

# Low-Latency AI Inference Tooling for Personalized Product Recommendations at Checkout

Udit Agarwal

[udit15@gmail.com](mailto:udit15@gmail.com)

## Abstract:

This technical report addresses the architectural and tooling requirements necessary to achieve low-latency inference for personalized product recommendations within high-traffic e-commerce checkout flows. The temporal constraints imposed by the point-of-sale environment are exceptionally severe; empirical evidence suggests that system delays as minor as 100 milliseconds (ms) can correlate directly with a 1% loss in sales. Consequently, meeting a stringent Service Level Objective (SLO) for the worst-case tail latency (P99) of less than 500 ms is non-negotiable for maximizing conversion rates.

Modern recommendation services, predominantly based on Deep Learning Recommendation Models (DLRMs), face inherent difficulties due to their reliance on vast embedding tables (EMTs) that can exceed a Terabyte in size. This challenge necessitates an algorithm-system co-design approach. The analysis identifies two primary tooling mechanisms for mitigating these bottlenecks: Model Quantization and Inference Compilation. Quantization, specifically employing INT4 precision, has demonstrated the capacity to reduce DLRM model size from 12.58 GB to 1.57 GB on Terabyte datasets while maintaining competitive accuracy. This size reduction directly addresses memory bandwidth constraints. Concurrently, Inference Compilation, using runtimes like ONNX, delivers typical inference speedups ranging from 2x to 10x. These optimized models must be deployed within a fault-tolerant Microservices architecture that leverages low-overhead communication protocols such as gRPC to ensure high throughput and minimize inter-service communication friction. The synergistic deployment of these strategies is essential for operationalizing personalized recommendations without compromising the integrity of the real-time checkout experience.

**Keywords:** Recommendation Systems, Deep Learning Recommendation Models (DLRM), Low Latency, Quantization, Inference Compilation, Microservices, E-commerce Checkout, P99 Latency.

## I. INTRODUCTION: THE CRITICALITY OF REAL-TIME RECOMMENDATIONS AT CHECKOUT

### I. A. Contextualizing Recommendation Systems in E-commerce

Recommendation systems have emerged as vital information filtering techniques designed to alleviate the difficulties inherent in "information overload," guiding users toward suitable products, content, or services. These systems have evolved significantly, progressing from simpler rule-based and nearest-neighbor designs to sophisticated deep learning approaches capable of highly personalized outputs.

A key distinction must be made regarding the deployment location of the recommendation engine within the e-commerce journey. **In-store recommenders** guide users while they browse or are on the main product pages. In contrast, **checkout recommenders** are applications specifically engineered to suggest items just before the user finalizes the payment transaction. This placement—at the culmination of the buyer journey—imposes unique and acute requirements for both speed and reliability. The AI response must be delivered quickly enough that it enhances the purchase without introducing friction, which is highly detrimental at the final decision point.

## I. B. The Low-Latency Imperative: Conversion and Revenue

The relationship between system performance and commercial outcomes is thoroughly documented. Latency, defined as the time lag between a shopper's action and the system's subsequent response, introduces friction and frustration, analogously to a pause encountered in a physical checkout queue. This delay is not merely an inconvenience; it represents a significant, measurable threat to revenue. Research indicates that every 100 milliseconds (ms) of delay can result in a corresponding 1% loss in sales. This establishes a clear, immediate financial incentive for achieving superior performance.

The consequences of poor performance are observed directly in abandonment rates. Slow page loads and unresponsive processes contribute significantly to cart abandonment, a phenomenon already averaging 70.22% across e-commerce platforms, with this figure surging higher during periods of system degradation or server overload. Since the recommendation occurs during the checkout process itself, the AI inference component must operate within a microscopic latency budget. The goal is to ensure the recommendation adds value without incurring a revenue-damaging delay that causes the system to exceed human perception thresholds (which typically demand sub-second response times). The necessary architectural and tooling optimizations detailed in this report are focused explicitly on achieving sub-second performance in AI inference, as fast responses are proven to boost customer satisfaction, conversion rates, and overall trust in the brand.

## II. DEFINING PERFORMANCE BENCHMARKS AND METRICS

In high-stakes, real-time environments, performance success is not measured by average speed, but by consistency across all user experiences. This demands a rigorous adherence to percentile-based latency objectives and specialized AI metrics.

### A. Quantified Latency Objectives (P50 and P99)

E-commerce microservices require stringent Service Level Objectives (SLOs) focused on eliminating poor user experiences, which are quantified using tail latency metrics, specifically the 99th percentile (P99). The P99 latency measures the performance experienced by the slowest 1% of users, and minimizing this tail is crucial for maintaining a consistently positive interaction.

For core e-commerce transactions, such as the checkout process itself, an explicit SLO target often cited is **"99th percentile checkout latency < 500 ms"**. This benchmark ensures that nearly all customers experience a smooth process, preventing the multi-second waits that inevitably lead to cart abandonment. Furthermore, for specialized real-time APIs, such as those powering conversational commerce or instant nudges, the targets are even more demanding: P50 (median latency) must be kept below 50 ms, and the P99 (worst-case latency) must remain below 300 ms. If the P99 latency is allowed to exceed 2 seconds, high user frustration and the perception of a sluggish application become guaranteed outcomes. The optimization strategies must therefore be calibrated to meet these tight, non-linear performance thresholds.

### B. Critical AI Inference Metrics

Beyond general API benchmarks, real-time AI inference requires dedicated metrics to measure both responsiveness and capacity.

1. **Time to First Token (TTFT):** TTFT measures the initial time a user must wait before the model's output begins to stream. This metric is particularly relevant for recommendation systems that now incorporate sophisticated components like conversational search or detailed product analysis. A low TTFT is essential for perceived responsiveness. TTFT inherently includes request queuing time, network latency, and the model's prefill time. The prefill phase, which involves computing the input sequence and generating the key-value (KV) cache for sequence generation, means that longer input prompts (such as a large cart history or extensive user context) will result in a larger TTFT.

2. **End-to-End Request Latency (E2E Latency):** This comprehensive metric captures the total time elapsed from the submission of the user query until the full response is received. It encompasses the performance of various systemic elements, including batching mechanisms, queuing delays, and network latencies. Optimizing E2E latency directly minimizes the risk of catastrophic cart abandonment caused by perceived sluggishness.

3. **Throughput:** Throughput defines the system's capacity, quantifying the number of tasks or queries the model can successfully handle within a specific time frame. For recommendation systems, which serve a potentially vast number of concurrent users, high throughput is essential for handling traffic spikes and maintaining consistent service quality.

Table 1 summarizes the required performance envelope for real-time commerce microservices, demonstrating the tight constraints that inform all subsequent architectural decisions.

Table 1: Latency Benchmarks for E-commerce Microservices

API Type/Metric	P50 (Median) Goal	P99 (Worst-Case) Goal	Business Implication
E-commerce Checkout Latency	< 200 ms	< 500 ms (Targeted SLO)	Every 100 ms delay costs approximately 1% of sales, establishing a clear ROI for optimization.
Real-Time APIs (General)	< 50 ms	< 300 ms	Benchmark necessary for conversational components and instant real-time nudges.
End-to-End Latency (E2E)	System response delay (ms)	Incorporates queuing, prefill time, and network latencies.	Minimizing risk of catastrophic cart abandonment (70%+ baseline rate) due to perceived sluggishness.
Time to First Token (TTFT)	Time until initial output is visible.	Larger for longer input sequences due to KV-cache computation.	

The need for highly constrained TTFT, even in systems traditionally focused on numerical prediction, confirms that modern personalization often involves complex models that must manage dynamic context lengths efficiently (e.g., the shopper's current cart and history) while rigidly adhering to the P99 latency budget. This necessitates a solution that fundamentally reduces the computational cost of the models themselves.

### III. ARCHITECTURAL CHALLENGES OF DEEP LEARNING RECOMMENDATION MODELS

The inherent complexity and scale of state-of-the-art Deep Learning Recommendation Models (DLRMs) present the most significant physical barriers to achieving low-latency serving. Optimizing these models requires addressing fundamental memory and bandwidth constraints.

#### A. The DLRM Memory and Bandwidth Bottleneck

DLRMs are pervasive in contemporary cloud services, representing a cornerstone of personalized user experience across social networks, streaming platforms, and e-commerce. The demand they place on infrastructure is substantial; Meta reported that DLRMs consume over 60% of their production AI inference cycles.

The architectural challenge stems primarily from the need to encode high-dimensional categorical features, such as user IDs or item attributes, into dense vector representations using embedding layers. The resultant Embedding Tables (EMTs) are enormous, frequently reaching tens of gigabytes or scaling up to Terabytes in real-world production settings. During inference, DLRMs execute multiple embedding lookups and reductions, resulting in frequent and irregular memory accesses. This process is highly intensive on memory capacity and bandwidth, creating memory bottlenecks that directly increase operational latency.

### ***B. Hybrid Architectures and Communication Friction***

To manage EMTs that exceed the capacity of a single GPU, existing studies frequently adopt a **CPU-GPU hybrid architecture**. In this setup, the massive EMTs reside in the CPU memory, while the GPUs are reserved for highly parallel computations.

While solving the capacity problem, this architecture introduces a new performance constraint: unavoidable communication overhead between the CPU and the GPU. This inter-device transfer can cause GPU stalls, ultimately compromising the overall real-time performance necessary for checkout recommenders. Furthermore, maintaining the accuracy of these models depends crucially on the frequency of updates. Large model size introduces significant latency challenges for propagating updates across geo-distributed servers. Solutions are required to minimize the required write bandwidth for rapid personalization, providing accuracy comparable to that of a fully fresh model.

### ***C. Algorithmic Efficiency: The Two-Stage Approach***

Because the infrastructure cost of iterating over the entire catalog for every user request is prohibitive, algorithmic strategies are crucial for optimization. The two-stage retrieval and ranking funnel has emerged as a key technique for improving inference efficiency.

This method employs a **lightweight retriever** that quickly processes the input (which may include sequential or incremental data) to narrow down the pool of candidate items. The subsequent, computationally heavier **ranker** then processes only this significantly reduced candidate set. This funnel dramatically cuts the computational cost and latency compared to applying a complex model directly across the full item catalog. This process is vital because relying solely on increasing hardware capacity to solve the latency problem is often financially and architecturally infeasible due to the terabyte-scale memory demand and the fundamental bandwidth limitations of hybrid serving. Therefore, optimization must be a co-design process involving both algorithmic simplification and physical model compression.

## **IV. MODEL OPTIMIZATION TOOLING FOR ACCELERATED INFERENCE**

Achieving the required P99 latency for checkout recommendations mandates the application of specialized AI inference tooling designed to optimize models below their native production fidelity. These tools provide fundamental reductions in both data size and execution time.

### ***A. Precision Reduction via Model Quantization***

Quantization is the process of reducing the numerical precision of a neural network's weights and activations. This mapping uses a scale factor and a zero-point to represent continuous values within a smaller, discrete range.

1. **Benefits and Techniques:** Quantization delivers significant gains by reducing memory footprint, accelerating inference execution, and lowering energy consumption. The primary driver of speedup is the ability to leverage hardware-accelerated integer arithmetic.
  - **Post-Training Quantization (PTQ):** The model is converted after it has been fully trained, affecting only the inference state. PTQ requires a calibration dataset to pre-compute activation statistics (e.g., minimum and maximum values) used to define the quantization parameters. Advanced PTQ methods, such as SmoothQuant, address quantization difficulty by migrating activation outliers to weights, enabling efficient 8-bit weight, 8-bit activation quantization.
  - **Quantization-Aware Training (QAT):** This technique simulates quantization effects during the training process, allowing the model to adapt and minimize accuracy loss, generally yielding superior results when precision loss is a concern.
2. **Quantifiable Impact on DLRMs:** Quantization is particularly critical for recommendation systems because it streamlines the massive, daily computations involved in processing millions of user-item interactions. For DLRMs, which are memory-bandwidth constrained, the benefit is dramatic: quantized models (DQRM) demonstrated a massive reduction in size on the Terabyte dataset, shrinking from 12.58 GB to just **1.57 GB**. This nearly 8x reduction in model size directly alleviates the DLRM memory bottleneck and is indispensable for rapid deployment in memory-constrained serving environments.

### B. Inference Compilation and Runtime Acceleration

Model compilation provides a second layer of optimization, focusing on optimizing the execution graph of the model for target hardware, ensuring the fastest possible execution path.

- 1. Compiler Tools and Techniques:** Tools such as ONNX Runtime are used to convert models from training frameworks (e.g., PyTorch) into an optimized, platform-agnostic format. This process generates an efficient graph suitable for high-speed inference. Further optimizations, such as `torch.compile` and CUDA graph capture, can be applied, although they require careful management to address dynamic memory behaviors.
- 2. Performance Gains:** Compilation provides substantial, quantifiable speedups in inference time. Users commonly report **2x to 10x faster inference** when deploying models via ONNX Runtime on CPU, with even greater acceleration possible when utilizing GPU execution providers like TensorRT.
- 3. Synergistic Optimization:** The most effective approach for achieving extreme low-latency targets involves combining quantization and compilation. Quantization addresses the data size and bandwidth constraints, while compilation optimizes the execution pathway, resulting in a synergistic strategy that simultaneously minimizes memory usage and maximizes execution speed.

Table 2 highlights the dual importance of these optimization tools, demonstrating how compression (quantization) and execution efficiency (compilation) together contribute to the required performance envelope.

Table 2: Low-Latency Optimization Techniques for Recommendation Models

Technique	Description	Primary Benefit	Observed Gains/Impact
Quantization	Reduces numerical precision of weights/activations (FP32 to 8-bit integers).	Reduced Memory Footprint; Faster Compute	DLRM size reduction from 12.58 GB (FP32) to 1.57 GB (INT4).
Inference Compilation (ONNX, torch.compile)	Optimizes graph execution and leverages hardware acceleration.	Inference Speedup; Execution Efficiency	Typical gains of 2x–10x faster inference on CPU.
Specialized Serving Servers	Optimize memory usage, efficient batching, and GPU support (e.g., NVIDIA Triton, TorchServe).	High Throughput; Low Latency	Provides a proprietary engine for maximum performance; specific platforms show up to 2.3x faster speed.

Model optimization is thus observed as a layered process: quantization addresses the fundamental data size constraint, while compilation ensures execution path efficiency. Both layers are indispensable for achieving the maximal reduction in End-to-End latency required to hit the stringent real-time targets.

## V. SCALABLE SERVING ARCHITECTURE AND DEPLOYMENT

Model-level optimizations must be supported by a robust, cloud-native architecture capable of managing high-volume, concurrent traffic while maintaining resilience. The shift to a microservices architecture is essential, but it introduces communication overheads that must be aggressively mitigated.

### A. Microservices for Scalability and Resilience

The deployment of low-latency AI inference components is best supported by a microservices architecture. This design pattern provides several critical advantages for high-volume e-commerce platforms:

1. **Fault Tolerance:** The decoupled nature of microservices is crucial for minimizing the impact of potential service failures. If the recommendation engine encounters an error, other core functions, such as payment gateways, continue to operate seamlessly. This resilience is paramount in e-commerce, where system downtime results in significant revenue loss and damage to customer trust.
2. **Agility and Modularity:** Each AI model can be encapsulated within a container, allowing it to be updated, replaced, or experimented with independently of other services. This modularity reduces the risk associated with updates and facilitates the rapid deployment of new algorithms.
3. **Independent Scaling:** Components can scale independently based on demand, which is essential for handling unpredictable traffic spikes typical of major e-commerce events (e.g., Black Friday) and maintaining the required high throughput.

### ***B. Overcoming Communication Overhead***

While microservices provide resilience and scaling benefits, they introduce complexity, notably inter-service communication overhead, which can compromise the stringent low-latency goals. The latency budget is razor-thin, demanding sophisticated infrastructure choices.

1. **Low-Overhead Protocols:** To mitigate communication friction, lightweight protocols such as **gRPC** must be employed. gRPC is critical for optimizing communication efficiency and reducing the time spent on data serialization compared to traditional protocols like REST/HTTP, ensuring that the computational gains achieved through model optimization are not lost in network transfer time.
2. **Specialized Serving Engines:** To maximize throughput under latency constraints, specialized inference servers are necessary. Tools such as NVIDIA Triton Inference Server, TensorFlow Serving, or TorchServe are specifically designed for low-latency inference, optimizing memory usage and efficiently batching requests, often with direct GPU/TPU support. The deployment of these optimized serving layers is often streamlined through **Inference as a Service (IaaS)** platforms, which abstract the complexities of GPU provisioning and infrastructure management, ensuring reliable, instant scalability regardless of the workload. IaaS provides a competitive edge by allowing teams to focus on improving model quality rather than managing underlying hardware.

### ***C. Data Flow Orchestration and Continuous Feedback***

An effective low-latency system must be fed by real-time data to ensure relevance and timeliness. This requires robust data flow orchestration.

Handling real-time data, potentially sourced from interaction logs or immediate user actions, necessitates timely synchronization mechanisms, often implemented using streaming platforms like Apache Kafka, to prevent latency issues caused by stale information. Furthermore, the architecture must support continuous learning and refinement. User feedback, such as item ratings, must be integrated into the data access layer and used to iteratively adjust and refine the recommendation algorithm, thus continuously improving relevance and accuracy.

The foundational architectural choice of microservices introduces a systemic trade-off: resilience and flexibility are gained, but communication overhead is incurred. The successful implementation of low-latency systems relies on aggressively mitigating this trade-off using protocol optimization (gRPC) and specialized serving engines, thereby translating technological sophistication into reliable real-time performance.

## **VI. CONCLUSION**

The reliable provision of personalized product recommendations at the critical e-commerce checkout juncture requires an integrated, multi-layered optimization strategy that addresses both the algorithmic constraints of the recommendation model and the operational realities of a distributed serving infrastructure. The evidence confirms that this challenge is fundamentally defined by the extreme sensitivity of e-commerce conversions to latency, where delays must be kept below the 500 ms P99 threshold to avoid measurable revenue erosion. The core technical impediment lies in the massive scale of Deep Learning Recommendation Models (DLRMs), whose Terabyte-scale embedding tables exceed traditional hardware capacity, resulting in memory and bandwidth bottlenecks within hybrid CPU-GPU architectures.

The comprehensive solution rests on three pillars of technological enhancement:

1. **Model Compression:** The adoption of **Quantization** techniques significantly reduces the model footprint (up to 8x for DLRMs), effectively alleviating memory access and bandwidth limitations inherent to DLRM serving.
2. **Execution Acceleration:** The use of **Inference Compilation** tools, such as ONNX Runtime, delivers guaranteed execution speedups, typically between 2x and 10x. This layer compounds the benefits of compression, ensuring maximum execution efficiency at the kernel level.
3. **Resilient Serving:** Deployment within a **Microservices architecture** ensures fault tolerance and independent scaling. Crucially, communication friction must be minimized using lightweight protocols like **gRPC** and specialized inference servers (e.g., NVIDIA Triton) to maximize throughput and sustain the required sub-second P99 latency targets.

By adopting these specific tooling and architectural strategies, e-commerce platforms can reliably transition from traditional, sluggish serving models to high-speed AI inference capable of delivering personalized, conversion-boosting recommendations during peak-traffic periods, ultimately securing a consistently superior customer experience.

### Future Directions

While the immediate focus is on latency and throughput, optimizing the entire AI lifecycle also involves considering efficiency and sustainability. Further research must quantitatively explore the trade-offs between model performance and resource consumption, specifically evaluating **Resource Utilization** to maximize cost-efficiency and assessing the quantitative **environmental impact** of increasingly complex deep learning algorithms utilized in production recommender systems.

### REFERENCES:

1. P. Resnick and H. R. Varian, "Recommender systems," *Commun. ACM*, vol. 40, no. 3, pp. 56-58, 1997.
2. P. Resnick, N. Iacovou, M. O. Hailes, P. Herlocker, and J. A. Riedl, "An architecture for collaborative filtering of netnews," in *Proceedings of the 1994 ACM Conference on Computer Supported Cooperative Work*, 1994, p. 94.
3. M. Marcus, "Immigrant Voters in Massachusetts: Implications for Political Parties," Harvard Kennedy School, 2007. 34
4. H. T. T. Naumov et al., "Deep learning recommendation model for personalization and recommendation systems," *arXiv:1908.08284 [cs.LG]*, 2019. 5
5. P. Covington, J. Adams, and E. Sargin, "Deep neural networks for youtube recommendations," in *Proceedings of the 10th ACM Conference on Recommender Systems*, 2016.
6. P. Gupta et al., "GSP: General Search Platform," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020.
7. Elgafy and K. Lafdi, "Nanoparticles and fiber walls interactions during nanocomposites fabrication," *Journal of Scientific Conference Proceedings*, vol. 2, no. 1, pp. 15–23, 2010. 35
8. M. S. Shareef, A. Ojo, and T. Janowski, "Exploring digital divide in the Maldives," in J. Berleur, M. D. Hercheui, & L. M. Hilty (Eds.), *What Kind of Information Society? Governance, Virtuality, Surveillance, Sustainability, Resilience*, pp. 51–63, Springer, 2010. 35
9. K. B. Jang, "Deconstructing the opposition of natural/arbitrary in Coleridge's theory of language," [Paper presentation], *NASSR 2019: Romantic Elements*, Chicago, IL, United States, 2019. 35
10. (DQRM Technical Report), "We show that combining gradient sparsification and quantization together significantly reduces the amount of communication," *arXiv:2410.20046*, 2024. 8
11. Z. Wu et al., "Designing Microservices Architectures for Scalable AI in E-commerce Applications," *ResearchGate*, 2024.
12. Y. Yang et al., "Revisiting the optimization landscape for DLRM inference," *IEEE Micro*, 2023.
13. T. Chen et al., "DeepRecInfra: An End-to-End Modeling Infrastructure for Neural Personalized Recommendation," *arXiv:2001.02772*, 2020.

14. T. Chen et al., "DeepRecSched: A Recommendation Inference Scheduler for Maximizing Latency-bounded Throughput," in Proceedings of the 47th International Symposium on Computer Architecture (ISCA), 2020. 36
15. Y. Zhong, S. Liu, J. Chen, J. Hu, Y. Zhu, X. Liu, X. Jin, and H. Zhang, "DistServe: Disaggregating Prefill and Decoding for Goodput-optimized Large Language Model Serving," arXiv:2401.09670, 2024. 28
16. H. Peng, H. Xu, S. H. Song, J. Zhou, B. Wang, M. Chen, G. Zong, and J. Wu, "Graph-Convolved Factorization Machines for Personalized Recommendation," IEEE Transactions on Knowledge and Data Engineering, 2021. 37
17. J. McAuley et al., (Paper on TransFM and S-BPR optimization), RecSys 2018. 38
18. D. G. T. B. J. A. Yu, Z. Jiang, D.-D. Chen, S. Feng, D. Li, Q. Liu, and J. Yi, "Leveraging Tripartite Interaction Information from Live Stream E-Commerce for Improving Product Recommendation," in Proceedings of the 27th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2021. 39
19. Y. Zhong, et al., "SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models," arXiv:2312.07104 [cs.AI], 2023.