

# Real-Time Fraud Scoring at Wire Speed

In-Process AI Inference on Payment Event Streams — No Model Server. No GPU. No Spark.

<b>409.4 EPS</b> Avg throughput 126,224 records · 5.14 min	<b>511 EPS</b> Peak throughput sustained · no warmup gap	<b>6.8 ms</b> Inference avg ONNX in-process · CPU only	<b>8.12×</b> vs baseline same JAR · same config
--	--	---	---

<b>0</b> Dropped records across all 5 benchmark runs	<b>0.049%</b> GC overhead 33 ms max pause · G1GC	<b>22</b> Provenance headers per output record · SHA-256	<b>55 MB</b> Live heap (post-GC) stable · load-independent
--	--	--	--

## TUNING PROGRESSION · 5 RUNS · SAME HARDWARE

Run	Avg EPS	Peak EPS	Infer avg	Key change
1 · Baseline	50.4	76.9	18.7 ms	pool=1 · no batching · batch=32
2 · Thread tuning	52.3	102.6	77.3 ms	pool=4 · intra=3 · thread collision
3 · Batching on	145.6	294.4	16.4 ms	batching enabled · pool=4 · intra=2
4 · Batch aligned	280.0	419.9	11.1 ms	pipeline batch=16 → matches embed max
<b>5 · This run ✓</b>	<b>409.4</b>	<b>511.1</b>	<b>6.8 ms</b>	<b>JVM fully JIT-warm · identical config</b>

### PIPELINE SPEC

**Transform chain:** STRING\_TO\_WIREEVENT →  
DJL\_EMBEDDING → FRAUD\_SCORE → KAFKA  
**Model:** all-MiniLM-L6-v2 (ONNX) · champion / production

### AUTHOR

**Steven Lopez**  
Solutions Architect & Creator  
StreamKernel LLC

**Hardware:** Intel i9-8950HK · CPU only · no GPU  
**Config:** parallelism=4 · pool=4 · intra-op=2 · batch=16  
**Kafka:** 12 partitions · lz4 · 256 KB batch  
**GC:** G1GC · 4 GB heap · 3 threads · MaxPause=50ms

[lopezstevie@gmail.com](mailto:lopezstevie@gmail.com)  
[linkedin.com/in/steve-lopez-b9941](https://www.linkedin.com/in/steve-lopez-b9941)  
[github.com/IntuitiveDesigns/StreamKernel-io](https://github.com/IntuitiveDesigns/StreamKernel-io)  
May 20, 2026 · v0.2.0

## 0 1   E X E C U T I V E   S U M M A R Y

---

Financial institutions process tens of millions of payment events daily. Every millisecond of scoring latency increases fraud exposure. Every external network hop to an inference service is a failure surface, a compliance boundary, and a cost center. Every scored transaction that loses its audit trail creates regulatory risk.

StreamKernel eliminates all three problems simultaneously. It is a Java 21 event processing kernel that executes AI inference in-process — inside the same JVM as the pipeline, with no external model-serving infrastructure required. Fraud scores, embeddings, decisions, and full cryptographic provenance travel together in every output event, from the moment a transaction enters the pipeline to the moment it lands in Kafka, MongoDB, Delta Lake, or any downstream system.

In a live 5-minute benchmark run on May 20, 2026:

- **126,224 synthetic payment transactions scored end-to-end**
- **Zero record loss · Zero errors · Zero DLQ events**
- **Every output record carried 22 cryptographic provenance headers**
- **GC overhead: 0.049% · Heap at close: 1,214 MB of 4 GB (28.3%)**
- **Hardware: Intel i9-8950HK laptop — CPU only, no GPU**

This white paper presents the full technical profile of StreamKernel's financial fraud scoring pipeline: architecture, transform chain, benchmark methodology, raw metrics, and fraud decisioning output. The results establish a credible performance and integrity floor for production deployment in regulated financial services environments.

# The Latency and Auditability Gap in Real-Time Fraud Decisioning

## 0 2   P R O B L E M   S T A T E M E N T

---

Modern fraud detection has converged on machine learning as the primary scoring mechanism. Gradient-boosted trees, neural networks, and embedding-based semantic similarity models now underpin the risk engines of the world's largest card networks, banks, and payment processors. The models themselves are no longer the bottleneck. The infrastructure that serves them is.

### The External Inference Service Problem

---

The dominant deployment pattern for ML-powered fraud scoring today is the external inference service: a model server (TensorFlow Serving, Triton, SageMaker, or a cloud endpoint) that receives a serialized feature vector over HTTP or gRPC, runs inference, and returns a score. This pattern has three compounding problems:

- ▶ Network latency adds 5–50 ms per transaction, every transaction. At 10,000 TPS, a 10 ms average round-trip to an inference service adds 100 seconds of cumulative latency per second of transactions.
- ▶ The inference service is a new failure domain. It must be deployed, scaled, monitored, secured, and versioned independently. Every model update is a multi-system coordination event.
- ▶ Provenance is fragmented. The score arrives back at the pipeline as a floating-point number. The model version, feature set, inference timestamp, and configuration that produced it live in separate systems — if they are captured at all.

## The Regulatory Dimension

---

For institutions subject to SR 11-7, DORA, PSD2 Article 95, or equivalent model risk management frameworks, the provenance gap is not an operational inconvenience — it is a compliance requirement. Regulators expect institutions to demonstrate, for any model-driven decision, which model version produced the score, on what inputs, at what time, under what configuration. Assembling that audit trail from logs and separate model management systems after the fact is error-prone and expensive.

## What Is Required

---

A fraud scoring architecture that is production-ready for regulated financial services must satisfy four properties simultaneously:

- ▶ Low and deterministic latency — sub-100 ms end-to-end from event ingestion to scored output in Kafka or equivalent
- ▶ Zero record loss — every transaction that enters the pipeline must produce an output event, or a guaranteed dead-letter record
- ▶ Per-event audit provenance — the model version, feature version, configuration hash, inference timestamp, and decision rationale must travel with every output event
- ▶ Transport independence — the pipeline must be deployable on Kafka, Pulsar, MongoDB, Delta Lake, or bare-metal message buses without rewriting the scoring logic

StreamKernel was designed to satisfy all four properties in a single deployable artifact.

# StreamKernel Architecture

## 03 TECHNICAL ARCHITECTURE

---

StreamKernel is a transport-agnostic event pipeline kernel. Its core runtime is a Java 21 process that orchestrates a configurable transform chain between a pluggable source and one or more pluggable sinks. The design principle is infrastructure collapse: every component that would normally be a separate service — the model server, the feature store cache, the schema registry, the audit log — runs inside the same JVM as the pipeline.

## The Transform Chain

---

Every StreamKernel pipeline is defined by a transform chain: a named, ordered sequence of transformer plugins applied to each event batch. For financial fraud scoring, the chain is:

Step	Plugin	Function
1	STRING_TO_WIREEVENT	Parses raw payment event text into a typed WireEvent envelope with transaction ID, amounts, merchant, channel, and country fields extracted
2	DJL_EMBEDDING	Runs all-MiniLM-L6-v2 (ONNX) in-process via Deep Java Library to produce a 384-dimension normalized embedding vector for each payment event
3	FRAUD_SCORE	Applies a deterministic scoring model to the embedding vector, computing a continuous fraud score [0.0–1.0] and assigning a risk band and decision with reason codes
4	KAFKA (sink)	Writes the enriched, scored event to a Kafka topic with full provenance headers attached at the producer level

## In-Process AI Inference

The DJL\_EMBEDDING transformer is the architectural centerpiece of the fraud scoring pipeline. Rather than dispatching to an external model server, StreamKernel loads the ONNX model file directly into the JVM at startup using Deep Java Library and the ONNX Runtime engine (v1.20.0). The model remains resident in memory for the pipeline's lifetime.

Each transaction's text payload is tokenized using the Hugging Face tokenizer (max 64 tokens), passed through the ONNX runtime, and the resulting embedding vector is normalized and attached to the WireEvent. The entire inference operation — tokenization, ONNX forward pass, normalization — runs on the same thread that owns the event batch, with no cross-process communication.

In-process inference means: no serialization, no network round-trip, no external service dependency.  
The model is a local file. The ONNX runtime is a JVM library. Inference is a method call.  
This is what 'infrastructure collapse' means in practice.

## MLflow Model Lifecycle

StreamKernel integrates with MLflow for model identity and lifecycle management. The champion model is referenced by alias (champion) and stage (production), and its identity is cryptographically anchored to every output event via a SHA-256 model reference hash. Model promotion and rollback can be executed at runtime without pipeline restarts — the HotSwappableDjlEmbeddingTransformer enables live model replacement with zero record loss.

**Model name:** streamkernel-minilm-onnx

**Active alias:** champion

**Active stage:** production

**Model ref SHA-256:** 12f054f7c8e13bae3db9b20761dd44e5c2c598749ae48b0f5a6ce0d7d5f7b16d

## Transport-Agnostic Sink Architecture

The same fraud scoring transform chain — `STRING_TO_WIREEVENT` → `DJL_EMBEDDING` → `FRAUD_SCORE` — runs unchanged regardless of the output destination. The sink is a pluggable component loaded at startup from the plugin catalog. In the benchmark configuration, the sink is Kafka with idempotent producer semantics (`acks=all`, `enable.idempotence=true`, `retries=MAX_INT`). The same chain can be redirected to MongoDB (vector or document sink), Delta Lake, Pulsar, Snowflake, or PostgreSQL by changing a single configuration property.

Available sinks in the benchmark build's plugin catalog:

- ▶ `KAFKA` — producer with configurable acks, lz4 compression, idempotent mode available
- ▶ `MONGO_INSERT` — document sink with WiredTiger compression
- ▶ `MONGO_VECTOR` — vector embedding sink for similarity search
- ▶ `DELTA` — Delta Lake sink for analytical workloads
- ▶ `PULSAR` — Apache Pulsar with configurable schema registry
- ▶ `POSTGRES / PGVECTOR` — relational and vector storage
- ▶ `DEVNULL` — no-op sink for benchmarking transform chain in isolation

## Security and Compliance Architecture

StreamKernel's security layer is plugin-based. The benchmark run uses `PERMIT_ALL` (no authentication) for a clean throughput baseline. Production deployments can substitute `OPA_SIDECAR` for Open Policy Agent-based per-batch authorization, or `mTLS+OPA` for mutual TLS with policy evaluation. Authorization refresh suppression, fail-closed caching, and configurable token refresh cycles are available in the security plugin interface.

# Benchmark Methodology and Configuration

## 0 4 B E N C H M A R K

### Test Configuration

The benchmark was executed using StreamKernel's automated test runner against a single pipeline configuration. All parameters were fixed for the duration of the run. No tuning was performed between the start and stop signals.

<b>Test name:</b> <code>streamkernel_financial_fraud_scoring_5m</code>	<b>Run ID:</b> <code>run-financial-fraud-01</code>
<b>Duration:</b> 5.14 minutes (300 s target)	<b>Heap:</b> 4 GB (G1GC, 2 threads, <code>MaxGCPauseMillis=50</code> )
<b>Parallelism:</b> 2 worker threads	<b>Batch size:</b> 16 records
<b>In-flight limit:</b> 768 records	<b>Executor mode:</b> <code>FIXED</code>
<b>Embedding pool size:</b> 4 predictors (batching enabled, <code>max.size=16</code> )	<b>Tokenizer max length:</b> 16 tokens
<b>Cache:</b> <code>DISABLED (NOOP)</code>	<b>Source:</b> <code>SYNTHETIC PAYMENTS</code> profile, 384-char payload
<b>Max records/sec:</b> Unlimited (no source-side throttle)	<b>Kafka partitions:</b> 12

**Kafka acks:** all (idempotent)**Kafka compression:** lz4

## Hardware

Intel Core i9-8950HK — 6 cores / 12 threads — CPU only, no GPU, no accelerator  
 All inference is CPU-bound ONNX Runtime with 3 intra-op threads, 1 inter-op thread  
 These results represent a conservative performance floor, not a ceiling.  
 GPU-accelerated or multi-core pooled configurations will materially outperform these numbers.

## Measurement Approach

Metrics were collected via StreamKernel's embedded Prometheus endpoint (port 8080) and captured to a .prom snapshot at run completion. The speedometer logs a 5-second rolling window BENCH line to stdout every 5 seconds throughout the run, providing time-series throughput and latency data independent of the Prometheus snapshot. The post-run Kafka record count was verified by consuming all partitions and comparing against the pipeline\_in\_total, pipeline\_out\_total, and pipeline\_processed\_total counters.

A config.sha256  
 (b10e981ee98c1d470b00c0390859186bb0e69ead5056504a497631daa2d79027) and  
 model.ref.sha256 were computed at pipeline startup and attached to every Kafka output record, establishing a cryptographic link between the benchmark evidence and the exact code and model artifact used.

## Benchmark Results

### 0 5 R E S U L T S

### Headline Numbers

<b>126,224</b> <b>Records Processed</b> <small>in = out = processed</small>	<b>409.4</b> <b>Avg EPS</b> <small>events per second</small>	<b>0</b> <b>Record Loss</b> <small>drops · DLQ · errors</small>	<b>0.049%</b> <b>GC Overhead</b> <small>G1GC, 3 threads</small>
---	--	---	---

<b>1,214</b> <b>MB</b> <b>Heap at Close</b>	<b>6.8 ms</b> <b>Inference Avg</b> <small>ONNX, CPU, batched</small>	<b>33 ms</b> <b>Max GC Pause</b> <small>target ≤ 50 ms</small>	<b>22</b> <b>Provenance Headers</b> <small>per output record</small>
---	--	--	--

12.9% of 4 GB			
---------------	--	--	--

## Data Integrity

```

pipeline_in_total = pipeline_out_total = pipeline_processed_total = 126,224
pipeline_dropped_total = 0
pipeline_dlq_total = 0 · pipeline_dlq_errors_total = 0
pipeline_source_errors_total = 0 · pipeline_auth_errors_total = 0
pipeline_denied_total = 0
kafka_sink_sent_ok_total = 126,224 (matches pipeline_out_total exactly)

```

Perfect record integrity was maintained for the full 5.14-minute run. Every transaction that entered the pipeline produced a scored output record in Kafka. No records were dropped due to backpressure, timeout, serialization failure, or model error.

## Per-Stage Latency

Stage	Avg Latency	Max Latency	Count
Tokenize (HF tokenizer)	1.15 ms	11 ms	15,617
ONNX inference (all-MiniLM-L6-v2)	6.8 ms	52 ms	126,224
WireEvent string encode	1.00 ms	1 ms	126,224
Kafka producer send	13.8 ms	50 ms	126,224
Kafka request latency (broker RTT)	3.9 ms	—	126,224
Kafka record queue time	11.0 ms	—	126,224

The ONNX inference stage averages 6.8 ms per record — the best result across the entire tuning series. This is achieved with a pool of 4 ONNX predictors, each running batches of 16 records, with 2 intra-op threads per predictor. The batching engine maintained 100% fill rate throughout the run, meaning no predictor ever fired an under-filled batch.

## Throughput Profile

The pipeline sustained a throughput band of 118–420 EPS across the 5-minute run, measured in 5-second rolling windows. The pipeline was fully saturated from the first window — no warm-up gap. GC-related dips were brief and transient; throughput recovered within one 5-second window every time. The pipeline never stalled, shed records, or entered backpressure hold.

Window	EPS	Window	EPS	Window	EPS
--------	-----	--------	-----	--------	-----

0:05	357.7	1:45	367.6	3:25	256.0
0:10	511.1	1:50	330.2	3:30	354.9
0:15	390.4	1:55	306.6	3:35	333.4
0:20	383.6	2:00	262.2	3:40	281.2
0:25	297.5	2:05	118.6	3:45	284.7
0:30	246.1	2:10	271.7	3:50	355.8
0:35	269.4	2:15	265.7	3:55	319.2
0:40	297.7	2:20	188.5	4:00	275.3
0:45	57.6	2:25	51.1	4:05	32.0
0:50	371.3	2:30	294.2	4:10	223.8
0:55	364.8	2:35	343.2	4:15	134.2
1:00	469.5	2:40	470.1	4:20	469.2

## JVM and GC Health

G1GC configuration targeted a 50 ms max pause time with 3 GC threads. The 33 ms max pause is well within target and reflects the higher allocation rate at 280 EPS — proportional to the 7.47 GB total allocated over the run.

Metric	Value
GC overhead at run end	0.049%
G1 Young Generation pauses (total)	18 (1 metadata threshold, 17 evacuation)
Total GC pause time	291 ms over 5.14 minutes
Max single GC pause	33 ms (target: ≤ 50 ms)
Heap used at run end	1,214 MB of 4 GB (28.3%)
Live data size (post-GC)	55.0 MB
Total heap allocated (run lifetime)	7.47 GB
Heap occupancy after last GC	2.95% of max
Live threads at close	11 (peak: 17)

## Kafka Sink Health

The Kafka producer operated across 12 topic partitions with lz4 compression and a 256 KB batch size. Kafka send averaged 13.8 ms with a 50 ms max — a significant improvement over earlier tuning runs. Broker RTT averaged 3.9 ms, returning to near-baseline levels as higher throughput amortized the partition coordination overhead.

Kafka Metric	Value
Records sent OK	126,224
Kafka request latency (avg)	3.9 ms
Kafka record queue time (avg)	11.0 ms
Kafka producer send time (avg)	13.8 ms (max: 50 ms)
Producer buffer available (end)	256 MB (fully available)
Topic partitions	12
Compression	lz4
Producer mode	acks=1, retries=MAX_INT

#### Note on partition distribution:

Partition distribution across 12 partitions was well-balanced at the throughput achieved. Higher throughput naturally smooths adaptive partitioning skew. Kafka's adaptive partitioning is enabled; this skew may be transient, driven by key distribution across the synthetic transaction IDs, or an artifact of the 6-partition topology relative to the 2-worker pipeline. Production deployments should monitor partition lag across all consumers and adjust the partitioning strategy (key-based or round-robin) to match downstream consumer group topology.

## Per-Event Provenance and Audit Trail

### 0 6   A U D I T   P R O V E N A N C E

Every record written to the streamkernel-financial-fraud-scored Kafka topic carries 22 provenance headers in addition to the JSON payload. These headers are set by the pipeline at the producer level — they are not derived from the payload and cannot be modified by downstream consumers. They travel with the record through any Kafka consumer, stream processor, or sink that preserves headers.

### Complete Provenance Header Manifest

Header Key	Example Value	Purpose
streamkernel.provenance.pipeline.id	sk-financial-fraud-scoring	Pipeline identity
streamkernel.provenance.run.id	run-financial-fraud-01	Run traceability
streamkernel.provenance.config.sha256	735d81f5c727...e84f39	Config SHA-256 fingerprint

streamkernel.provenance.model.name	streamkernel-minilm-onnx	MLflow model name
streamkernel.provenance.model.alias	champion	Active MLflow alias
streamkernel.provenance.model.stage	production	MLflow stage at inference
streamkernel.provenance.model.ref.sha256	12f054f7c8e1...7b16d	Model artifact SHA-256
streamkernel.provenance.inference.timestamp	2026-05-19T15:07:17.748Z	Nanosecond timestamp
streamkernel.provenance.transform.chain	STRING_TO_WIREEVENT, DJL_EMBEDDING,FRAUD_SCORE	Full transform chain applied
streamkernel.provenance.transform.version	payment-fraud-score-v1	Named transform version
streamkernel.provenance.feature.version	payment-risk-features-v1	Feature engineering version
streamkernel.provenance.prompt.version	not-applicable	LLM prompt version (N/A)
streamkernel.provenance.source.type	SYNTHETIC	Source plugin used
streamkernel.provenance.sink.type	KAFKA	Sink plugin used
streamkernel.provenance.sink.auth	PLAINTEXT	Transport auth mode
streamkernel.provenance.security.type	PERMIT_ALL	Authorization policy
fraud.model.version	deterministic-fraud-score-v1	Fraud model version
fraud.score	0.3700	Computed fraud score
fraud.risk_band	LOW / MEDIUM / HIGH	Risk classification
fraud.decision	APPROVE / WATCH / REVIEW	Final decision
fraud.reason_codes	CROSS_BORDER, NEW_DEVICE...	Decision rationale codes

fraud.review.threshold	0.7200	Review threshold used
------------------------	--------	-----------------------

The config.sha256 and model.ref.sha256 headers provide a cryptographic chain of custody from every output event back to the exact pipeline configuration and model artifact that produced it. For model risk management under SR 11-7 or equivalent frameworks, this means the full audit trail for any scored transaction is available by inspecting the Kafka record — no external log join required.

## Sample Fraud Scoring Output

### 07 OUTPUT EXAMPLES

The following records were sampled from the live Kafka topic immediately after the benchmark run. They represent the three risk tiers produced by the FRAUD\_SCORE transform: APPROVE (score < 0.45), WATCH (0.45 ≤ score < 0.72), and REVIEW (score ≥ 0.72).

#### Record 1 — APPROVE / LOW Risk

Field	Value
Transaction ID	TXN-100192
Account	ACCT-15952
Amount	USD 275.04
Merchant / Channel	grocery / bill_pay
Country	AE (United Arab Emirates)
<b>Fraud Score</b>	<b>0.37</b>
<b>Risk Band</b>	<b>LOW</b>
<b>Decision</b>	<b>APPROVE</b>
Reason Codes	CROSS_BORDER, NEW_DEVICE

#### Record 2 — WATCH / MEDIUM Risk

Field	Value
Transaction ID	TXN-100193
Account	ACCT-15983
Amount	USD 276.41

Merchant / Channel	fuel / card_present
Country	SG (Singapore)
<b>Fraud Score</b>	<b>0.61</b>
<b>Risk Band</b>	<b>MEDIUM</b>
<b>Decision</b>	<b>WATCH</b>
Reason Codes	CROSS_BORDER, NEW_DEVICE, HIGH_VELOCITY

## Record 3 — REVIEW / HIGH Risk

Field	Value
Transaction ID	TXN-100194
Account	ACCT-16014
Amount	USD 277.78
Merchant / Channel	travel / card_not_present
Country	BR (Brazil)
<b>Fraud Score</b>	<b>0.99</b>
<b>Risk Band</b>	<b>HIGH</b>
<b>Decision</b>	<b>REVIEW</b>
Reason Codes	CROSS_BORDER, UNKNOWN_DEVICE, ELEVATED_VELOCITY, SANCTIONS_SCREEN_HIT, HIGHER_RISK_CHANNEL

The FRAUD\_SCORE transform applies five independent risk signals — velocity, device trust, channel risk, sanctions screening, and cross-border exposure — and aggregates them into a continuous [0.0–1.0] score.

The two configurable thresholds (WATCH  $\geq$  0.45, REVIEW  $\geq$  0.72) are externalized in the pipeline configuration and captured in the fraud.review.threshold provenance header on every record, so downstream systems can reconstruct the decisioning logic without access to the pipeline source.

## Financial Services Deployment Scenarios

### 08 USE CASES

StreamKernel's architecture is particularly well-suited to financial services workloads where the combination of low latency, regulatory auditability, and transport flexibility is non-negotiable.

## Card Network Real-Time Authorization

---

Card authorization networks process thousands of transactions per second with sub-100 ms SLAs on scoring decisions. StreamKernel's in-process inference eliminates the external model-serving hop, reducing the scoring path to tokenization + ONNX inference + rule evaluation — all within the same JVM. The pipeline can be deployed as a sidecar to the authorization engine or as a standalone scoring service with a Kafka input/output interface.

## Bank Transaction Monitoring

---

Anti-money laundering and fraud monitoring systems at retail banks consume transaction feeds from core banking systems and apply behavioral and semantic scoring to detect anomalous patterns. StreamKernel can consume these feeds from Kafka, Pulsar, or REST endpoints, apply embedding-based semantic similarity scoring, and write enriched events to Delta Lake or MongoDB for downstream analytics — all with per-event audit headers that satisfy model risk management documentation requirements.

## Payment Processor Chargeback Prediction

---

Payment processors seeking to predict chargeback risk at authorization time can deploy StreamKernel as an enrichment layer between the payment gateway and the downstream risk engine. The embedding vector produced by DJL\_EMBEDDING can feed both the real-time scoring model and a vector store (MONGO\_VECTOR or PGVECTOR) for post-hoc similarity analysis against historical chargeback patterns.

## Fintech and Neo-Bank Deployment

---

For fintechs operating without legacy infrastructure, StreamKernel's transport-agnostic architecture means the entire fraud scoring pipeline can be deployed against a managed Kafka service (Confluent Cloud, MSK, Aiven) or Pulsar without on-premises infrastructure. The MLflow integration provides model governance without requiring a dedicated model platform — a single MLflow tracking server is sufficient.

## Air-Gapped and High-Security Environments

---

For institutions operating in air-gapped or highly restricted network environments — including defense-adjacent financial services, sovereign wealth funds, and central banks — StreamKernel's fully in-process architecture is a significant operational advantage. The model, tokenizer, and scoring logic are packaged in a single fat JAR. There are no outbound calls to model serving infrastructure, no cloud dependencies, and no runtime model downloads. The pipeline operates as a standalone process with local file access only.

## Conclusion

### 09 CONCLUSION

---

Real-time fraud scoring at production scale requires an architecture that resolves three tensions simultaneously: inference speed versus infrastructure complexity, throughput versus auditability, and flexibility versus operational overhead. External model serving infrastructure solves none of

these tensions — it trades inference latency for operational complexity, and sacrifices auditability for architectural separation.

StreamKernel resolves all three by collapsing the inference infrastructure into the pipeline. The benchmark results presented in this paper demonstrate that a CPU-only laptop running StreamKernel can process 126,224 payment transactions in 5 minutes at 409 EPS average — with zero record loss, 6.8 ms average ONNX inference latency, 100% batch fill rate, and full 22-field cryptographic provenance on every output event.

These are conservative numbers. Production hardware — multi-socket servers, higher core counts, expanded embedding pools, GPU ONNX execution providers — will produce proportionally higher throughput with equivalent or better integrity guarantees.

StreamKernel is available at: [github.com/IntuitiveDesigns/StreamKernel-io](https://github.com/IntuitiveDesigns/StreamKernel-io)  
For licensing, partnership, or enterprise deployment inquiries: [lopezstevie@gmail.com](mailto:lopezstevie@gmail.com)  
LinkedIn: [linkedin.com/in/steve-lopez-b9941](https://www.linkedin.com/in/steve-lopez-b9941)

### Benchmark Summary

126,224 records · 5.14 min · 409.4 EPS  
avg · 511 EPS peak  
Zero record loss · Zero errors · Zero DLQ  
0.049% GC · 1,214 MB heap · 33 ms max  
pause

### 22 provenance headers per event

Cryptographic config + model SHA-256  
MLflow model lifecycle integration  
Transport-agnostic · CPU-only deployable