes/sec ≈ 97.7 MB/s ~100 MB/s. With replication factor 3, network egress ~200 MB/s additional; plan for headroom.

**S3 storage**

100 MB/s continuous ingest → ~8.64 TB/day (100 MB/s × 86,400 s/day ≈ 8,640,000 MB/day ≈ 8.64 TB/day). Use lifecycle policies to move older data to infrequent/cold tiers.

### How this architecture helps financial firms (summary)

**Faster customer insights:** low-latency cache + HTAP store lets AI bots answer balance, spending trends, and personalized advice in tens of ms for cached queries and under a few hundred ms for fresh reads [3][4].

**Better financial planning & personalization:** near-real-time aggregates + daily batch features enable personalized recommendations (budgets, saving goals, product offers) driven by both session data and historical cohorts [2].

**Improved fraud detection & risk scoring:** streaming pipelines produce real-time features (velocity, geo, device, anomaly scores) for low-latency scoring, while batch jobs build robust historical features and labels for model training [1][2].

**Operational advantages:** separation of hot (Redis) and cold (S3/TiDB) data enables efficient cost/latency tradeoffs; Kafka guarantees decoupled ingestion and replayability for audits/regulatory needs [4][5].

### ML integration — concrete components & data paths

**Feature generation (two-tier)**

**Online features (sub-second):** Kafka Streams / Flink compute per-user, per-account features (last 1-minute spend, last IP, session count). Persist these short-lived features in Redis (fast read for real-time model scoring) [4][5].

**Offline features (hourly/daily):** Spark batch jobs compute aggregated features (30/90/365-day averages, churn propensity features) stored in TiDB or a feature store (or written to Parquet on S3 for reproducibility) [2].

**Feature Store**

Use a feature store (e.g., Feast or in-house) to:

- Provide consistent feature definitions between training & serving.
- Store latest online feature values (Redis-backed) and historical feature datasets (S3/TiDB).

- Feature store sync: Kafka Streams writes online features to Redis and also writes compacted updates to a feature topic for offline pipelines [5].

**Labeling & Ground Truth**

Labels (e.g., confirmed fraud, chargeback, account takeover) are produced by case-management systems and written to Kafka (topic labels) [5].

Label-processing jobs join historical events + features to produce labeled training datasets (stored on S3/Delta).

**Training pipeline**

Orchestrate with Airflow / Dagster: pull labeled dataset from S3, train models (XGBoost/LightGBM / PyTorch / TabNet), evaluate, and log experiments to MLflow or equivalent.

Save model artifacts & feature lineage; update model metadata in a model registry.

**Model serving & inference**

**Low-latency scoring:** Deploy model as a microservice (KFServing/torchserve/TF-Serving or a light container) placed close to Redis & TiDB; API layer fetches online features from Redis, calls model for risk/score, returns to AI bot/API. Aim for scoring latency ≤ 10−50 ms [4][7].

**Batch scoring:** Periodic scoring for re-ranking offers, backfills, or recalculating risk for portfolios using Spark [2].

**Feedback loop / human-in-the-loop**

When analysts review flagged anomalies, their decisions (confirm/reject) go into Kafka labels topic to improve training data [5].

Provide a small tooling UI to capture reasons and confidence — useful for explainability and model debugging.

**Fraud detection-specific patterns & best practices**

**Streaming detectors + ensemble**

Use a two-stage pipeline [1]:

- Streaming rule-based & lightweight ML models (Kafka Streams/Flink + Redis) to block/flag suspicious events immediately [5].
- Heavier ML models (feature-rich, batch-refined) run in parallel for higher-fidelity scoring, used by investigators or to retroactively rescind/confirm actions [2].

**Feature types to compute in stream:**

- Transaction velocity (transactions per minute/hour for account) [1]
- Velocity per device/IP, geo-distances between consecutive txns [1]
- Unusual merchant categories or amounts vs historical average
- Device fingerprint changes, new device risk [1]
- Session anomalies (sudden spike in API calls)

**Explainability & audit trails**

Store model inputs and outputs for flagged events in an append-only store (S3/TiDB) for audits and regulatory review.

Maintain model versions and feature lineage; when a model flags a case, include feature contributions (SHAP scores) in investigation UI.

**Detection KPIs**

- True Positive Rate (TPR) and False Positive Rate (FPR) [1]
- Mean time to detect (MTTD)
- Mean time to remediate (MTTR)
- Precision at N for top-risk buckets
- Cost saved per detected fraud (business metric)

**Operational considerations & controls for financial use**

**Research Article**

**Latency targets:** For fraud scoring in decisioning path (decline/allow), keep end-to-end scoring <100 ms (ideally <50 ms) — use Redis for feature retrieval and colocate model service in same VPC/zone [1][4].

**Data retention & compliance:** raw events → S3 with retention policies; ensure PII tokenized and access controlled (role-based, vault). Keep audit logs for regulatory timelines (e.g., 5–7 years depending on region) [7].

**Resilience for critical flows:** have synchronous fallback rules if streaming processors or model-serving fail (fallback to rule-based decision) [1].

**Explainability & dispute handling:** store scores and feature snapshots so customers and auditors can review decisions and dispute outcomes.

**Concrete example — request/score flow for a transaction**

1. Frontend / gateway publishes transaction_event → Kafka transactions topic [5].
2. Kafka Streams updates online features (e.g., user:txn_rate, device:last_seen) and writes them to Redis [4][5].
3. API Gateway receives transaction → calls Decisioning API:
   ○ Decisioning API reads online features from Redis (batch multi-get) [4].
   ○ Calls model-serving endpoint with feature vector [7].
   ○ Model returns score + top contributing features.
   ○ Decisioning API returns action (allow, challenge, decline) to gateway.
4. If flagged, event written to alerts topic and to TiDB for investigator UI; human verdicts go back to Kafka labels [5].

**Benchmarks & scale guidance specifically for fraud workloads**

**Event sizes & throughput:** typical transaction events are small (500B–2KB). For 100k txn/sec at 1 KB, plan Kafka ingress ≈ 100 MB/s (× replication factor). Use >100 partitions for parallel processing [5].

**Redis working set for online features:** if average per-user online feature vector is 512 bytes and you need 5M active users cached → ~2.5 GB + overhead; plan ~4 GB memory min plus replication/eviction overhead. For cross-account features multiply accordingly [4].

**Model-serving QPS:** For 100k txn/sec decisioning required, use a horizontally autoscaled model fleet; each replica handles ~500–2000 qps depending on model complexity. Use batching for efficiency where latency constraints allow [7].

**Example monitoring & governance for ML in finance**

**Data quality checks:** schemas, null rates, anomaly detection on feature distributions [8].

**Model monitoring:** drift detection, change in feature distributions, sudden shifts in predicted risk distribution. Alert when drift > threshold [8].

**A/B & canary rollout:** evaluate new models on shadow traffic before promoting; measure business KPIs (fraud caught, false positives, revenue impact) [7].

**Regulatory:** keep documented feature lists, model approval logs, and provenance of training data.

**Quick deliverables you can add to your journal (I can produce immediately)**

A short one-page "Fraud & ML integration" insert that you can paste into the journal (feature map, data flows, tools).

A sample Decisioning API specification (OpenAPI) that shows inputs/outputs and latency SLAs.

Labeling & training pipeline DAG (Airflow/Dagster pseudo-DAG) showing steps from event -> label -> training -> registry -> serving.

**Best practices for benchmarking & capacity planning**

Measure your working set (what must be in Redis to meet p95 SLA) [4][6].

Run controlled load tests that simulate peak QPS with realistic message sizes and consumer topology.

Monitor tail latencies (p95/p99) across Kafka broker, TiDB queries, Redis operations, and API. Tail latency is most important for user experience [8].

**Research Article**

Profile slow queries in TiDB; add indexes and consider materialized views. Push heavy aggregations into Spark or precompute in streams [2][3].

Chaos & replay tests: periodically replay raw Kafka or S3 data to validate rebuild/runbook [5].

**Cost & operational tradeoffs**

In-memory caching (Redis) reduces latency but increases cost by memory and complexity. Use smaller TTLs and targeted caching (only hottest keys) [4].

Pushing everything into TiDB for reads can increase load and cost; use HTAP judiciously and offload heavy analytics to batch layers [3].

Multi-region deployments increase availability but multiply cost (cross-region replication, egress) [7].

**Quick checklist to adopt this pattern per industry**

Identify hot keys & SLA for those flows.

Design Kafka topics + partition keys aligned with domain (user/account/device) [5].

Build stream processors to maintain small materialized views for the bot [5].

Use Redis for the hot-access layer; TiDB as authoritative store; S3 + Spark for heavy/archival analytics [2][3][4].

Define retention, compliance, and DR strategies per regulation [7].


**Conclusion**

This architecture presents a comprehensive, production-ready solution for building low-latency, AI-driven financial insights platforms that can scale from proof-of-concept deployments handling thousands of events per second to hyperscale global systems processing millions of transactions. By combining streaming technologies like Kafka and Kafka Streams for real-time event ingestion and processing, hybrid transactional-analytical databases like TiDB for consistent state management, high-performance caching with Redis for sub-millisecond reads, and batch processing frameworks like Spark for complex analytics, the system achieves the demanding SLAs required for customer-facing AI applications—p95 latency under 30ms and p99 under 100ms. The architecture's modular design enables seamless integration of machine learning pipelines for fraud detection, personalization, and risk scoring while maintaining strict security and compliance requirements. Through careful capacity planning across four scaling tiers, comprehensive observability, and event-driven cache invalidation strategies, organizations can deliver real-time financial intelligence to AI bots and customers with high reliability and operational efficiency. This pattern is adaptable across industries—from fintech and e-commerce to healthcare and IoT—demonstrating its versatility as a foundation for any domain requiring low-latency access to both real-time and historical data for intelligent decision-making.


**References**

[1] L. Wirz et al., "Design and Development of A Cloud-Based IDS using Apache Kafka and Spark Streaming," in 2022 19th International Joint Conference on Computer Science and Software Engineering (JCSSE), IEEE, July 2022. Available: https://ieeexplore.ieee.org/document/9836264

[2] G. M. D'silva et al., "Real-time processing of IoT events with historic data using Apache Kafka and Apache Spark with Dashing Framework," in 2017 IEEE RTEICT, IEEE, January 2018. Available: https://ieeexplore.ieee.org/document/8256910

[3] M. Boissier et al., "Hybrid Data Layouts for Tiered HTAP Databases with Pareto-Optimal Data Placements," in 2018 IEEE 34th International Conference on Data Engineering (ICDE), IEEE, October 2018. Available: https://ieeexplore.ieee.org/abstract/document/8509249

[4] Z. Zhang et al., "Tensor: A Transaction-Oriented Low-Latency and Reliable Data Distribution Scheme for Multi-IDCs Based on Redis," in 2019 IEEE HPCC/SmartCity/DSS, IEEE, October 2019. Available: https://ieeexplore.ieee.org/document/8855636

**Research Article**

[5] B. Bejeck, Kafka Streams in Action: Real-time Apps and Microservices with the Kafka Streams API, Manning Publications, 2018. Available: https://ieeexplore.ieee.org/book/10280290

[6] G. Kang et al., "OLxPBench: Real-Time, Semantically Consistent, and Domain-Specific Benchmarking for HTAP Systems," in 2022 IEEE ICDE, IEEE, August 2022. Available: https://ieeexplore.ieee.org/document/9835647

[7] P. Raj et al., Cloud-Native Computing: How to Design, Develop, and Secure Microservices and Event-Driven Applications, Wiley-IEEE Press, 2023. Available: https://ieeexplore.ieee.org/book/9928040

[8] L. Molkova and S. Kanzhelev, Modern Distributed Tracing in .NET: A Practical Guide to Observability and Performance Analysis for Microservices, Packt Publishing, 2023. Available: https://ieeexplore.ieee.org/book/10251329