

# Scaling Observability for Real-Time Personalization Engines

Sai Kumar Bitra

Principal Software Engineer - AT&T

## Why Observability Matters for Personalization



Personalization drives engagement and revenue: Companies that excel at personalization see up to *40% more revenue* from those efforts, as users respond to highly relevant content.



Real-time experience, zero room for error: Users expect personalized results in milliseconds - delays or failures directly hurt satisfaction and conversions.

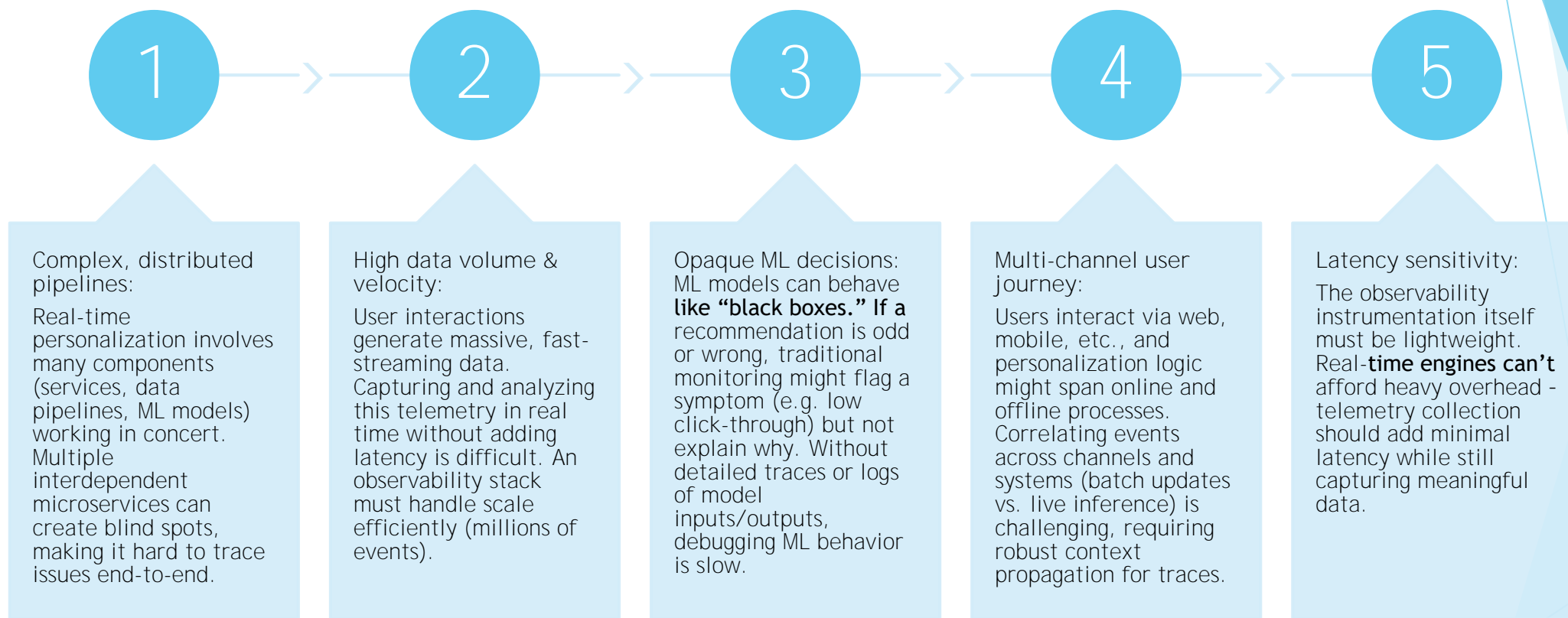


Complex ML decisions need validation: **Without observability, it's hard to verify if the** recommendation engine is doing the right thing at the right time. Lack of insight = lost user trust and missed opportunities.



Competitive edge: Telemetry data (metrics, logs, traces) helps continuously improve algorithms and user experience, keeping your personalization engine (and product) a step ahead of the competition.

# Challenges in Observability for Real-Time ML Systems



# System Architecture Overview

- ▶ Ingestion Pipeline: Captures user events in real time (e.g. clickstream via Kafka or Kinesis) and feeds them into a streaming system for model updates. Ensures the model sees latest behavior.
- ▶ Decision Engine: The core personalization service (microservice or function) that computes recommendations per request. Integrates the ML model inference and business rules into the API flow.
- ▶ Delivery/API Layer: Exposes personalized content through REST/GraphQL APIs (or SDK) to the application. This is the entry point handling high QPS with low latency.
- ▶ Data & Feature Stores: Fast data sources (NoSQL, in-memory caches, feature stores) providing user profiles, content info, etc., for real-time decisioning.
- ▶ Telemetry Sources: Each component is instrumented - API, engine, model, and data pipeline emit logs, metrics, and traces. This feeds the Observability Pipeline (using collectors/agents) that aggregates data for analysis.

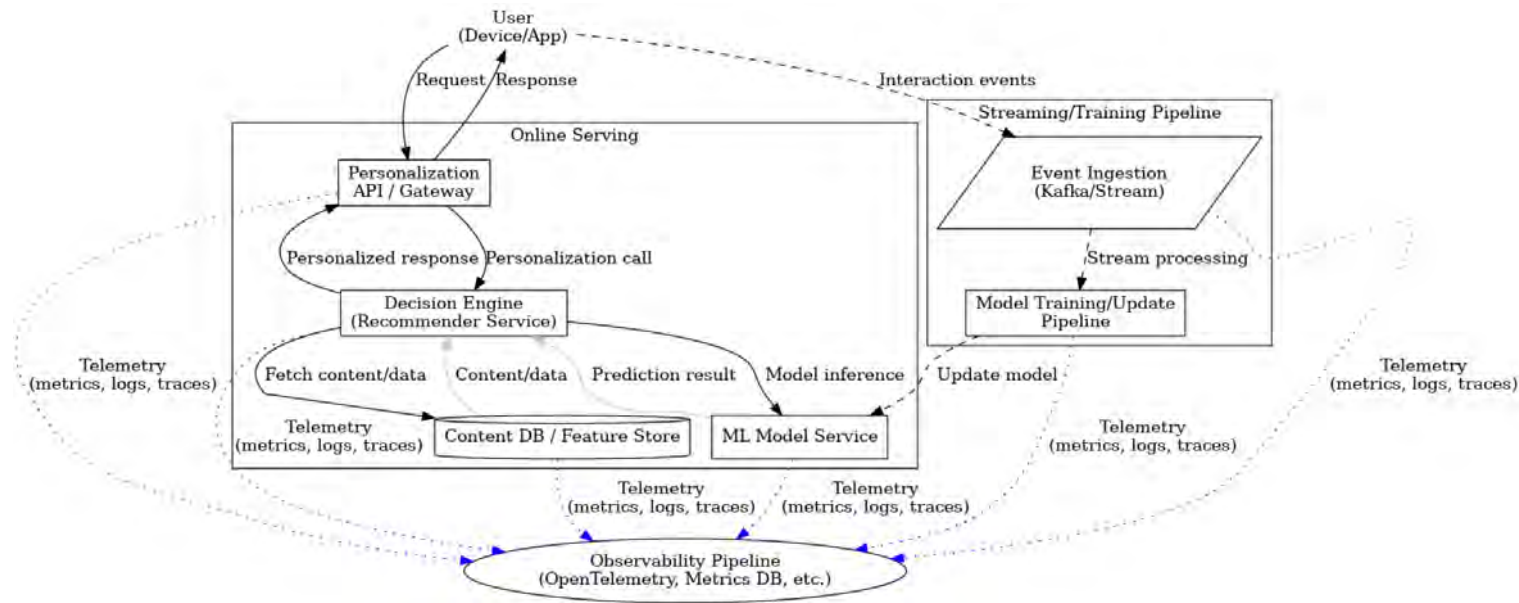


Figure: Representative architecture for a real-time personalization engine. The user's request flows through an API Gateway to a Decision Engine (recommendation service) that calls an ML Model Service and fetches data from a feature store. Meanwhile, user interaction events stream through an ingestion pipeline (e.g. Kafka) into a model update process. Telemetry (metrics, logs, traces) is emitted from each component into a centralized observability pipeline.

# Scaling Observability in a Real-Time Architecture

- ▶ End-to-end tracing: Implement distributed tracing from the user request through the personalization service and back. Every request carries a trace ID through the API, decision engine, model calls, and database queries, enabling end-to-end visibility of each personalized response.
- ▶ Unified instrumentation: Adopt OpenTelemetry (OTel) SDK across services for standard metrics, logs, and trace collection. OTel provides vendor-agnostic APIs to instrument code once and export telemetry to multiple backends.
- ▶ High-throughput metrics: Use a time-series DB (e.g. Prometheus) to scrape and store metrics at scale. Metrics (QPS, latency, recommendation counts, etc.) are aggregated to avoid overwhelming storage. For example, compute percentile latencies and sample fine-grained events.
- ▶ Log aggregation: Stream logs (JSON structured logs including user/session IDs, model decisions) into a central system (ELK stack or cloud logging). Structured logs indexed by request or trace ID allow filtering for specific user sessions or error scenarios.
- ▶ Scalable storage & retention: Employ a back-end that can handle high data volume - e.g. a scalable observability platform or data lake for long-term analysis. Ensure retention policies and downsampling are in place so historical trends can be analyzed without infinite storage growth.

# Key Metrics That Matter

- ▶ Latency and throughput: Measure how fast and how many recommendations are served. Track p95/p99 latency of the personalization API and model inference time. Even small delays can frustrate users in real time, so tight SLOs (e.g. <100ms) are common.
- ▶ Error rates and timeouts: Monitor error rates in the personalization pipeline - e.g. fallback rates (when the model fails and a default is used), API 5xx errors, timeouts calling the model or data store. Spikes indicate issues in the ML service or data layer.
- ▶ Engagement metrics: Tie observability to business KPIs. Track user engagement outcomes of personalization: click-through rates, conversion rates, dwell time, etc., in real time. A drop in engagement might signal a personalization issue (e.g. model drift or bad recommendations) beyond just system health.
- ▶ Model performance indicators: Custom metrics from the ML model itself - e.g. the confidence score distribution of recommendations, frequency of each model variant being used, drift metrics comparing **live input data to training data. These help detect when the model's quality is degrading.**
- ▶ Resource utilization: Keep an eye on system metrics (CPU, memory of the engine, throughput of Kafka topics, queue lengths) to ensure the infrastructure can handle load. Sudden changes could explain performance outliers.

# Tracing the Personalization Journey

- ▶ Distributed trace example: When a user visits and gets personalized content, a trace follows the request from the frontend API through the decision engine, into the model service and database, and back. Each segment (span) of the trace records the duration and result of that step.
- ▶ OpenTelemetry for tracing: Using OTel, insert trace instrumentation at critical points - e.g. mark a span around the “*GenerateRecommendations*” function, another for “*ModelInference*”, another for “***DB Fetch***”. These spans then appear in tracing UIs (Jaeger, Zipkin) to visualize the call flow.
- ▶ Bottleneck identification: Traces help pinpoint where latency comes from. For example, a trace might show the API spent 5ms, the decision engine 20ms, but a DB call took 80ms due to a slow query. Such insights are hard to get from metrics alone and enable targeted optimization.
- ▶ Correlating with user actions: Tag traces with user or session IDs (carefully, to avoid high cardinality) or experiment IDs. This allows correlating a poor experience (e.g. user saw a stale recommendation) with the exact trace of the decision that led to it, providing valuable context for debugging.
- ▶ Sampling strategy: In high QPS systems, trace every request may be infeasible. Implement trace sampling (e.g. sample 1% of requests, or always sample errors). This keeps overhead low while ensuring you collect traces for anomaly scenarios and a representative slice of traffic.

# Logging for ML Decision-Making

- ▶ Rich, structured logs: Logging in a personalization engine should capture key inputs and outputs of the ML decisions. For each request, log details like user ID, relevant features (age bracket, segment, etc.), the recommendations returned, and maybe an explanation score or reason code if available.
- ▶ Enable root-cause analysis: When something goes wrong (e.g., irrelevant content shown), logs can **reveal what the model saw. For instance, a log might show “Input user\_segment=null” which explains a** bad default recommendation. These clues in logs are critical for debugging data or logic issues.
- ▶ Trace-log correlation: Integrate logs with tracing context. Include the trace ID or request ID in every log entry. This way, when inspecting a distributed trace, you can jump to the associated logs across services for that same request. It provides both high-level flow and detailed event data together.
- ▶ Anonymize and protect data: Personalization logs can include user data - ensure compliance by anonymizing or redacting sensitive info (GDPR concerns). An observability platform should enforce data handling policies while still giving engineers enough info to troubleshoot.
- ▶ Use log levels wisely: For real-time systems, use INFO level to log normal decision info (with sampling if needed), WARN for unusual situations (e.g. model fallback used), and ERROR for actual failures. This **allows alerting on specific log patterns (like “fallback model used too often” or exceptions).**



# Scaling Observability: Patterns & Best Practices

- ▶ Avoid metric overload: Be mindful of metric cardinality. Rather than a metric labeled with every user or item ID (which would explode), use aggregated metrics (e.g. counts by category or percentile latency). High-cardinality metrics can overwhelm Prometheus and memory.
- ▶ Trace and log sampling: As mentioned, sample traces to balance detail vs. cost. Similarly, consider logging at high detail only for a sample of requests or when errors occur. This adaptive telemetry reduces noise and cost while retaining diagnostic power.
- ▶ Batching and buffering: Use collectors (like OpenTelemetry Collector or Kafka) to buffer telemetry **data. For example, batch trace exports and send in chunks, so the application isn't blocked on sending** each span. Buffer logs to avoid disk I/O becoming a bottleneck.
- ▶ Scalable storage & querying: Implement a tiered storage for observability data. Recent data stays in fast, queryable stores (for quick troubleshooting), while older data is down-sampled or moved to cheaper storage. This ensures the observability platform scales without monstrous costs.
- ▶ Resilience of observability systems: Treat your monitoring pipeline as a critical part of the system. Scale out your metrics and logging backends, and set up alerts on the observability systems themselves (e.g. if the log forwarder lags or Prometheus is behind). An observability outage can be especially painful during a production incident.

# Incident: Observability in Action (Debugging Story)

- ▶ The scenario - silent model drift: Imagine engagement gradually dropped over a week for a personalized feed. No single component crashed, but the recommendations became less relevant. Without observability, this went unnoticed until business metrics fell significantly.
- ▶ Detection through metrics and alerts: With a strong observability setup, an alert fired when CTR (click-through rate) dipped 10% below baseline. At the same time, a custom model drift metric showed the distribution of inputs had shifted from training data. This correlation hinted the model was no longer tuned to current user behavior.
- ▶ Using traces to pinpoint impact: Engineers pulled up distributed traces for user sessions with low engagement. The traces showed normal latency, but a pattern emerged: many users received a default recommendation path (logged in the decision engine spans). This indicated the model was unsure and falling back.
- ▶ Root cause via logs: By querying logs (filtered by trace IDs from those sessions), the team discovered the **feature “user\_age\_bucket” was frequently null for affected users. Upstream, an ETL bug had stopped populating that feature for new users - crucial information the model needed. The model’s outputs drifted as a result.**
- ▶ Resolution and takeaways: Thanks to observability, the team identified the root cause in a few hours (versus days). They fixed the ETL and retrained the model. Post-mortem showed that prior to their observability tools, a similar issue took 3× longer to diagnose. The incident reinforced how metrics tied to business KPIs and detailed telemetry make a decisive difference in debugging.

# Tools & Stack in Action

- ▶ OpenTelemetry: The open-source standard for instrumentation - used to collect metrics, traces, and logs in a unified way. OTel SDKs were integrated into the personalization microservices, allowing easy export of telemetry. This avoids vendor lock-in and ensures compatibility with many backends.
- ▶ Prometheus & Grafana: Prometheus serves as the metrics backend, scraping instrumented services for stats. Grafana provides dashboards and visualizations, from request rates to model latency histograms, and even combined views (Grafana can display both Prom metrics and trace data together). Alerts are configured in Prometheus (e.g. high error rate) and Grafana (for business metrics thresholds).
- ▶ Distributed Tracing Systems: Jaeger or Zipkin is deployed for tracing, receiving spans from **OpenTelemetry. Engineers can open Jaeger's UI to see a waterfall of a user request trace across all services**, or use Grafana Tempo for a seamless metrics-to-trace drill-down.
- ▶ Logging Pipeline: An ELK stack (Elasticsearch, Logstash, Kibana) or cloud log service aggregates the structured logs. Kibana dashboards enable searching for specific user IDs or error codes. In this case, logs are also integrated with trace IDs - e.g. using Kibana to filter all logs for a given trace.
- ▶ APM and Advanced Tools: In addition to open-source tools, the team leveraged an APM solution (e.g. Dynatrace or DataDog) for features like anomaly detection on metrics and AI-assisted root cause analysis. These tools can complement the DIY stack by catching subtle issues (like a sudden change in user behavior patterns) automatically.

# Tools Enabling Observability at Scale

**Purpose:** Highlight the actual tools and technologies used across telemetry layers (logs, metrics, traces) in a scalable, ML-powered personalization system.

## Tools Used at Scale

### Instrumentation & Telemetry Collection

#### •OpenTelemetry (OTel)

Unified standard for metrics, logs, and traces. Widely adopted and vendor-agnostic.

- SDKs integrated into microservices
- Supports trace context propagation and auto-instrumentation

### Metrics Collection & Dashboards

#### •Prometheus

Time-series DB scraping application metrics (latency, throughput, error rates)

#### •Grafana

Dashboards visualizing system health, model KPIs, business metrics (e.g., CTR)

### Distributed Tracing

#### •Jaeger or Grafana Tempo

Visualizes end-to-end traces of user requests across services

- Identify slow spans
- Trace model decisions and database latency

### Log Aggregation & Search

#### •ELK Stack (Elasticsearch, Logstash, Kibana) or Loki

Structured logs (JSON) searchable by trace/session ID

- Root cause analysis of user journeys
- Cross-service error correlation

### Advanced Monitoring & APM

#### •Dynatrace, DataDog, or New Relic

Used for:

- Anomaly detection on engagement or latency
- Real-time dashboards across ML, API, and infra layers

### Trace/Log Storage at Scale

#### •Kafka or OpenTelemetry Collector

Used as a buffering layer for telemetry before storage

#### •ClickHouse / Parquet on Data Lake

Used for querying long-term observability data and joining with offline features

# Platform Impact of Scaled Observability



Dramatic reduction in outages: With robust observability, issues are caught and resolved faster. One case study saw mean time to resolution improve by 75% after adopting modern observability practices. Faster fixes mean less user impact and higher uptime.



Improved performance and reliability: The team can identify bottlenecks proactively. For example, by monitoring latency per component, they optimized a slow database query before it became a customer issue. Overall system reliability (measured by error rates and uptime) increased noticeably.



Higher confidence in deployments: Observability data gave developers and stakeholders confidence to **push updates faster. They know that if a new model or feature causes regressions, they'll see it in** dashboards or alerts within minutes. This improved the velocity of experimentation in personalization.



Better alignment with business goals: Because telemetry was tied into engagement metrics, the platform team now speaks the same language as product owners. They can demonstrate how a latency improvement raised conversion, or how a quick rollback (informed by an alert) saved revenue. Observability became a business enabler rather than just an ops tool.



Cost vs. benefit managed: Initially, adding so much instrumentation raised concerns about data volume and **cost. However, by refining what's collected (sampling, aggregation), they kept overhead reasonable. The** insights gained far outweighed the costs - in fact, better efficiency and fewer incidents saved money in the long run.

# Lessons Learned & Best Practices

- ▶ Build in observability from Day 1: **Don't bolt on monitoring later** - design the personalization system with observability in mind. Define key metrics and events upfront (for each new feature or model), and **instrument as you develop. It's harder to retrofit once issues occur.**
- ▶ Treat logs, metrics, traces as first-class citizens: All three telemetry types provide unique insights. Use all of them in tandem - metrics for high-level health, traces for flow analysis, logs for details. Invest in tools that unify these views to avoid data silos.
- ▶ Design for visibility, not just performance: Sometimes you might choose a slightly less complex design **because it's easier to observe. For example, using HTTP calls between services with trace context** instead of an opaque binary protocol can make troubleshooting easier. *Observable architecture* is a feature, not an afterthought.
- ▶ Empower the team with data: Train engineers and data scientists to use the observability dashboards and tracing tools. Encourage a culture where hypothesis-driven debugging (using telemetry data) is the norm. This makes incident response more collaborative and effective.
- ▶ Iterate and improve: **Observability is not “set and forget.”** Continuously refine what you monitor. As the personalization engine evolves (new models, features), update dashboards and alerts. Periodically do game days or incident drills to ensure your telemetry actually helps resolve issues quickly.

# Final Thoughts and Call to Action

- ▶ Observability is product-critical: In a real-time ML environment, monitoring isn't just an ops concern - it directly influences user experience and trust in your product. A personalization engine is only as effective as your ability to observe and adjust it in real time.
- ▶ Trust through transparency: **By illuminating the “black box” of your personalization system with metrics and traces**, you build confidence in the recommendations both internally (for developers, data scientists) and externally (for users, through reliability). As one industry example showed, modern observability practices not only improved MTTR but also cut costs and boosted user satisfaction.
- ▶ Take action: Evaluate your current observability maturity. Are you collecting the right data at the right points? Consider integrating open standards like OpenTelemetry and setting up a robust pipeline with tools like Prometheus and Grafana. Start with a small slice of your system, instrument it, and expand from there.
- ▶ Continuous improvement: Scaling observability is an ongoing journey. Use telemetry insights to drive improvements in both the ML models and the system architecture. Every outage prevented or performance win discovered via observability is a direct win for your users and your business.
- ▶ Call to action: Embrace the mindset that **if you can't measure it, you can't improve it**. Bring clarity to the complexity of real-time personalization by observing it smarter. Your personalization engine - and your users - will thank you for it.

# Thank You

Sai Kumar Bitra

<https://www.linkedin.com/in/sai-kumar-bitra-a4054121/>