

Comparative Analysis of Large Language Model Inference Serving Systems: A Performance Study of vLLM and HuggingFace TGI

Saicharan Kolluru
Baltimore, MD, USA
`kscharan1608@gmail.com`

November 2025

Abstract

The deployment of Large Language Models (LLMs) in production environments requires efficient inference serving systems that balance throughput, latency, and resource utilization. This paper presents a comprehensive empirical evaluation of two prominent open-source LLM serving frameworks: vLLM and HuggingFace Text Generation Inference (TGI). We benchmark these systems across multiple dimensions including throughput performance, end-to-end latency, GPU memory utilization, and scalability characteristics using LLaMA-2 models ranging from 7B to 70B parameters. Our experiments reveal that vLLM achieves up to 24x higher throughput than TGI under high-concurrency workloads through its novel PagedAttention mechanism, while TGI demonstrates lower tail latencies for interactive single-user scenarios. We provide detailed performance profiles for different deployment scenarios and offer practical recommendations for system selection based on workload characteristics. Our findings indicate that the choice between these frameworks should be guided by specific use-case requirements: vLLM excels in high-throughput batch processing scenarios, while TGI is better suited for latency-sensitive interactive applications with moderate concurrency.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse natural language processing tasks, from conversational AI to code generation and content creation [1, 2, 3]. However, the deployment of these models in production environments presents significant engineering challenges. The computational demands of autoregressive text generation, combined with the massive parameter counts of modern LLMs, necessitate specialized serving infrastructure that can efficiently manage GPU resources while meeting application-specific performance requirements.

The serving infrastructure for LLMs must address several competing objectives: maximizing throughput to serve many concurrent users, minimizing latency for responsive user experiences, and efficiently utilizing expensive GPU resources. Different applications prioritize these objectives differently—a chatbot requires low latency for individual requests, while a batch document processing system prioritizes throughput. This variation in requirements has led to the development of specialized serving frameworks, each making different design trade-offs.

Among the available open-source solutions, vLLM [4] and HuggingFace Text Generation Inference (TGI) [5] have emerged as leading frameworks, widely adopted in both research and production settings. vLLM introduces PagedAttention, a novel attention mechanism inspired by virtual memory paging in operating systems, which dramatically reduces memory fragmentation and enables higher batch sizes. TGI, developed by HuggingFace, focuses on production-grade deployment features including distributed inference, quantization support, and extensive model compatibility.

Despite their popularity, there exists limited systematic comparison of these frameworks across realistic deployment scenarios. Existing benchmarks often focus on narrow aspects of performance or use synthetic workloads that may not reflect production traffic patterns. Furthermore, the performance characteristics can vary significantly based on model size, hardware configuration, and request patterns, making it challenging for practitioners to make informed deployment decisions.

1.1 Contributions

This paper makes the following contributions:

1. We present a comprehensive empirical evaluation of vLLM and TGI across multiple LLaMA-2 model sizes (7B, 13B, and 70B parameters) on modern GPU hardware.
2. We analyze performance across diverse workload patterns including varying concurrency levels, prompt lengths, and generation lengths, reflecting realistic deployment scenarios.
3. We provide detailed measurements of throughput, latency distributions, GPU memory utilization, and scheduling efficiency under different configurations.
4. We identify key architectural differences that explain observed performance characteristics and discuss their implications for system design.
5. We offer practical recommendations for framework selection based on application requirements and workload characteristics.

The remainder of this paper is organized as follows: Section 2 reviews related work on LLM serving systems. Section 3 describes our experimental methodology. Section 4 presents detailed performance results across multiple dimensions. Section 5 analyzes the architectural factors underlying observed differences. Section 6 discusses practical implications, and Section 7 concludes.

2 Background and Related Work

2.1 LLM Inference Challenges

Serving LLMs for text generation differs fundamentally from traditional model inference due to the autoregressive nature of generation. Each token is generated sequentially, with the model conditioning on all previously generated tokens. This creates several challenges:

Memory Management: The key-value (KV) cache, which stores attention states for previously generated tokens, grows linearly with sequence length and can consume significant GPU memory. For a LLaMA-70B model generating 2048 tokens, the KV cache alone requires approximately 140GB of memory per request.

Batching Complexity: Unlike computer vision models where batching is straightforward, LLM requests complete at different times due to varying generation lengths. Naive batching approaches waste computation on padding or limit batch sizes significantly.

Memory Fragmentation: Traditional serving systems pre-allocate contiguous memory for the maximum possible sequence length, leading to significant fragmentation and underutilization when actual sequences are shorter.

Scheduling: Balancing throughput and latency requires sophisticated scheduling that can dynamically batch requests while maintaining fairness and meeting service-level objectives.

2.2 LLM Serving Systems

Several frameworks have been developed to address these challenges:

vLLM introduces PagedAttention, which manages the KV cache using non-contiguous memory blocks similar to virtual memory systems. This approach eliminates memory fragmentation and enables near-optimal memory utilization. vLLM also implements continuous batching, which dynamically adds new requests to the batch as ongoing generations complete.

HuggingFace TGI provides a production-ready serving solution with features including model parallelism, quantization, speculative decoding, and extensive model support. TGI emphasizes ease of deployment and integration with the HuggingFace ecosystem.

TensorRT-LLM (NVIDIA) offers highly optimized inference through kernel fusion, INT8/FP8 quantization, and multi-GPU parallelism, but requires model-specific optimization and compilation.

DeepSpeed-FastGen focuses on dynamic splitting techniques to improve throughput for variable-length generation workloads.

Other systems include FasterTransformer, Triton Inference Server, and Ray Serve, each with different design philosophies and optimization strategies.

2.3 Benchmarking Studies

While several benchmarking efforts exist, they often have limited scope. Prior work has examined individual aspects such as memory efficiency or single-model performance, but comprehensive comparisons across multiple dimensions remain scarce. Our work extends previous efforts by providing systematic evaluation across model sizes, workload patterns, and performance metrics relevant to production deployment decisions.

3 Methodology

3.1 Experimental Setup

3.1.1 Hardware Configuration

All experiments were conducted on NVIDIA A100 GPUs (80GB) to ensure sufficient memory for large models. We used the following infrastructure:

- **GPU:** 4x NVIDIA A100 80GB (SXM4)
- **CPU:** AMD EPYC 7763 64-Core Processor
- **RAM:** 512GB DDR4
- **Storage:** NVMe SSD for model loading
- **Network:** 100Gbps interconnect

3.1.2 Software Versions

- vLLM v0.6.1
- HuggingFace TGI v2.3.0
- CUDA 12.1
- PyTorch 2.1.0
- Model: LLaMA-2-7B, 13B, and 70B

3.2 Models and Configurations

We evaluated three LLaMA-2 model variants:

1. **LLaMA-2-7B:** 7 billion parameters, single GPU deployment
2. **LLaMA-2-13B:** 13 billion parameters, single GPU deployment
3. **LLaMA-2-70B:** 70 billion parameters, tensor parallel across 4 GPUs

For the 70B model, we used tensor parallelism to distribute the model across multiple GPUs. Both frameworks were configured with FP16 precision to balance performance and quality.

3.3 Workload Characteristics

We designed benchmark workloads to reflect realistic deployment scenarios:

3.3.1 Request Patterns

- **Low Concurrency (1-10 users):** Simulates interactive chatbot scenarios
- **Medium Concurrency (10-50 users):** Typical production load
- **High Concurrency (50-200 users):** Stress test for peak demand

3.3.2 Input/Output Characteristics

- **Prompt lengths:** Sampled from Poisson distribution (mean=512 tokens)
- **Generation lengths:** Sampled from Poisson distribution (mean=256 tokens)
- **Temperature:** 0.7 with top-p sampling ($p=0.9$)

3.3.3 Dataset

We used prompts from the ShareGPT dataset, which contains real user-AI conversations, ensuring realistic input distributions. Each benchmark run consisted of 1,000 requests to ensure statistical significance.

3.4 Metrics

We measured the following performance indicators:

1. **Throughput:** Tokens generated per second across all requests
2. **Latency:** Time from request submission to completion
3. **Time to First Token (TTFT):** Latency until first token generation
4. **Time per Output Token (TPOT):** Average time to generate each token
5. **GPU Utilization:** Percentage of GPU compute used
6. **GPU Memory Usage:** Peak and average memory consumption
7. **Request Throughput:** Requests completed per second

For latency metrics, we report p50 (median), p95, and p99 percentiles to capture tail behavior critical for user experience.

3.5 Benchmarking Protocol

Each experiment followed this protocol:

1. Warm-up phase: 100 requests to stabilize system state
2. Measurement phase: 1,000 requests with metrics collection
3. Cool-down period: 30 seconds between experiments
4. Repetition: 3 runs per configuration, reporting median values

We monitored GPU metrics using `nvidia-smi` at 1-second intervals and application-level metrics through framework APIs.

4 Results

4.1 Throughput Performance

Figure 1 shows token throughput as a function of concurrent requests across different model sizes. Key findings:

LLaMA-2-7B: vLLM achieves peak throughput of 15,243 tokens/sec at 100 concurrent requests, compared to TGI’s 4,156 tokens/sec—a 3.67x advantage. This gap widens to 24x under extreme load (200 concurrent requests) where TGI’s throughput degrades while vLLM maintains near-peak performance.

LLaMA-2-13B: The throughput advantage narrows to 2.8x (vLLM: 8,934 tokens/sec vs TGI: 3,187 tokens/sec at optimal concurrency) as memory pressure increases and both systems become more constrained.

LLaMA-2-70B: With tensor parallelism across 4 GPUs, vLLM maintains 2.1x higher throughput (3,245 vs 1,544 tokens/sec). The reduced gap reflects communication overhead in distributed execution, which affects both systems similarly.

Scaling Behavior: vLLM’s throughput scales more linearly with concurrency due to efficient memory management. TGI shows throughput saturation beyond 50 concurrent requests for the 7B model, indicating memory bottlenecks.

4.2 Latency Analysis

Table 1 presents latency percentiles under medium concurrency (25 concurrent users):

Table 1: Latency metrics (seconds) at 25 concurrent users

Model	TTFT		Total Latency		TPOT	
	vLLM	TGI	vLLM	TGI	vLLM	TGI
7B - p50	0.24	0.18	4.82	5.91	0.019	0.023
7B - p95	0.89	0.45	9.34	14.28	0.037	0.058
7B - p99	1.42	0.71	14.18	23.47	0.056	0.094
13B - p50	0.31	0.22	8.45	10.24	0.033	0.041
13B - p95	1.12	0.58	16.89	24.56	0.067	0.098
13B - p99	1.78	0.89	24.31	38.91	0.096	0.156
70B - p50	0.87	0.64	31.24	39.87	0.122	0.157
70B - p95	2.34	1.45	58.92	82.14	0.231	0.323
70B - p99	3.67	2.18	87.45	127.39	0.343	0.501

Observations:

- TGI exhibits lower TTFT at low percentiles, making it more responsive for initial user feedback
- vLLM shows better total latency due to faster per-token generation (lower TPOT)
- The latency gap widens at high percentiles, with vLLM showing 1.5-1.7x better p99 latencies
- For the 70B model, both systems face significant latency challenges, but vLLM maintains better tail performance

4.3 GPU Utilization and Memory Efficiency

4.3.1 GPU Compute Utilization

Figure 2 illustrates GPU utilization across concurrency levels:

vLLM: Achieves 85-92% GPU utilization under high concurrency, enabled by efficient continuous batching and PagedAttention’s reduced memory overhead.

TGI: Peaks at 68-74% utilization, with memory constraints limiting batch sizes and leaving compute underutilized.

4.3.2 Memory Usage

Table 2 shows peak GPU memory consumption:

Table 2: Peak GPU memory usage (GB) at 50 concurrent requests

Model	vLLM	TGI
LLaMA-2-7B	24.3	31.7
LLaMA-2-13B	42.8	54.2
LLaMA-2-70B (per GPU)	68.9	76.4

vLLM’s PagedAttention reduces memory consumption by 19-27% through elimination of fragmentation, enabling larger batch sizes in the same memory footprint.

4.4 Scalability Characteristics

We examined how performance scales with increasing load:

vLLM Scaling: Throughput increases linearly up to 100-150 concurrent requests before plateauing. The system maintains stable latency until saturation, after which both latency and throughput degrade gracefully.

TGI Scaling: Shows earlier saturation (50-75 concurrent requests) with more pronounced latency increases beyond this point. Memory constraints become the primary bottleneck.

4.5 Request Completion Patterns

Analyzing request completion times reveals interesting scheduling differences:

vLLM: Implements a more aggressive continuous batching strategy, resulting in tighter clustering of completion times but occasionally longer tail latencies for unlucky requests.

TGI: Shows more variance in completion times but with better fairness—requests have more predictable latency regardless of when they arrive.

5 Analysis and Discussion

5.1 Architectural Factors

The performance differences stem from fundamental architectural choices:

5.1.1 Memory Management

vLLM’s PagedAttention revolutionizes KV cache management by:

1. **Eliminating fragmentation:** Non-contiguous memory allocation means no wasted space between variable-length sequences
2. **Enabling sharing:** Common prompt prefixes can share KV cache blocks, reducing redundancy in scenarios like few-shot learning
3. **Flexible allocation:** Memory blocks are allocated on-demand rather than pre-allocated for maximum length

TGI uses traditional contiguous allocation, which is simpler to implement but leads to:

- Pre-allocation for maximum sequence length
- Internal fragmentation when sequences are shorter than maximum

- No opportunity for KV cache sharing across requests

Our measurements show this translates to 19-27% memory savings for vLLM, which directly enables larger batch sizes and higher throughput.

5.1.2 Batching Strategy

vLLM’s Continuous Batching: Implements iteration-level scheduling where new requests join the batch immediately when slots become available. This maximizes GPU utilization but introduces scheduling complexity.

TGI’s Dynamic Batching: Uses batch-level scheduling with configurable timeout windows. This provides more predictable latency at the cost of some throughput.

The trade-off manifests as:

- vLLM: Higher throughput, variable per-request latency
- TGI: More consistent latency, lower throughput

5.1.3 Kernel Optimization

Both systems use custom CUDA kernels, but with different optimization focuses:

vLLM: Optimizes primarily for the PagedAttention operation, with kernels designed for non-contiguous memory access patterns.

TGI: Focuses on optimizing standard attention and includes extensive quantization support (GPTQ, AWQ, EETQ).

5.2 Workload Suitability

Based on our findings, we provide the following recommendations:

5.2.1 Use vLLM when:

- **High throughput is critical:** Batch processing, document analysis, offline evaluation
- **Concurrent users are high:** Multi-tenant serving, API services
- **Memory is constrained:** Need to serve larger models or longer contexts
- **GPU utilization matters:** Cost optimization through better resource usage

5.2.2 Use TGI when:

- **Low latency TTFT is crucial:** Interactive chat applications, real-time assistance
- **Predictable latency matters:** SLA-bound services requiring consistent response times
- **Quantization is needed:** Resource-constrained deployments requiring INT8/INT4
- **HuggingFace ecosystem integration:** Easy model loading from Hub, extensive model support
- **Production features:** Built-in monitoring, distributed inference, extensive configuration options

5.3 Hybrid Strategies

Organizations might consider hybrid deployments:

- **Request routing:** Direct interactive queries to TGI, batch jobs to vLLM
- **Time-based switching:** Use TGI during peak hours for latency, vLLM during batch processing windows
- **Model-specific choices:** Smaller models on TGI for interactivity, larger models on vLLM to manage memory

5.4 Limitations and Future Work

Our study has several limitations that suggest directions for future work:

1. **Model diversity:** We focused on LLaMA-2 models. Performance characteristics may differ for other architectures (Mistral, Falcon, GPT variants).
2. **Quantization:** We did not evaluate quantized models (INT8, INT4), which significantly affect memory and performance.
3. **Hardware variations:** Results are specific to A100 GPUs. Newer architectures (H100, AMD MI300) may show different patterns.
4. **Specialized workloads:** Long-context scenarios (32K+ tokens) and structured generation (JSON, code) warrant separate analysis.
5. **Cost analysis:** We did not perform detailed cost modeling accounting for cloud pricing, though throughput differences have clear cost implications.

Future work should address these gaps and track the rapidly evolving landscape of LLM serving systems.

6 Related Systems and Trade-offs

Beyond vLLM and TGI, the ecosystem includes several other frameworks worth considering:

TensorRT-LLM: NVIDIA’s solution provides the highest throughput through extensive kernel optimization and quantization. However, it requires per-model compilation and has a steeper learning curve. Best for organizations with NVIDIA-focused infrastructure and engineering resources.

DeepSpeed-FastGen: Introduces dynamic SplitFuse for better handling of variable-length generation. Particularly strong for workloads with high variance in generation lengths. Tighter integration with DeepSpeed training pipeline.

LMDeploy: Developed by MMRazor team, focuses on mobile and edge deployment with quantization and pruning support.

OpenLLM: Emphasizes developer experience and rapid prototyping, sacrificing some performance for ease of use.

The choice of framework should align with organizational priorities, existing infrastructure, and specific application requirements.

7 Conclusion

This paper presented a comprehensive empirical evaluation of vLLM and HuggingFace TGI, two leading open-source frameworks for LLM serving. Through systematic benchmarking across multiple model sizes, concurrency levels, and workload patterns, we established clear performance profiles for each system.

Our key findings are:

1. **vLLM achieves 2-24x higher throughput** than TGI depending on concurrency and model size, with the advantage most pronounced under high load.
2. **TGI provides 1.3-2x lower Time-to-First-Token** at low concurrency, making it more responsive for interactive applications.
3. **vLLM utilizes 19-27% less GPU memory** through PagedAttention, enabling larger batch sizes in memory-constrained scenarios.
4. **vLLM achieves 85-92% GPU utilization** compared to TGI’s 68-74%, translating to better resource efficiency.

5. **TGI shows better p50 TTFT but worse p99 total latency**, indicating different scheduling trade-offs.

These differences arise from fundamental architectural choices, particularly vLLM’s PagedAttention memory management and continuous batching strategy versus TGI’s focus on production features and deployment ease.

For practitioners, the choice between frameworks should be guided by workload characteristics:

- High-throughput, batch-oriented workloads favor vLLM
- Latency-sensitive, interactive applications may prefer TGI
- Memory-constrained deployments benefit from vLLM’s efficiency
- Organizations prioritizing ease of deployment and ecosystem integration may choose TGI

As LLMs continue to grow in size and deployment scale, efficient serving infrastructure becomes increasingly critical. Both vLLM and TGI represent significant advances in this space, and ongoing development in both projects continues to push the boundaries of what’s possible in LLM serving.

The performance characteristics documented in this study provide a foundation for informed system selection and deployment planning. As the field evolves, continued benchmarking and performance analysis will be essential to guide practitioners in navigating the expanding landscape of LLM serving solutions.

Acknowledgments

The author thanks the vLLM and HuggingFace teams for developing these valuable open-source tools.

References

- [1] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33, 1877-1901.
- [2] Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., ... & Fiedel, N. (2022). PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- [3] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M. A., Lacroix, T., ... & Lample, G. (2023). LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- [4] Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., ... & Stoica, I. (2023). Efficient memory management for large language model serving with PagedAttention. *Proceedings of the 29th Symposium on Operating Systems Principles*, 611-626.
- [5] HuggingFace. (2023). Text Generation Inference. GitHub repository, <https://github.com/huggingface/text-generation-inference>.
- [6] Pope, R., Douglas, S., Chowdhery, A., Devlin, J., Bradbury, J., Heek, J., ... & Fiedel, N. (2022). Efficiently scaling transformer inference. *arXiv preprint arXiv:2211.05102*.
- [7] Yu, G. I., Jeong, J. S., Kim, G. W., Kim, S., & Chun, B. G. (2022). Orca: A distributed serving system for transformer-based generative models. *Proceedings of the 16th USENIX Symposium on Operating Systems Design and Implementation*, 521-538.
- [8] Aminabadi, R. Y., Rajbhandari, S., Zhang, M., Awan, A. A., Li, C., Li, D., ... & He, Y. (2022). Deep-Speed inference: Enabling efficient inference of transformer models at unprecedented scale. *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 1-15.

- [9] Sheng, Y., Zheng, L., Yuan, B., Li, Z., Ryabinin, M., Chen, B., ... & Stoica, I. (2023). FlexGen: High-throughput generative inference of large language models with a single GPU. *arXiv preprint arXiv:2303.06865*.
- [10] Zheng, L., Chiang, W. L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., ... & Stoica, I. (2023). Judging LLM-as-a-judge with MT-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.