# Microservices architecture for real-time adaptive learning at scale

Modern recommendation systems at companies like Netflix, Uber, and Pinterest achieve **sub-100ms latency while serving millions of predictions per second** through a convergent architectural pattern: event-driven microservices paired with dual-store feature systems, ANN-powered retrieval, and sophisticated auto-scaling. Netflix alone generates **75-80% of viewing hours through algorithmic recommendations**, processing 1.25 million requests per second at peak. This report provides implementation-ready guidance on building these systems, covering the critical intersection of event-driven architectures, feature stores, and latency optimization.

## Event-driven patterns separate training from inference

The fundamental insight driving ML microservices architecture is that training (write-heavy, batch-oriented) and inference (read-heavy, real-time) have fundamentally incompatible requirements. **CQRS (Command Query Responsibility Segregation)** addresses this by maintaining separate models for each workload.
( Martin Fowler )

In production ML systems, the command side handles training triggers, user feedback capture, and feature updates, while the query side optimizes exclusively for low-latency prediction serving. ( Martin Fowler ) Vispera, an image recognition company processing millions of images monthly, implemented this pattern with EventStoreDB: a single user task generates approximately **18 events on average**, with Protocol Buffers encoding reducing event size to just 25% of JSON equivalents. ( kurrent ) ( Kurrent ) Their architecture projects events to multiple read-optimized stores—MongoDB for task management, PostgreSQL for operator monitoring, and a dedicated training data repository for model improvement.

**Event sourcing** complements CQRS by storing every state change as an immutable event, enabling complete audit trails and replay capability essential for ML debugging. The typical event flow for recommendation systems includes:

- **UserInteractionEvent**: Captures clicks, views, and purchases with timestamps

- **PredictionEvent**: Logs model version, prediction, and latency per request

- **FeedbackEvent**: Correlates actual outcomes with prior predictions

- **ModelAccuracyEvent**: Aggregates true/false positive rates for monitoring

For ML pipeline coordination, the choice between **choreography and orchestration** depends on complexity. Choreographed pipelines, where services react independently to events without central coordination, work well for real-time inference with fewer than five services— ( DEV Community )Netflix uses this pattern for streaming recommendations. However, complex training workflows benefit from orchestration through tools like Apache Airflow or Temporal, which provide visibility into pipeline state and built-in transaction management. ( Microsoft Learn ) The recommended hybrid approach uses orchestration for training (visibility, error handling) and choreography for inference (low latency, independent scaling).

Regarding streaming platforms, **Apache Kafka delivers 2x write throughput compared to Pulsar** with p99 latencies around 5ms at high load. However, Pulsar provides more consistent 5-15ms p99 latencies and **60% faster historical data reads** for catch-up scenarios. Netflix processes **700 billion Kafka messages daily** across 36 clusters, while Spotify uses Pulsar for real-time user behavior streaming. Choose Kafka for maximum throughput and ecosystem maturity; choose Pulsar for multi-tenant environments or when consistent latency matters more than peak performance.

## Feature stores bridge batch training and real-time serving

Every major ML platform converges on a **dual-store architecture**: an offline store optimized for batch training (S3, BigQuery, Snowflake) and an online store optimized for inference (Redis, DynamoDB, Cassandra). The offline store maintains months of historical data for point-in-time correct training, while the online store retains only the latest feature values for sub-millisecond retrieval. (Feast) (Aerospike)

Benchmark data reveals significant performance differences between implementations. For **50 features at 10 requests per second**, Hopsworks achieves **6-8ms p99 latency** compared to Feast with Redis at **20-25ms p99** and Feast with Cassandra at **35-40ms p99**. This gap widens dramatically at scale: retrieving 70 features for 30 entities under load sees Hopsworks maintain **30-40ms p99** while Feast/Cassandra degrades to **700-800ms p99**.

| Feature Store | Online Store Backend | P99 Latency (50 features) | Key Differentiator |
|---|---|---|---|
| **Hopsworks** | RonDB (built-in) | 6-8ms | Integrated KV store, fastest benchmarks |
| **Feast** | Redis/DynamoDB (pluggable) | 20-25ms | Modularity, no vendor lock-in |
| **Tecton** | DynamoDB (managed) | 15-20ms* | End-to-end automation, monitoring |
| **AWS SageMaker FS** | DynamoDB/In-memory | ~10ms | AWS integration, Iceberg support |

*With 90% cache hit rate

**Uber's Michelangelo/Palette** pioneered the lambda architecture for features, combining Spark batch pipelines with Kafka/Flink streaming. (Featureform) (Pantelis) Their feature hierarchy (Domain → Feature Group → Feature Name → Join Key) enables **FeatureServingGroups** that guarantee specific latency SLOs by routing to appropriately provisioned Cassandra/Redis clusters. Airbnb's Zipline follows a similar pattern with a declarative configuration language that reduced feature engineering from months to days. (Slideshare)

For streaming features, DoorDash's **Riviera framework** built on Flink demonstrates the production pattern: YAML configuration defines the source-transformation-sink pipeline, with Flink SQL serving as the DSL for

feature transformations. (DoorDash) Their optimizations include using **Redis Hashes** for entity-grouped features (2.85x CPU reduction), string hashing for feature names (reducing 27-byte strings to 32-bit integers), and client-side caching delivering **70% performance improvement**.

## Achieving sub-100ms recommendations requires optimization at every layer

The **100ms latency threshold** represents a critical UX boundary—Amazon found every 100ms of slowness causes 1% conversion loss. (Emanuele) Achieving this at scale requires coordinated optimization across ANN search, model serving, caching, and request handling.

**Approximate Nearest Neighbor search** benchmarks show significant variation. Qdrant delivers the highest requests per second with best filtering capabilities. ScaNN provides **2x accuracy improvement** over competitors for inner-product similarity through anisotropic vector quantization. FAISS excels at billion-scale with GPU acceleration. For vector databases, the critical insight is that performance metrics only matter with recall attached—comparing "10ms at 90% recall" versus "50ms at 99% recall" is meaningless.

| Scale | Recommended Approach |
| --- | --- |
| <1M vectors | HNSW (hnswlib) — simple, fast |
| 1-10M vectors | FAISS IVF + HNSW hybrid |
| 10M-1B vectors | Milvus/Qdrant with disk-based indexes |
| >1B vectors | FAISS with GPU + sharding |

**Model serving infrastructure** choices significantly impact latency. NVIDIA Triton delivers **p50 latency of 5.86ms and p99 of 8.09ms** for ResNet50 on H100, with the Bing team reporting **2x GPU utilization** improvements while holding latency flat. TensorFlow Serving achieves similar numbers (~6ms p50) for TensorFlow-native deployments. For LLMs, **continuous batching delivers 23x throughput improvement** over static batching— (Anyscale) vLLM's PagedAttention implementation is the current standard.

**Caching strategies** prove critical at scale. Uber's CacheFront achieves **99.9% cache hit rate handling 40+ million requests per second** through adaptive timeouts, negative caching, and cache warming. (Clickit Tech) For embeddings, exact-match caching (input hash → embedding) works well for static content, while semantic caching using similarity thresholds can return cached results for semantically similar queries at the cost of careful threshold tuning.

**Latency budgeting** allocates the total budget across microservice components. Whatnot reduced their ML pipeline from **700ms to 120ms p99** (5.8x improvement) through parallel execution of independent operations and aggressive optimization:

| Component | Typical Budget | Optimization Focus |
|---|---|---|
| Network/Transport | 5-15ms | gRPC, connection pooling |
| Feature Retrieval | 10-20ms | Redis caching, colocation |
| ANN Search | 5-15ms | Index optimization, quantization |
| Model Inference | 20-50ms | Batching, model pruning |
| Post-processing | 5-10ms | Async operations |
| Buffer | 10-20% | Absorb tail latency |

Migrating from REST to gRPC alone provides **40-60% latency reduction** due to binary protobuf encoding (3-7x smaller, 5-10x faster parsing). (Nexastack) Connection pooling eliminates per-request overhead—Google recommends targeting fewer than 50 concurrent RPCs per gRPC channel. (Google Cloud)

## Production systems reveal convergent architectural patterns

Major tech companies have independently converged on similar architectures, validating these patterns at billion-user scale.

**Netflix** serves 260M+ subscribers with a three-tiered architecture (offline, nearline, online). Their Keystone pipeline processes **trillions of events daily** through 36 Kafka clusters. The Hydra multi-task learning system handles homepage ranking, search, and notifications in unified models, while their Foundation Model learns member preferences centrally and distributes embeddings to downstream models. EVCache delivers **95-99% cache hit rates** for personalization data, storing 100-200KB per profile. At peak, Netflix handles **~750 million API calls every 10 minutes**.

**Uber's Michelangelo** demonstrates end-to-end ML platform design, making **millions of predictions per second** across hundreds of deployed use cases. Their Transformer/Estimator pattern ensures identical model behavior online and offline, eliminating training-serving skew. The unified DSL enables feature transformations to be defined once and executed identically in Spark (batch) and Samza (streaming). Canvas (Michelangelo 2.0) extends this with a Python DSL for streamlined model development.

**Pinterest's Pixie** achieves remarkable efficiency: a **bipartite graph of 7 billion nodes and 100+ billion edges** runs on single machines with terabyte-scale RAM, delivering **60ms p99 latency at 1,200 requests per second per server**. Their graph pruning strategy removes low-quality boards, reducing graph size by 6x while improving recommendation quality by 58%. The Pixie random walk algorithm with early stopping provides the same quality in ~50% fewer steps compared to naive approaches.

**DoorDash's Sibyl** prediction service handles **millions of predictions per second with <100ms latency**, performing approximately 1,000 feature lookups per request from their Redis feature store containing billions of

records. Their computational graph approach compiles Python ensemble definitions to C++ serving code, achieving **12x CPU time reduction**. The two-layer ETA architecture separates base probabilistic predictions from business-objective optimization.

**Spotify** powers 713 million listeners through a three-model approach: collaborative filtering for behavioral patterns, NLP for text analysis, and audio analysis for sonic features. Their Cassandra infrastructure spans **100+ clusters with 3,000+ nodes**, achieving **<5ms average latency at the 95th percentile** for feature reads handling 50K operations per second.

## Scaling patterns balance throughput, latency, and resilience

Horizontal scaling for ML inference differs fundamentally for stateless traditional models versus stateful LLM services with KV caches. For **stateless inference**, standard Kubernetes HPA works well when configured with queue size metrics—Google recommends starting with a target queue depth of 3-5 and increasing until latency targets are met. (google)

For **stateful LLM services**, session-aware routing is essential. (Kubernetes) KubeAI's Consistent Hashing with Bounded Loads achieved **87.4% cache hit rates with sub-400ms latencies** by routing requests to pods most likely to have relevant cached content. (Kubeai) Their PrefixHash strategy delivers **4.3x throughput improvement and 71% lower time-to-first-token** compared to round-robin. The recommended pattern externalizes conversation state to Redis while sending full token history each turn, letting vLLM's prefix caching avoid redundant computation. (DevTechTools)

**Latency-based auto-scaling** should target p95 as the primary SLO with p99 monitoring for tail anomalies. (OneUptime) Predictive scaling approaches like LA-IMR achieve **20.7% p99 latency reduction** over reactive baselines by proactively scaling before queues build up. (arXiv) For GPU workloads, KEDA enables scaling based on Redis queue depth or Prometheus GPU utilization metrics, with scale-to-zero capability for cost optimization. (DEV Community)

**Multi-region deployment** requires careful handling of model consistency. (Tecton) The AWS SageMaker pattern uses S3 bucket replication to secondary regions, with CodePipeline extended for cross-region deployment and Route 53 for traffic routing. (AWS) For feature stores, Tecton uses DynamoDB Global Tables for automatic replication, (Tecton) while Hopsworks recommends Active/Active for online stores (lowest latency) and Active/Passive for offline stores (simpler operations).

**Graceful degradation** implements three-level fallback: primary model service → lightweight backup model → rule-based defaults. Circuit breakers with configurable thresholds (typically 3 failures in 60 seconds before opening) prevent cascading failures. (Markaicode) Quality-versus-latency tradeoffs under load include falling back to smaller models, reducing context windows, or returning cached responses. (Mhtechin) The key principle: partial degraded responses beat complete failures.

## Synthesis: building production adaptive learning systems

These patterns interact in production through careful orchestration. A typical request flow traverses:

1. **Gateway** receives request, applies rate limiting (5-10ms)

2. **Feature service** retrieves entity features from online store (10-20ms)

3. **Candidate retrieval** performs ANN search on embedding index (5-15ms)

4. **Ranking service** scores candidates through model inference (20-50ms)

5. **Post-processing** applies business rules and diversity constraints (5-10ms)

Events flow bidirectionally: user interactions stream to Kafka topics, feeding both real-time feature updates (Flink → online store) and batch training pipelines (Spark → offline store → model training → deployment). Model updates propagate through versioned registries, with traffic gradually shifted via weighted routing.

The key architectural decisions are:

- **Event sourcing over mutable state** for complete audit trails and replay capability

- **Dual-store feature architecture** with Redis/RonDB online and S3/BigQuery offline

- **Choreographed inference with orchestrated training** for the right tradeoffs in each domain

- **Latency budgets defined per component** with p95/p99 monitoring and predictive scaling

- **Circuit breakers with multi-level fallback** to maintain availability under failures

Companies achieving sub-100ms latency at scale combine all these patterns: Netflix's three-tiered computation, Uber's unified online/offline features, Pinterest's efficient graph algorithms, DoorDash's compiled ensemble serving, and Spotify's algotorial hybrid of algorithms and editorial curation. The consistent theme is aggressive optimization at every layer—from Protocol Buffer encoding to GPU-accelerated inference to cache-aware load balancing—while maintaining resilience through event-driven decoupling and graceful degradation.