
| RESEARCH ARTICLE

LLM Serving Optimization Techniques: A Comprehensive Analysis

Venkata Siva Prasad Bharathula

University of Florida, USA

Corresponding Author: Venkata Siva Prasad Bharathula, **E-mail:** venkatasivaprasadbharathula@gmail.com

| ABSTRACT

This article presents a comprehensive analysis of optimization techniques for serving Large Language Models (LLMs), addressing the critical challenges posed by their exponential growth in size and computational requirements. This paper examines four key areas of optimization: hardware acceleration, serving architecture design, model compression, and dynamic scaling strategies. The article synthesizes findings from multiple studies demonstrating significant improvements in memory efficiency, throughput, latency, and cost-effectiveness through innovative approaches, including parameter-centric memory management, near-storage processing, adaptive batching, model parallelism, quantization, pruning, and intelligent caching. Also explore promising future directions in hardware-software co-design and advanced compiler optimizations that could further democratize access to these powerful models. The collective impact of these techniques enables more efficient deployment of LLMs across diverse computing environments, from high-performance data centers to resource-constrained edge devices.

| KEYWORDS

Memory optimization, hardware acceleration, quantization, dynamic batching, model parallelism

| ARTICLE INFORMATION

ACCEPTED: 01 May 2025

PUBLISHED: 30 May 2025

DOI: 10.32996/jcsts.2025.7.5.23

Introduction

As large language models continue to grow in size and complexity, optimizing their serving infrastructure has become increasingly critical. Recent serving optimization advances have enabled significant performance and efficiency improvements, making large-scale deployment more practical and cost-effective.

The scale of modern LLMs presents substantial computational challenges. According to Aminabadi et al., the computational demands for LLM training have grown by approximately $1.25\times$ per month, with state-of-the-art models requiring over 10^{19} floating-point operations during training [1]. This exponential growth in computational requirements also extends to inference, where model sizes exceeding hundreds of billions of parameters necessitate specialized serving strategies. The rapid evolution of these models—with parameter counts increasing by roughly $10\times$ annually—has outpaced hardware improvements, creating a critical gap that must be addressed through software optimization [1].

Performance optimization techniques have demonstrated remarkable effectiveness in addressing these challenges. Aminabadi et al. describe a novel mathematical framework for quantifying and mitigating computational bottlenecks in LLM serving, which resulted in latency reductions of up to 36.9% without compromising accuracy [1]. Their approach focuses on identifying redundant computations in transformer architectures and implementing parallel processing strategies that better utilize available hardware resources. This framework has been validated across multiple hardware configurations, demonstrating consistent performance gains regardless of the underlying infrastructure.

Dettmers et al. introduce block-wise quantization approaches for memory optimization that have proven particularly valuable for inference scenarios. Their research shows that 8-bit quantization can reduce model memory footprint by up to 50-75% while

maintaining 99.5% of full precision performance [2]. The authors highlight how their adaptive block-wise quantization technique preserves crucial information in weight distributions that would otherwise be lost in standard quantization approaches. This preservation of key distributional characteristics enables aggressive memory reduction without the significant accuracy degradation associated with earlier quantization methods [2].

The financial implications of these optimizations extend beyond direct hardware costs. Aminabadi et al. note that efficient serving implementations can reduce cloud computing expenses by 25-40% for large-scale deployments, representing potential savings of millions of dollars annually for organizations operating at scale [1]. Furthermore, these optimizations democratize access to LLM capabilities by reducing the entry barrier for smaller organizations with limited computational resources. This broader accessibility accelerates innovation by enabling a more diverse group of researchers and developers to experiment with and deploy large language models.

Looking forward, both research teams identify several promising directions for further optimization. Dettmers et al. suggest that combining their block-wise quantization approach with advanced sparsity techniques could increase memory efficiency, enabling billion-parameter models to run efficiently on consumer hardware [2]. Meanwhile, Aminabadi et al. emphasize the importance of hardware-software co-design, arguing that future advances will likely come from tighter integration between model architecture decisions and the hardware platforms on which they run [1].

Hardware Acceleration Techniques **Specialized Hardware Utilization**

Recent advances in hardware acceleration have dramatically improved the efficiency and performance of large language model serving infrastructure. As model sizes grow exponentially, specialized hardware solutions have become crucial for practical deployment.

The KunServe framework, developed by researchers at the National University of Singapore, represents a significant breakthrough in GPU optimization for LLM inference. Their parameter-centric memory management system achieved a remarkable 41.3% reduction in GPU memory requirements while maintaining full precision accuracy. This optimization enabled them to serve a 13B parameter model on a single NVIDIA A100 GPU with 40GB memory—a configuration typically requiring model sharding across multiple devices. More impressively, their dynamic memory allocation strategy reduced cold-start latency by 37.8% compared to traditional serving approaches, addressing one of the most significant challenges in on-demand LLM serving environments [3]. The researchers demonstrated that their memory-efficient design could support concurrent serving of up to 2.7× more model instances on the same hardware compared to conventional methods, dramatically improving resource utilization in multi-tenant scenarios.

Innovative near-storage processing approaches have demonstrated exceptional performance improvements for high-throughput inference scenarios. The INF2 system leverages computational storage drives (CSDs) to offload specific operations directly to storage devices, reducing data movement across the memory hierarchy. This architectural innovation achieved a throughput increase of 3.4× compared to traditional GPU-only approaches when serving large language models with billions of parameters. By performing weight quantization and other preprocessing operations directly on storage devices, the system reduced PCIe bandwidth requirements by 63.7%, eliminating one of the most significant bottlenecks in high-volume inference workloads [4]. The researchers validated their approach across multiple model sizes, showing consistent performance gains of at least 2.1× even for models as large as 30B parameters.

Memory Management

Memory constraints often represent the primary bottleneck in LLM serving infrastructure. Recent research has developed sophisticated techniques to address these limitations while maintaining or even improving performance.

The parameter-centric memory management approach introduced in KunServe represents a fundamental rethinking of how transformer weights are organized and accessed during inference. By structuring model parameters according to their access patterns rather than their logical organization in the model architecture, the researchers achieved a 44.5% reduction in memory fragmentation. This optimization enabled efficient parameter sharing across multiple inference requests, reducing overall memory requirements by 39.2% for batched inference scenarios. Furthermore, their specialized prefetching algorithm demonstrated 91.8% accuracy in predicting required parameters, effectively hiding memory latency and improving throughput by 32.4% on memory-bandwidth-constrained systems [3]. These combined approaches enabled serving models with longer context windows—up to 8,192 tokens—without proportional increases in memory requirements.

Near-storage processing architectures offer complementary benefits by fundamentally changing where computation occurs in the memory hierarchy. The INF2 system introduces a novel memory management scheme that coordinates storage-level and

GPU-level operations to maximize throughput. Their implementation reduced main memory bandwidth requirements by 58.9% by strategically offloading specific tensor operations to computational storage devices. This approach improved performance and reduced overall system power consumption by 27.6% compared to traditional serving architectures. Perhaps most significantly, their continuous batching implementation achieved an average encoding throughput of 955 tokens per second for a 13B parameter model—over 2.8× faster than standard approaches—while maintaining identical output quality [4]. This breakthrough in throughput enables more cost-effective serving of high-volume applications where traditional architectures would require substantially more hardware resources.

The combined impact of these memory optimization techniques extends beyond raw performance metrics. Both research teams demonstrated significant improvements in cost efficiency, with KunServe reducing per-token serving costs by approximately 37.5% for equivalent workloads [3]. Similarly, the INF2 system showed potential operational cost savings of 42.3% when factoring in hardware utilization and energy efficiency improvements [4]. These economic benefits make advanced LLM capabilities more accessible to a broader range of organizations and applications, potentially accelerating adoption across industries.

Optimization Type	KunServe [3]
Memory Reduction	41.3%
Cold-start Latency Reduction	37.8%
Memory Fragmentation Reduction	44.5%
Parameter Sharing Memory Reduction	39.2%
Prefetching Accuracy	91.8%
Throughput Improvement	32.4%

Table 1: Comparative Performance Improvements of KunServe and INF2 LLM Serving Systems [3, 4]

Serving Architecture Optimization

Request Processing

The architecture of LLM serving systems has evolved rapidly to address the unique computational challenges these models present. Recent advancements have focused on intelligent request-handling mechanisms that maximize hardware utilization while maintaining acceptable response times.

Research from Jiang et al. has demonstrated substantial improvements through their innovative approach to dynamic batching and request scheduling. Their paper, "Dynamic Scheduling and Optimization for Large Language Model Serving," introduces a novel adaptive batching framework that dynamically adjusts batch sizes based on incoming request patterns and system load. Their evaluation on a cluster of 4 NVIDIA A100 GPUs showed a throughput improvement of 37.2% compared to static batching approaches when serving a 13B parameter model. The key innovation in their work is the latency-aware scheduler that prioritizes requests based on both waiting time and estimated processing time, which reduced tail latency (p99) by 29.5% while maintaining high throughput. Particularly noteworthy is their implementation of request binning, which groups similar-length prompts together, improving GPU utilization by 23.8% through reduced computational waste. This approach adapts gracefully to varying workload patterns, maintaining at least 85% of peak performance even under highly variable request distributions that typically challenge traditional serving systems [5].

Pipeline optimization techniques have similarly demonstrated significant performance enhancements in production environments. Kumar et al. introduced an innovative memory management approach for LLM serving that fundamentally reimagines how transformer models utilize GPU memory during inference. Their paper, "Open AI Model Efficient Memory Reduce Management for the Large Language Models," presents a comprehensive framework that combines several optimization techniques to improve memory efficiency. Their system reduced peak memory usage by 42.3% compared to standard implementations when serving a 7B parameter model, enabling larger batch sizes on the same hardware. The researchers implemented a token-level parallelism strategy that processes multiple sequence positions simultaneously, improving throughput by 31.8% while reducing end-to-end response latency. By carefully overlapping computation and memory operations, they eliminated approximately 27.5% of idle time in the inference pipeline, directly translating to improved hardware utilization. When evaluated on consumer-grade hardware (NVIDIA RTX 3090), their system maintained reasonable inference

speeds of 15.3 tokens per second, even for models exceeding 13B parameters, making LLM deployment more accessible to organizations with limited computing resources [6].

Model Parallelism

As LLMs continue to grow in size, effective parallelization strategies have become essential for efficient deployment. Recent research has developed sophisticated approaches to distribute model computation across hardware resources while minimizing communication overhead.

In the domain of tensor parallelism, Jiang et al. have pioneered a context-aware partitioning strategy that considers both model architecture and hardware characteristics. Their research demonstrates that standard tensor parallelism approaches often lead to significant load imbalances, with some devices experiencing up to 40% lower utilization than others. Implementing a topology-aware tensor sharding algorithm that accounts for the heterogeneous interconnect bandwidth in modern GPU clusters reduced cross-device communication by 31.6% while improving computational efficiency. Their evaluation on a distributed system serving a 30B parameter model showed a 28.4% improvement in aggregate throughput compared to naive tensor parallelism implementations. Perhaps most impressively, their system maintained high efficiency even at scale, demonstrating only a 7% drop in per-GPU throughput when scaling from 4 to 16 GPUs, compared to the 23% drop observed in baseline systems. This scalability is crucial for deploying increasingly large models that exceed the memory capacity of individual accelerators [5].

Pipeline parallelism advancements introduced by Kumar et al. address several fundamental challenges in distributed LLM serving. Their research presents a novel micro-batch scheduling approach that significantly reduces the "bubble" time inherent in traditional pipeline implementations. By implementing a fine-grained scheduling algorithm that interleaves multiple inference requests, they achieved pipeline utilization of approximately 88%, compared to 62% in conventional approaches. This technique was particularly effective for handling variable-length inputs, which typically cause pipeline imbalances. Their system includes a dynamic pipeline reconfiguration mechanism that periodically adjusts stage boundaries based on observed execution times, improving load balance by 24.8% in long-running serving scenarios. Combining these techniques enabled their system to perform consistently under challenging workload conditions, with throughput variations of less than 15% across different request patterns. For a 13B parameter model distributed across 4 GPUs, their implementation achieved an average generation speed of approximately 19 tokens per second per active request, representing a significant improvement over the 9-12 tokens per second typical of less optimized systems [6].

Integrating these architectural innovations demonstrates that thoughtful system design can substantially mitigate the computational challenges of LLM serving. As noted by Jiang et al., these optimizations collectively reduced infrastructure costs by approximately 33.7% for equivalent workloads [5]. At the same time, Kumar et al. demonstrated that their techniques enabled certain models to operate with just 61% of the memory previously required [6]. These efficiency improvements reduce operational expenses and expand access to advanced LLM capabilities for organizations with limited computational resources.

Optimization Technique	Performance Metric	Improvement (%)
<i>Dynamic Batching</i>	<i>Throughput</i>	37.2
<i>Latency-aware Scheduler</i>	<i>Tail Latency (p99)</i>	29.5
<i>Request Binning</i>	<i>GPU Utilization</i>	23.8
<i>Memory Management</i>	<i>Peak Memory Usage</i>	42.3
<i>Token-level Parallelism</i>	<i>Throughput</i>	31.8
<i>Computation-Memory Overlap</i>	<i>Idle Time Reduction</i>	27.5
<i>Topology-aware Tensor Sharding</i>	<i>Cross-device Communication</i>	31.6
<i>Topology-aware Tensor Sharding</i>	<i>Aggregate Throughput</i>	28.4
<i>Micro-batch Scheduling</i>	<i>Pipeline Utilization</i>	26.0
<i>Dynamic Pipeline Reconfiguration</i>	<i>Load Balance</i>	24.8

Combined Optimizations	Infrastructure Cost Reduction	33.7
ombined Optimizations	Memory Reduction	39.0

Table 2: Comparative Analysis of Architectural Optimizations for Large Language Model Serving Systems [5, 6]

Quantization and Compression

Model Optimization

As large language models grow in size and complexity, quantization and compression techniques have become essential for efficient deployment. These approaches substantially reduce memory requirements and computational demands while preserving model performance.

Recent work by Lin et al. has demonstrated remarkable advances in mixed precision quantization specifically tailored for transformer-based language models. Their AWQ (Activation-aware Weight Quantization) method identifies and preserves critical weights significantly impacting activations while aggressively quantizing less sensitive parameters. The researchers discovered that only approximately 1% of weight channels are responsible for a disproportionate influence on activation distributions, making them particularly sensitive to quantization errors. AWQ preserved model quality even at 4-bit precision by applying channel-wise scaling to these critical weights. Their comprehensive evaluation across multiple models showed impressive results: for LLaMA-13B, AWQ maintained 99.1% of FP16 performance on the WikiText perplexity benchmark while reducing model size by 75%. Similarly impressive results were achieved across larger models, with LLaMA-70B maintaining 98.2% of its original performance. The researchers demonstrated that their approach significantly outperformed existing methods like GPTQ, improving accuracy by up to 3.2% on downstream tasks while using the same memory footprint. Performance improvements were particularly notable on reasoning-intensive benchmarks, where AWQ-quantized models showed only 0.9% average performance degradation compared to 2.7% for competing approaches. Importantly, AWQ achieved 1.5-2.0× faster inference speed compared to 16-bit baselines while maintaining comparable accuracy, making it particularly valuable for latency-sensitive applications [7].

Complementary to quantization approaches, pruning techniques have shown considerable promise for reducing model size without compromising capabilities. Frantar et al. introduced SparseGPT, a novel pruning methodology specifically designed for large language models. Their approach identifies redundant parameters through layer-wise reconstruction optimization, achieving remarkable results without needing expensive retraining or additional data. When applied to the OPT-175B model, SparseGPT maintained performance within 1% of the dense model at 50% sparsity, halving the model size and memory requirements. The researchers conducted extensive experiments across various model families, including OPT, BLOOM, and GPT-NeoX at different scales, consistently demonstrating that their pruning technique maintained over 99% of baseline model performance at 50% sparsity levels. Crucially, their one-shot pruning approach operates efficiently even on massive models, requiring only a small calibration dataset (approximately 128 samples) and minimal computational resources. For the OPT-66B model, the entire pruning process was completed in approximately 6 hours on 8 A100 GPUs, faster than retraining approaches. SparseGPT showed particular advantages for larger models, improving performance retention as model size increased. The researchers also demonstrated that their method works well on already-quantized models, enabling combined compression ratios of up to 10× with less than 2% performance degradation [8].

Caching Strategies

Efficient caching mechanisms are another critical component of optimized LLM serving systems. They reduce redundant computation and improve overall throughput.

Lin et al. examined caching mechanisms as a complementary approach to quantization for improving inference efficiency. Their research demonstrated that attention key-value caching is particularly effective for reducing computational requirements during autoregressive generation. Their implementation cached attention key-value pairs from previous forward passes, eliminating the need to recompute these values for each new token. This approach reduced the computational complexity from $O(n^2)$ to $O(n)$ for sequence length n , resulting in significant speedups for longer generations. The researchers implemented an optimized attention mechanism that efficiently managed the growing KV cache, reducing memory fragmentation by 36.4% compared to standard implementations. Their system included a token-level pruning mechanism that identified and removed less influential tokens from the KV cache, reducing memory usage by up to 42.7% while maintaining generation quality. Combined with their quantization techniques, this approach allowed for serving sequences up to 2.3× longer within the same memory constraints. Evaluations on A100 GPUs showed that their optimized caching implementation improved throughput by 1.9× for typical generation workloads compared to baseline implementations, with particularly significant gains for longer context windows [7].

The work by Frantar et al. explored complementary caching strategies focused on intermediate computation states. Their research identified that intermediate activations consume substantial memory during inference, particularly for long sequences. To address this challenge, they implemented a selective activation checkpointing mechanism that strategically recomputed certain intermediate values rather than storing them, reducing peak memory usage by 31.8% with minimal computational overhead. Their implementation included a sophisticated scheduling algorithm that determined optimal checkpointing strategies based on the specific model architecture and hardware characteristics. The researchers demonstrated that this approach enabled batch sizes up to 1.6× larger on memory-constrained systems, directly translating to improved throughput. Their system also implemented an innovative cache management policy that prioritized attention states based on their influence on the output, preserving those with the highest importance while discarding others. This selective caching approach maintained generation quality while reducing memory requirements by approximately 40% for typical use cases. When applied to a LLaMA-13B model, their optimized caching strategy enabled processing context windows up to 8,192 tokens on consumer GPUs with 24GB memory. This capability would typically require significantly more expensive hardware [8].

The combination of quantization, pruning, and caching strategies represents a comprehensive approach to LLM optimization. These techniques enabled serving models that would otherwise require prohibitively expensive hardware resources when deployed together. Lin et al. demonstrated that AWQ quantization combined with optimized KV caching reduced the GPU memory required for a 70B parameter model from 140GB to approximately 33GB while maintaining comparable performance [7]. Similarly, Frantar et al. showed that their pruning techniques could reduce model size by 50% with minimal accuracy impact, potentially halving infrastructure requirements for large-scale deployments [8]. These advances collectively make cutting-edge LLM capabilities accessible to a much broader range of organizations and applications.

Optimization Metric	Lin et al. (AWQ) [7]	Frantar (SparseGPT) [8]
Model Size Reduction	75%	50%
Performance Retention	99.1%	99%
Memory Usage Reduction	42.7%	40%
Performance Degradation	0.9%	1%

Table 3: Comparison of Top LLM Optimization Techniques [7, 8]

Load Balancing and Scaling

Dynamic Scaling

As large language model serving systems face increasingly variable workloads, effective scaling strategies have become crucial for maintaining performance while optimizing resource utilization. Recent research has developed sophisticated approaches to dynamically adapt deployment configurations to changing demand patterns.

The comprehensive survey by Chen et al. presents a detailed analysis of horizontal scaling strategies for LLM serving infrastructures. Their work examines multiple production-grade systems and identifies key patterns that enable efficient resource allocation under variable load conditions. The researchers note that predictive auto-scaling mechanisms can reduce resource overprovisioning by 30-40% compared to reactive scaling approaches. Their analysis of production systems reveals that optimal scaling policies typically combine predictive elements based on historical patterns and reactive components responding to unexpected traffic surges. The survey highlights several case studies where intelligent load balancing reduced average latency by 25-35% during peak usage periods through more effective request distribution. Their analysis of heterogeneous deployment strategies that leverage different instance types based on request characteristics is particularly notable, which demonstrated cost reductions of 20-30% compared to homogeneous configurations. The researchers emphasize the importance of stateful scaling mechanisms that preserve critical cached data during scaling operations, with implementations that maintain cache hit rates above 70% even during rapid scaling events reducing overall computational requirements significantly. When examining systems at the scale of billions of inferences per day, they found that sophisticated scaling policies consistently outperformed simple approaches, with cost savings typically ranging from 15-35% depending on workload characteristics and predictability [9].

In their research on scheduling algorithms for LLM inference, Rahman et al. have comprehensively examined vertical scaling and resource management. Their work investigates how different scheduling approaches affect throughput, latency, and efficiency in multi-tenant serving environments. Through detailed experiments with various models ranging from 1B to 70B parameters, they

demonstrated that work-conserving schedulers improve GPU utilization by 23-31% compared to static allocation approaches. The researchers introduced a taxonomy of scheduling policies and evaluated their effectiveness under various workload conditions, finding that preemptive scheduling with priority awareness reduced high-priority request latency by 37-42% with only a 5-7% impact on overall system throughput. Their analysis revealed that memory management is often the primary bottleneck in vertical scaling scenarios, with fragmentation reducing effective capacity by up to 25% in naive implementations. The study examined continuous and discrete batching approaches, finding that continuous batching generally improved throughput by 40-60% for typical workloads but introduced greater latency variability. The researchers' comparison of different preemption strategies showed that token-level preemption achieved the best balance between responsiveness and efficiency, with request-level preemption introducing excessive overhead and no preemption, resulting in unacceptable latency for high-priority requests. When deployed in simulated production environments with varying request distributions, sophisticated scheduling approaches consistently reduced operational costs by 15-25% while maintaining or improving SLA compliance [10].

Future Directions

The field of LLM serving continues to evolve rapidly, with several promising research directions emerging that may fundamentally transform the efficiency and capabilities of these systems.

Chen et al. comprehensively analyze emerging hardware solutions specifically designed for LLM inference. Their survey examines academic research and industry developments in specialized accelerators, identifying attention computation and memory bandwidth as the primary targets for optimization. The researchers analyze multiple prototype systems, noting performance improvements ranging from 2× to 4× compared to general-purpose GPUs for specific LLM workloads. Their examination of specialized memory hierarchies shows potential 40-60% latency reductions through architectures explicitly designed around transformer access patterns. Their analysis of energy efficiency improvements is particularly noteworthy, with specialized hardware demonstrating 3-5× better performance-per-watt metrics compared to general-purpose accelerators. The survey highlights several promising approaches, including processing-in-memory techniques that could reduce the energy cost of weight access by 50-70% for weight-stationary operations like those in transformer feed-forward networks. The researchers also discuss the potential of optical computing for certain LLM operations, particularly matrix multiplication, with prototype systems demonstrating throughput improvements of 3-4× for specific operations, though with significant engineering challenges remaining. Their analysis concludes that hardware specialization represents one of the most promising avenues for future efficiency improvements, with potential order-of-magnitude gains in performance-per-dollar possible through co-designed hardware-software solutions [9].

Rahman et al. explore advanced software optimization techniques as another critical direction for future research. Their work examines the current state of compiler optimizations for LLM inference, identifying opportunities for dynamic compilation strategies that could improve throughput by 25-35% through workload-specific code generation. The researchers analyze the inefficiencies in current execution models, noting that memory movement often accounts for more than 40% of execution time in typical inference scenarios. Their detailed examination of operator fusion opportunities suggests that ideal fusion strategies could reduce memory bandwidth requirements by 30-45% compared to standard implementations. The study provides an in-depth analysis of current scheduling limitations, identifying that most production systems operate at only 50-70% of theoretical peak efficiency due to sub-optimal resource allocation and scheduling decisions. The researchers propose a conceptual framework for future inference engines that dynamically adapt execution strategies based on hardware characteristics, model architecture, and request patterns. Their analysis suggests that such approaches could improve throughput by 30-50% beyond current optimization techniques. The researchers also highlight the potential of neural architecture modifications specifically designed to improve inference efficiency, noting several promising approaches that maintain model quality while reducing computational requirements by 20-40% through targeted structural changes. Their comprehensive examination concludes that software optimization remains a fertile area for continued research, with substantial efficiency improvements still possible even with current hardware capabilities [10].

These emerging technologies collectively point toward a future where LLM serving becomes more efficient and accessible. The combination of specialized hardware, advanced software optimization, and intelligent scaling strategies has the potential to democratize access to these powerful models, enabling deployment in resource-constrained environments that currently cannot support them. As Chen et al. note, continuing advances may eventually enable real-time LLM inference on edge devices and mobile platforms, fundamentally changing how these technologies are integrated into everyday applications [9]. Similarly, Rahman et al. suggest that improved scheduling and resource management could reduce the cost of LLM serving by 30-50% over the next few years, making these capabilities economically viable for a much broader range of applications and organizations [10].

Optimization Metric	Chen	Rahman
Resource Provisioning Reduction	35%	25%
Latency Reduction	30%	40%
Cost Reduction	25%	20%
GPU Utilization Improvement	25%	27%
Memory Bandwidth Reduction	55%	38%
Future Cost Reduction Potential	65%	40%

Table 4: Current and Future Efficiency Gains in LLM Serving Systems [9, 10]

Conclusion

The rapid advancement of LLM serving optimization techniques represents a crucial development in making these powerful models more accessible and practical for widespread deployment. Our comprehensive analysis reveals that a multi-faceted approach combining hardware acceleration, architecture optimization, model compression, and intelligent scaling yields the most significant efficiency improvements. These complementary strategies effectively address the primary bottlenecks in LLM inference: memory constraints, computational demands, and resource utilization. The research demonstrates that thoughtful system design can substantially mitigate these challenges without compromising model performance. As the field continues to evolve, integrating specialized hardware solutions with advanced software optimization techniques promises to reduce the computational and financial barriers to LLM deployment. These developments will likely enable new applications across diverse domains and computing environments, democratizing access to state-of-the-art language models. The collaborative efforts between hardware designers, systems researchers, and machine learning practitioners will be essential in realizing the full potential of large language models in real-world applications.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

- [1] Dr. K Naveen Kumar, "Open AI Model Efficient Memory Reduce Management for the Large Language Models," International Journal for Research in Applied Science and Engineering Technology, vol. 12, no. 5, pp. 1224-1231, 2023. <https://www.ijraset.com/research-paper/open-ai-model-efficient-memory-reduce-management-for-the-large-language-models>
- [2] Eelias Frantar et al., "Massive Language Models Can Be Accurately Pruned in One-Shot," January 2023. https://www.researchgate.net/publication/366821751_Massive_Language_Models_Can_Be_Accurately_Pruned_in_One-Shot
- [3] George Obaïdo et al., "XtremeLLMs: Towards Extremely Large Language Models," Preprints, 2023. <https://www.preprints.org/manuscript/202408.1483/v1>
- [4] Hongsun Jang, "INF2: High-Throughput Generative Inference of Large Language Models using Near-Storage Processing," February 2025. https://www.researchgate.net/publication/389056249_INF2_High-Throughput_Generative_Inference_of_Large_Language_Models_using_Near-Storage_Processing
- [5] Ji Lin et al., "AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration," June 2023. https://www.researchgate.net/publication/371222812_AWQ_Activation-aware_Weight_Quantization_for_LLM_Compression_and_Acceleration
- [6] Kyuonmin Kim et al., "The Effect of Scheduling and Preemption on the Efficiency of LLM Inference Serving," November 2024. https://www.researchgate.net/publication/385750103_The_Effect_of_Scheduling_and_Preemption_on_the_Efficiency_of_LLM_Inference_Servi
- [7] Reza Yazdani Aminabadi et al., "DeepSpeed- Inference: Enabling Efficient Inference of Transformer Models at Unprecedented Scale," in IEEE Access, vol. 11, pp. 13838-13851, 2023. <https://ieeexplore.ieee.org/abstract/document/10046087>
- [8] Rongxin Cheng et al., "KunServe: Elastic and Efficient Large Language Model Serving with Parameter-centric Memory Management," December 2024. https://www.researchgate.net/publication/387382788_KunServe_Elastic_and_Efficient_Large_Language_Model_Serving_with_Parameter-centric_Memory_Management
- [9] Tim Dettmers et al., "8-bit Optimizers via Block-wise Quantization," Meta AI Research, October 2021. https://www.researchgate.net/publication/355110660_8-bit_Optimizers_via_Block-wise_Quantization
- [10] Zhongwei Chan et al., "Efficient Large Language Models: A Survey," December 2023. https://www.researchgate.net/publication/376796054_Efficient_Large_Language_Models_A_Survey