



6 streaming KPIs that make or break AI pipelines

A quick guide to the numbers that
keep real-time models fast,
accurate, and cost efficient.

Why these metrics matter



Batch dashboards only flag trouble after the fact. Track these streaming KPIs instead, and you'll know— instantly—if your feature pipes are healthy enough for production AI.



Building agents or automated decision-making? Think in events, not bulk streams. Event-driven architecture is how dynamic data reaches models safely and on time.



The KPI table you actually need

Metric	Why it matters for AI/ML	How to track Redpanda	Redpanda vs. Apache Kafka®
<p>Throughput (MB/s, msgs/s)</p>	Keeps feature stores fed during traffic spikes	rpk cluster quotas describe ¹	Thread-per-core design sustains up to 3× Kafka throughput on equal hardware ²
<p>P99 latency (ms)</p>	Drives inference speed; slow pipes mean stale predictions	Built-in latency dashboard in Redpanda Console	Up to 10× lower P99 vs. JVM-based Kafka ²
<p>Message durability (ack, ISR)</p>	Lost events create blind spots and retraining cost	Quorum replication + Tiered Storage = zero data loss	Redpanda's Raft avoids slow followers or lost state for superior data safety
<p>End-to-end lag (producer→consumer)</p>	Measures real-world delay from event to action	Rpk group describe lag column ³	Auto-balancing keeps lag flat as partitions grow
<p>Partition balance (% skew)</p>	Hot partitions throttle model inputs	Continuous Data Balancing enabled by default	Built-in partition balancing removes the need for Cruise Control, cutting a third-party dependency and freeing broker resources ⁴
<p>Error rate (fails/sec)</p>	Silent drops distort training data	Dead-letter queues + filter-by-error in Console	Idempotent producers detect and reject duplicate messages, preventing wasted retries and preserving order ⁵
<p>Cost per TB per year</p>	Determines ROI and footprint	Redpanda Cloud TCO tools	Up to 6x lower infra cost on average ⁶

Hit the targets with Redpanda



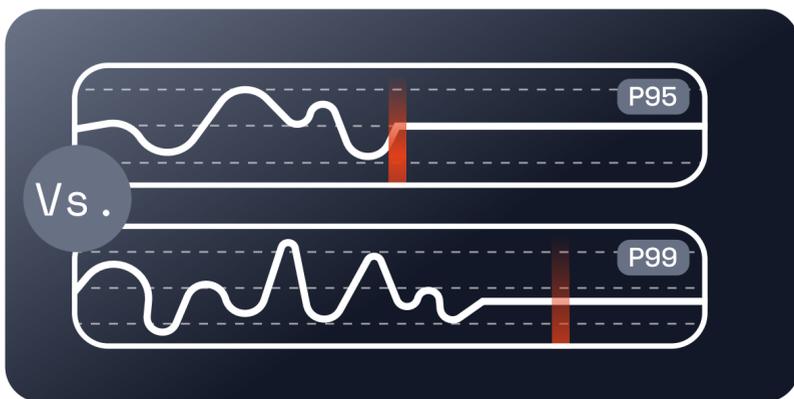
Keep latency low, lag flat, and data safe with one binary and built-in tools.

Make these 3 adjustments



1

Set an alert when lag exceeds **1.5 × steady state**.



2

Use the **ListKafkaConnections** API endpoint to identify and troubleshoot problematic clients.



3

Use **Tiered Storage** for long-tail training features.



Run it now in Redpanda Serverless



SERVERLESS QUICKSTART



TRY THE ENTERPRISE EDITION



Sources:

¹ Redpanda. [Manage throughput documentation](#). Accessed August 7, 2025.

² Redpanda. [Redpanda vs. Kafka with KRaft: Performance update](#). May 11 2023. Accessed August 7, 2025.

³ Redpanda. [Rpk_group describe documentation](#). Accessed August 7, 2025.

⁴ Redpanda. [High Availability documentation](#). Accessed August 7, 2025.

⁵ Redpanda. [Idempotent producers documentation](#). Accessed August 7, 2025.

⁶ Redpanda. [Redpanda vs. Apache Kafka \(TCO Analysis\)](#). Published October 18, 2022. Accessed August 7, 2025.