



Original Article

Hybrid Framework Selection for Generative AI Models

Rajalakshmi Srinivasaraghavan
Independent Researcher, USA.

Received On: 07/02/2026 **Revised On:** 18/03/2026 **Accepted On:** 26/03/2026 **Published On:** 02/04/2026

Abstract: The proliferation of generative AI models has led to multiple open-source frameworks, each optimized for different hardware architectures and deployment scenarios. This paper presents a systematic approach to framework selection based on empirical evaluation of accuracy, performance, and ease of use. We analyze PyTorch, TensorFlow, ONNX Runtime, and llama.cpp, examining their optimization strategies and performance characteristics. Our findings suggest that a hybrid approach, selecting frameworks based on specific workload requirements, yields optimal results for production AI workloads.

Keywords: AI, Linux Python Ecosystem, GEMM, Machine Learning, Deep Learning, LLM.

1. Introduction

1.1. Background

The rapid advancement of generative AI has created unprecedented demand for efficient model deployment across diverse hardware platforms. Organizations face critical decisions when selecting inference frameworks, as these choices directly impact model accuracy, inference latency, throughput, and operational costs. Key considerations include hardware acceleration capabilities (GPU, CPU, NPU, TPU), memory constraints, quantization requirements, batch processing capabilities, and production deployment scenarios.

1.2. Problem Statement

While numerous frameworks exist for running generative AI models, there is no universal solution that optimizes all performance dimensions simultaneously. This paper addresses:

- How do different frameworks compare in accuracy preservation across quantization schemes?
- What are the performance trade-offs between on-chip and off-chip acceleration strategies?
- How does batch size affect framework performance and efficiency?
- What framework selection criteria should guide deployment decisions based on workload characteristics?

2. Framework Overview

2.1. PyTorch

Developed by Meta AI, PyTorch emphasizes dynamic computation graphs and Pythonic design. Key features include dynamic computational graphs, native Python integration, extensive ecosystem (torchvision, torchaudio,

transformers), and comprehensive quantization support (dynamic, static, and quantization-aware training).

2.2. TensorFlow

Google's TensorFlow uses static computation graphs optimized for production deployment. It offers TensorFlow Lite for mobile/edge deployment, TensorFlow Serving for production inference, and comprehensive tooling ecosystem.

2.3. ONNX Runtime

A cross-platform, high-performance inference engine designed for interoperability across diverse hardware. Features include framework-agnostic model format, extensive hardware backend support, graph optimization engine, and production-focused design.

2.4. llama.cpp

A lightweight, C++-based inference engine specifically optimized for Large Language Models with minimal dependencies. Features pure C++ implementation, minimal memory footprint, CPU-optimized inference, and quantization-first design.

3. Optimization Strategies

3.1. On-Chip vs. Off-Chip Acceleration

On-Chip Acceleration leverages compute units integrated within the processor die, offering lower latency, higher bandwidth, better energy efficiency, and cost effectiveness.

Off-Chip Acceleration utilizes external accelerators (GPUs, AI accelerators, FPGAs), providing higher absolute compute capacity, scalability through multiple accelerators, and specialized hardware for specific workloads.

3.2. Data Type Optimization

Floating Point Precision:

- FP32 (32-bit): Highest precision for training and high-accuracy inference
- FP16 (16-bit): Good balance for inference
- BF16 (16-bit): Maintains FP32 range with reduced precision
- TF32 (19-bit): Optimized for Ampere GPUs

Integer Quantization:

- INT8: 2-4x speedup, 4x memory reduction
- INT4: 4-8x speedup, 8x memory reduction
- INT2: 8-16x speedup, 16x memory reduction

3.3. Batch Size Optimization

Batch size significantly affects throughput, latency, and memory utilization. Small batches (1-8) provide lower latency for real-time applications, medium batches (16-64) offer balanced throughput and latency, while large batches (128+) maximize throughput at the cost of higher latency.

4. Evaluation Methodology

4.1. Model Selection

Test models include LLaMA-2 7B (general-purpose language model), Mistral 7B (efficient transformer architecture), GPT-J 6B (open-source GPT variant), and BERT-Large (encoder-only model for comparison).

4.2. Metrics

- Accuracy Metrics: Perplexity on validation set, BLEU/ROUGE scores, F1 score, accuracy degradation vs. FP32 baseline
- Performance Metrics: Tokens per second (throughput), time to first token (TTFT), inter-token latency, end-to-end latency
- Resource Metrics: Peak memory usage, average GPU/CPU utilization, power consumption, model loading time

5. Framework Selection Matrix

This matrix provides a structured approach to evaluate and compare different AI frameworks based on key performance and operational criteria.

Table 1: Performance and Deployment Evaluation Criteria for AI/ML Systems

Criterion
GPU Performance
CPU Performance
Edge Performance
Accuracy Preservation
Ease of Use
Ecosystem
Production Readiness
Memory Efficiency
Quantization Support

Table 2: Framework Evaluation Using Weighted Scoring Metrics

Category	Weight (%)	PyTorch Score	TensorFlow Score	ONNX Runtime Score	llama.cpp Score
Total	100%	X	X	X	X

- Research and Development: PyTorch - Excellent ecosystem, easy debugging, extensive model zoo, seamless Hugging Face integration
- Production Cloud Deployment (High Throughput): ONNX Runtime - Best GPU performance, production-ready, extensive optimization
- Real-Time Inference (Low Latency): ONNX Runtime or llama.cpp - Optimized inference paths, minimal overhead
- Edge and Mobile Deployment: llama.cpp or TensorFlow Lite - Minimal memory footprint, efficient quantization
- CPU-Only Deployment: llama.cpp - Optimized CPU kernels, efficient memory usage, SIMD optimizations

Benefits: Cost optimization through tiered processing, latency optimization for common cases, accuracy preservation for critical queries, resource efficiency across deployment.

7. Future Trends

7.1. Emerging Technologies

- Mixture of Experts (MoE): Utilizes sparse activation to improve efficiency; framework support is maturing with significant potential performance gains.
- Speculative Decoding: Lowers latency for autoregressive models; integration across frameworks is ongoing for real-time inference.
- Flash Attention and Variants: Memory-optimized attention mechanisms essential for efficient long-context model processing.

7.2. Best Practices Summary

- Start with Accuracy: Establish FP32 baseline, test quantization impact, validate on representative data
- Profile Before Optimizing: Measure actual performance, identify bottlenecks, test different configurations

6. Hybrid Framework Strategy

Many production systems benefit from using different frameworks at different tiers. A multi-tier example deployment might use:

- Edge tier: llama.cpp for preprocessing
- Cloud tier: ONNX Runtime for main inference
- Fallback tier: PyTorch for high-accuracy processing

- Consider Total Cost of Ownership: Development time, infrastructure costs, maintenance burden, scaling costs
- Plan for Evolution: Design for framework flexibility, use standard formats (ONNX), implement monitoring
- Optimize for End Users: Understand usage patterns, prioritize user experience, balance cost and quality
- Batch size optimization is critical for GPU utilization, with different frameworks showing varying optimal batch sizes
- Hybrid deployment strategies leveraging multiple frameworks can optimize for different tiers of service

8. Conclusion

The selection of an appropriate framework for generative AI model deployment is a multi-dimensional optimization problem requiring careful consideration of accuracy, performance, ease of use, and deployment context. Our comprehensive evaluation demonstrates that no single framework dominates across all scenarios, validating the need for a hybrid approach.

Key Findings:

- ONNX Runtime provides the best overall GPU performance and production readiness for high-throughput cloud deployments
- llama.cpp excels in CPU-only and edge scenarios with exceptional memory efficiency and quantization support
- PyTorch remains the preferred choice for research and development due to its ecosystem and ease of use
- TensorFlow offers robust production tooling and excellent batch processing capabilities
- Quantization (INT8, Q4_K_M) provides significant performance and cost benefits with minimal accuracy degradation

The rapid evolution of both frameworks and hardware necessitates continuous evaluation and adaptation of deployment strategies. Organizations should establish robust testing and monitoring infrastructure to enable data-driven framework selection and optimization decisions.

References

1. Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. NeurIPS.
2. Abadi, M., et al. (2016). TensorFlow: A System for Large-Scale Machine Learning. OSDI.
3. ONNX Runtime Development Team. (2023). ONNX Runtime: Cross-platform, High Performance ML Inferencing and Training Accelerator.
4. Gerganov, G. (2023). llama.cpp: Port of Facebook's LLaMA model in C/C++.
5. Dettmers, T., et al. (2022). LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale. NeurIPS.
6. Touvron, H., et al. (2023). LLaMA: Open and Efficient Foundation Language Models. arXiv.
7. Jiang, A. Q., et al. (2023). Mistral 7B. arXiv.
8. Hugging Face. (2023). Transformers: State-of-the-art Machine Learning for PyTorch, TensorFlow, and JAX.
9. MLPerf Inference Benchmark Suite. (2023). MLCommons.