



Taming the Titans: A Survey of Efficient LLM Inference Serving

Anonymous ACL submission

Abstract

Large Language Models (LLMs) for Generative AI have achieved remarkable progress, evolving into sophisticated and versatile tools widely adopted across various domains and applications. However, the substantial memory overhead caused by their vast number of parameters, combined with the high computational demands of the attention mechanism, poses significant challenges in achieving low latency and high throughput for LLM inference services. This paper presents a comprehensive survey of recent advances in LLM inference optimization, categorized into: (1) fundamental techniques (model placement, request scheduling, decoding prediction, storage management, disaggregation, and multiplexing); (2) specialized architectures (MoE, LoRA, and speculative decoding); and (3) scenario-specific optimizations (long-context problem, RAG, Augmented LLMs, test-time reasoning, and multi-modal integration). Finally, we outline potential research directions to further advance the field of LLM inference serving.

1 Introduction

With the rapid evolution of open-source Large Language Models (LLMs), weekly updates to model architectures and capabilities have become the norm in recent years. The surging demand for these models is evident from Huggingface download statistics, which range from hundreds of thousands for models like Mistral-Small-24B-Instruct-2501 (Mistral, 2025), phi-4 (Abdin et al., 2024), and Llama-3.3-70B-Instruct (Grattafiori et al., 2024) to millions for DeepSeek-V3 (DeepSeek-AI et al., 2024) and DeepSeek-R1 (DeepSeek-AI et al., 2025) over recent months. However, when deploying these models, their large-scale parameters and attention mechanisms impose substantial demands on memory and computational resources, presenting significant obstacles to achieving the desired low latency and high throughput in processing requests. These

challenges have spurred extensive research across multiple domains of inference serving optimization to meet Service Level Objectives (SLOs).

This paper presents a systematic survey of LLM inference serving methods, organized hierarchically from fundamental techniques to specialized architectures and specific scenarios, as illustrated in Figure 1.

Fundamental Techniques optimization (§3) begins with model placement (§3.1), essential for distributing parameters across devices when single-GPU memory is insufficient. Subsequent request scheduling (§3.2) prioritizes batched processing through decoding length prediction (§3.3), where shorter requests are prioritized to reduce overall latency. Dynamic batch management then governs request insertion/eviction during iterative processing. While KV cache (§3.4) mitigates redundant computation, challenges persist in storage efficiency, reuse strategies, and compression. Due to the distinction between the prefill and decoding phases, the disaggregated architecture (§3.5) was introduced, facilitating the optimization of each phase. And how multiple LLMs efficiently share the same computing resources through multiplexing (§3.6).

Special Architecture (§4) including Mixture of Experts (MoE) (§4.1), which involves challenges such as expert placement, load balancing, and All-to-All communication. It also includes Low-Rank Adaptation (LoRA) (§4.2) and Speculative Decoding (§4.3), all of which require adaptability to address evolving demands.

Specific Scenario (§5) include advanced tasks such as Long Context processing (§5.1), Retrieval-Augmented Generation (RAG) (§5.2), Augmented LLMs (§5.2), Test-Time Reasoning (§5.4), and Multimodal (§5.5). These scenarios demand high adaptability to effectively address evolving requirements.

While previous work (Miao et al., 2023; Yuan et al., 2024; Zhou et al., 2024; Li et al., 2024a)

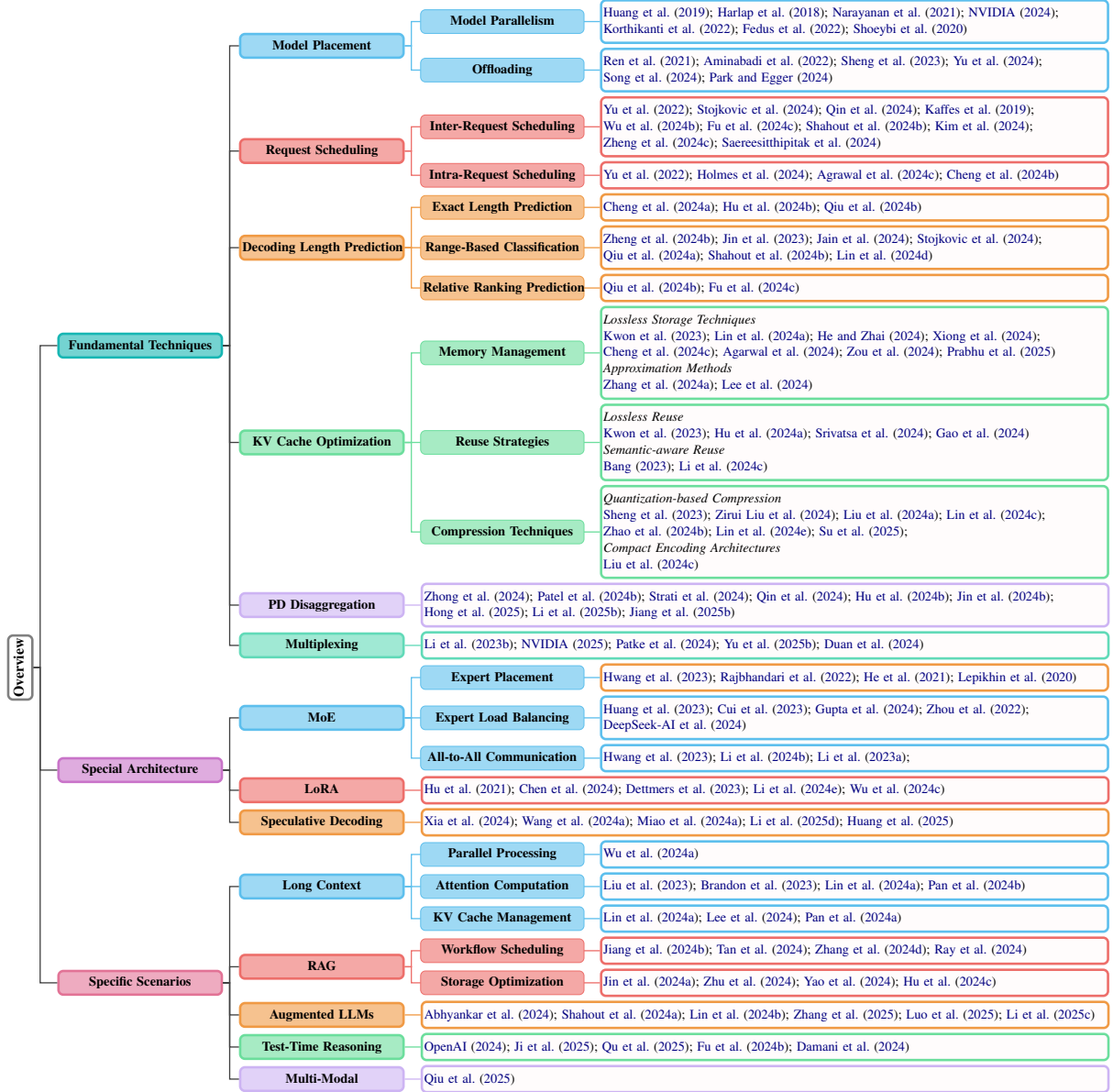


Figure 1: Taxonomy of approaches for LLM inference serving.

has built a strong foundation, the rapid development of the field calls for an urgent summary of emerging methods and scenarios to support deeper research. Our work integrates a wide range of both classical and cutting-edge methods, and adopts a forward-looking perspective and highlights several promising directions for future research. For completeness, we also include cloud-level topics (§A), miscellaneous areas (§B), and some inference framework (§C) in the Appendix.

2 Background

This section provides an overview of LLM fundamentals, aimed at enhancing the understanding of inference serving, along with the relevant evaluation metrics.

2.1 Transformer-based LLM

The LLM is primarily constructed on the foundation of the vanilla Transformer architecture, with a particular emphasis on its decoding component. The architecture is composed of multiple layers, primarily consisting of two key components: Multi-Head Self-Attention (MHA) and Feedforward Network (FFN), complemented by the LayerNorm operation.

The input representation $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ of the model is initially processed by tokenizing the user input and incorporating positional information. Subsequently, it is transformed through three learnable weight matrices, denoted as \mathbf{W}^Q , \mathbf{W}^K , and \mathbf{W}^V , to obtain the corresponding query (\mathbf{Q}), key (\mathbf{K}), and value (\mathbf{V}) vectors which are utilized as

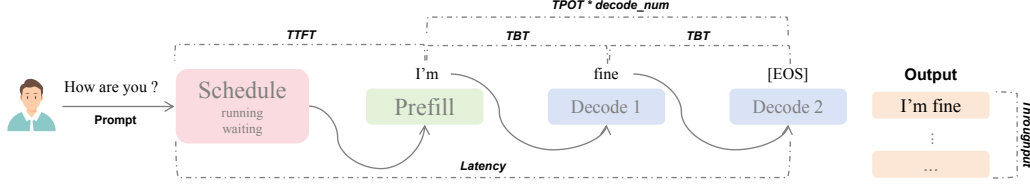


Figure 2: Illustration of some common evaluation metrics.

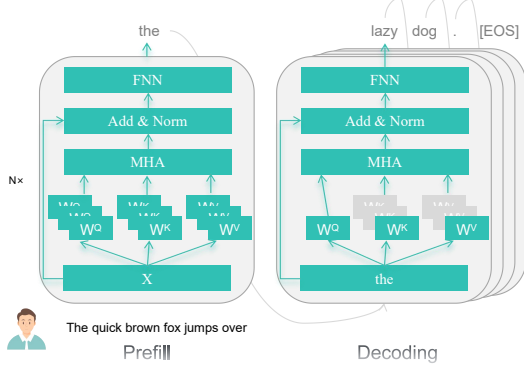


Figure 3: Illustration of the LLM Inference process.

inputs for the subsequent MHA:

$$\text{MHA}(\mathbf{Q}, \mathbf{W}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V} \quad (1)$$

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^Q; \mathbf{K} = \mathbf{X}\mathbf{W}^K; \mathbf{V} = \mathbf{X}\mathbf{W}^V$$

where d_k denotes the dimensionality of each attention head. It is evident that this constitutes the most time-consuming component, with a time complexity of $\mathcal{O}(n^2)$. The model processes m heads separately and concatenates the results:

$$\mathbf{O} = \text{Concat}(\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_m)\mathbf{W}^O \quad (2)$$

$$\mathbf{H}_i = \text{MHA}(\mathbf{Q}_i, \mathbf{W}_i, \mathbf{V}_i)$$

The FFN applies two linear transformations to its input, which is first processed by LayerNorm and residual connection:

$$\text{FFN}(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2 \quad (3)$$

2.2 Inference

LLM inference involves two phases: prefill and decoding, as illustrated in Figure 3. In prefill, the model processes the entire input in a compute-bound forward pass to produce the first token, while caching \mathbf{K} and \mathbf{V} (**KV cache**) to avoid recomputation. During decoding, tokens are generated sequentially using KV cache (gray blocks), which reduces the time complexity to $\mathcal{O}(n)$ at the cost of increased memory overhead. For each new token \mathbf{X}_{new} , the corresponding \mathbf{Q}_{new} , \mathbf{K}_{new} , and \mathbf{V}_{new} are computed, ensuring efficient generation and terminate at the [EOS] token.

2.3 Evaluation

These are the conventional metrics (Agarwal et al., 2023; Zhong et al., 2024; Qin et al., 2024; Yu et al., 2022) listed in Figure 2. In addition, Goodput (Zhong et al., 2024), or “effective throughput”, measures the maximum request rate that meets SLOs. Etalon (Agrawal et al., 2024a) is used to evaluate fluency to maintain smooth output during real-time interactions and its maximum output rate while preserving a certain level of fluency.

3 Fundamental Techniques

3.1 Model Placement

Due to the large number of parameters in LLMs, which exceed a single GPU’s capacity, distributing them across multiple GPUs or offloading them to CPUs has become a common practice.

Model Parallelism. This part focuses on several parallelism strategies. *Pipeline parallelism* (Huang et al., 2019; Harlap et al., 2018; Narayanan et al., 2021) distributes distinct model layers across multiple devices, enabling concurrent processing of sequential data to accelerate training/inference. *Tensor parallelism* (Shoeybi et al., 2020) splits individual operations or layers into smaller sub-tensors computed in parallel across devices, enhancing computational efficiency and enabling larger model dimensions. *Sequential parallelism* (Korthikanti et al., 2022) partitions LayerNorm and Dropout activations along the sequence dimension for long-context tasks. *Context parallelism* (NVIDIA, 2024) extends this by splitting all layers along the sequence dimension. *Expert parallelism* (Fedus et al., 2022) allocates sparse MoE components across GPUs, optimizing memory usage for sparse LLMs. More details can be seen in §4.1.

Offloading. When resources are limited, balancing GPU and CPU usage is essential. Some techniques (Ren et al., 2021; Aminabadi et al., 2022; Sheng et al., 2023) store most of a model’s weights in system memory or storage, loading only the necessary portions into GPU memory on demand.

PowerInfer (Song et al., 2024) computes hot neurons on the GPU and cold ones on the CPU, reducing both GPU memory usage and data transfer overhead. TwinPilots (Yu et al., 2024) integrates GPU and CPU with their hierarchical memories in an asymmetric multiprocessing framework. Park and Egger (2024) propose dynamic, fine-tuned workload allocation for efficient resource use.

3.2 Request Scheduling

For each instance, request scheduling directly impacts latency optimization. Here, we review relevant algorithms from both inter-request and intra-request scheduling perspectives.

Inter-Request Scheduling This section explores request batch prioritization under high load, focusing on execution order. Most LLM systems (Yu et al., 2022; Stojkovic et al., 2024; Qin et al., 2024) use First-Come-First-Served (FCFS), which can cause head-of-line blocking (Kaffes et al., 2019), such as a long request delaying a short one. Prioritizing shorter requests can reduce latency and better meet SLOs.

Advances in decoding length prediction (§3.3) have enabled smarter scheduling strategies. Fast-Serve (Wu et al., 2024b) uses a Skip-Join MLFQ to prioritize urgent or long-waiting requests while preempting long ones to speed up shorter tasks. Fu et al. (2024c) approximate Shortest Job First (SJF) by using predicted decoding times. Shahout et al. (2024b) improve Shortest Remaining Time First (SRTF) with dynamic length prediction and a preemption ratio, though frequent model calls introduce overhead. Prophet (Saareesitthipitak et al., 2024) applies SJF in prefill and Skip-Join MLFQ in decoding. Kim et al. (2024) shows that cost-guided preemption lowers GPU usage. BatchLLM (Zheng et al., 2024c), instead, focuses on maximizing global sharing.

Intra-Request Scheduling This section covers scheduling within concurrent request batches to enhance parallel decoding by handling variations in arrival, completion, and output length. Orca (Yu et al., 2022) offers iteration-level scheduling for dynamic request management per iteration, improving flexibility over inter-request methods. Dynamic SplitFuse (Holmes et al., 2024) and chunked-prefills (Agrawal et al., 2024c) break the prefill stage into smaller parts merged with decoding to cut delays from long prompts. SCLS (Cheng et al.,

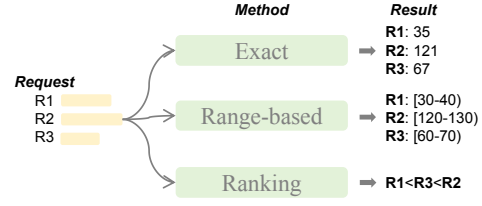


Figure 4: Overview of Length Prediction Methods.

2024b) divides generation into fixed-length slices, controlling service time and memory use precisely.

3.3 Decoding Length Prediction

Uncertain generation length complicates scheduling, but recent work (Figure 4) classifies prediction methods into three categories.

Exact Length Prediction predicts token counts using methods like BERT embeddings with random forest regression (Cheng et al., 2024a), small OPT models (Hu et al., 2024b), or simpler regression under constraints (Qiu et al., 2024b).

Range-Based Classification classifies requests into length bins. Some approaches (Zheng et al., 2024b; Jin et al., 2023; Jain et al., 2024; Qiu et al., 2024a; Stojkovic et al., 2024) predict length ranges directly from prompts. Notably, several works (Shahout et al., 2024b) leverage classifiers on token embeddings in real-time.

Relative Ranking Prediction predicts relative relationships between requests. Fu et al. (2024c) predicts relative relationships within the same batch, enhancing robustness and reducing overfitting.

Relative Ranking Prediction simplifies batch processing by only ordering requests, but requires re-ranking for carry-over requests, unlike other methods that avoid this overhead.

Other alternative approaches include SkipPredict (Shahout and Mitzenmacher, 2024), which uses quick classification for short/long tasks before detailed prediction, and BatchLLM (Zheng et al., 2024c) that predicts lengths through prompt analysis and statistical patterns.

3.4 KV Cache Optimization

KV cache cuts inference time from quadratic to linear but poses challenges in memory management, reuse, and compression. Optimizations for specialized storage are in §5.1 and 5.2.

Memory Management. *Lossless Storage Techniques* Kwon et al. (2023) propose PagedAttention and vLLM, using OS-style paging to nearly

Dimension	Lossless	Semantic-Aware
Matching Requests	Exact	Semantic similarity
Consistency	Fixed patterns, repeats	Diverse, open-ended
Overhead	✓✓✓	✓✓

Table 1: Comparison of KV cache reuse strategies.

eliminate memory fragmentation. FastDecode (He and Zhai, 2024) offloads cache to CPU memory through distributed processing, while LayerKV (Xiong et al., 2024) uses hierarchical allocation and offloading with layer-wise. KunServe (Cheng et al., 2024c) frees space for cache by removing model parameters, compensating via a pipeline mechanism from other instances. InstCache (Zou et al., 2024) enhances responsiveness through LLM-driven instruction prediction. PagedAttention’s non-contiguous layout adds complexity and overhead, while vAttention (Prabhu et al., 2025) reduces fragmentation and preserves virtual memory contiguity via CUDA-based memory decoupling.

Approximation Methods PQCache (Zhang et al., 2024a) uses Product Quantization to compress embeddings and cut computation. InfiniGen (Lee et al., 2024) improves performance through dynamic KV cache prefetching and reduced data transfer.

Reuse Strategies. *Lossless Reuse* PagedAttention (Kwon et al., 2023) enables multi-request cache sharing through page-level management. Radix tree-based systems (Hu et al., 2024a; Srivatsa et al., 2024) implement global prefix sharing with dynamic node deletion. CachedAttention (Gao et al., 2024) minimizes redundant computation in dialogues through cross-turn cache reuse.

Semantic-aware Reuse GPTCache (Bang, 2023) uses semantic similarity to cache and reuse LLM outputs, while SCALM (Li et al., 2024c) clusters queries to uncover meaningful semantic patterns. Compared with those two matching approaches in Table 1.

Compression Techniques. To reduce inference overhead, tensor quantization and compact representations are used to balance performance and efficiency (Wang et al., 2024b).

Quantization-based Compression This method reduces memory by lowering precision. (Sheng et al., 2023) applies group-wise 4-bit KV cache quantization without I/O overhead. AWQ (Lin et al., 2024c) highlights that quantizing non-salient

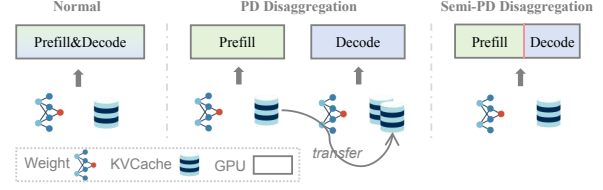


Figure 5: Comparing with normal, PD and Semi-PD Disaggregation.

weights reduces quantization loss. Kivi (Zirui Liu et al., 2024) uses per-channel/token cache quantization. MiniCache (Liu et al., 2024a) exploits inter-layer KV similarity for compression. Atom (Zhao et al., 2024b) adopts mixed-precision and fine-grained quantization. QServe (Lin et al., 2024e) co-designs W4A8KV4 quantization to boost GPU efficiency. OTT quantization (Su et al., 2025), which incorporates outlier token tracing, has made significant progress.

Compact Encoding Architectures It is also desirable to use smaller matrix representations instead of the previous heavy matrix. CacheGen (Liu et al., 2024c) employs a custom tensor encoder to compress KV cache into compact bitstreams, saving bandwidth with minimal decoding overhead.

3.5 PD Disaggregation

PD disaggregation separates computation-bound prefill from memory-bound decoding, enabling specialized optimization for each (Figure 5).

DistServe (Zhong et al., 2024) minimizes communication overhead by optimizing resource allocation and bandwidth-aware placement. Splitwise (Patel et al., 2024b) and HEXGEN-2 (Jiang et al., 2025b) explores homogeneous and heterogeneous device designs to optimize cost, throughput, and power. DéjàVu (Strati et al., 2024) eliminates pipeline stalls via microbatch swapping and state replication. Mooncake (Qin et al., 2024) uses a KVCache-centric disaggregated design, leveraging idle CPU, DRAM, and SSD with early rejection under load. TetriInfer (Hu et al., 2024b) applies two-level scheduling with resource prediction to avoid decoding hotspots. P/D-Serve (Jin et al., 2024b) tackles LLM deployment challenges via fine-grained prefill/decode organization, dynamic adjustments, on-demand request allocation, and efficient cache transmission.

KV Cache Transfer Problem. Semi-PD (Hong et al., 2025) separates the two stages on the same GPU by leveraging the streaming multiprocessor (SM) method, eliminating KV cache transfer over-

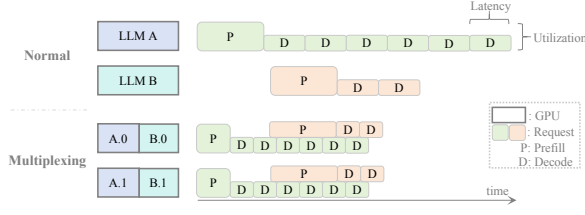


Figure 6: An example of Normal and Multiplexing.

head and enabling weight sharing. FlowKV (Li et al., 2025b) reduces KV cache migration time by minimizing communication, optimizing memory continuity, and using strategies like bidirectional alignment and merged transfers.

3.6 Multiplexing

Running a single LLM across multiple GPUs often leads to idle time due to varying model sizes and usage. To improve utilization, multiplexing multiple LLMs on a single GPU has gained attention (as shown in Figure 6).

AlpaServe (Li et al., 2023b) adopts a more rudimentary approach, resulting in suboptimal GPU utilization. To improve efficiency, MIG (NVIDIA, 2025) introduces hardware-level partitioning by dividing a single physical GPU into multiple isolated virtual instances, allowing parallel execution of multiple models. In contrast, MuxServe (Duan et al., 2024) partitions GPU Streaming Multiprocessors (SMs) via CUDA MPS and fully shares memory, requiring workload-dependent ratio tuning, they both well-suited for stable workloads with strong isolation requirements. Unlike these static methods, QLM (Patke et al., 2024) switches models on the GPU to handle different requests. Prism (Yu et al., 2025b) goes further by dynamically adjusting GPU resource allocation, enabling flexible and efficient multi-model sharing.

4 Special Architecture

4.1 MoE

MoE models, known for parameter sparsity, excel in LLMs (e.g., DeepSeek-V3 (DeepSeek-AI et al., 2024), Mixtral 8x7B (Jiang et al., 2024a)). Key inference latency challenges include expert parallelism, load balancing, and All-to-All communication (Figure 7), with Liu et al. (2024b) offering a comprehensive optimization survey.

Expert Placement. Tutel (Hwang et al., 2023) introduces switchable parallelism and dynamic pipelining without extra overhead, while

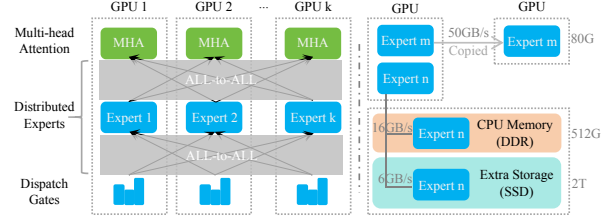


Figure 7: This figure illustrates a MoE architecture, highlighting expert placement, All-to-All communication (left), and load balancing (right). On the right, high-traffic Expert m and low-traffic Expert n are shown. For example, two strategies are presented: replicating m to a new GPU or offloading n to free space for m.

DeepSpeed-MoE (Rajbhandari et al., 2022) combines expert parallelism (He et al., 2021; Lepikhin et al., 2020) with expert-slicing. CoEL (Li et al., 2025a) leverages the sparse activation of MoE to enable resource collaboration both within and across devices. fMOE (Yu et al., 2025a) effectively guides expert prefetching, caching, and offloading decisions through its expert graph and semantically enhanced search mechanism.

Expert Load Balancing. Imbalanced token distribution causes device underutilization. Expert Buffering (Huang et al., 2023) allocates active experts to GPUs and others to CPUs, pairing high- and low-load experts using historical data. Brainstorm (Cui et al., 2023) dynamically assigns GPU units based on load, while Lynx (Gupta et al., 2024) adaptively reduces active experts. ExpertChoice (Zhou et al., 2022) selects top-k tokens per expert, rather than the reverse. High-load experts in DeepSeek-V3 (DeepSeek-AI et al., 2024) are identified using deployment statistics and periodically duplicated to optimize performance.

All-to-All Communication. Expert processing involves all-to-all exchanges for token dispatch and output gathering. Tutel (Hwang et al., 2023) uses a 2D hierarchical All-to-All algorithm, Aurora (Li et al., 2024b) optimizes token transmission order during All-to-All exchanges, and Lina (Li et al., 2023a) prioritizes All-to-All operations over concurrent All-Reduce whenever feasible, leveraging tensor partitioning to improve performance.

4.2 LoRA

LoRA (Hu et al., 2021; Chen et al., 2024; Dettmers et al., 2023) adapts LLMs to various tasks with small, trainable adapters. CaraServe (Li et al., 2024e) enables GPU-efficient, cold-start-free, SLO-

aware serving via model multiplexing, CPU-GPU coordination, and rank-aware scheduling. dLoRA (Wu et al., 2024c) dynamically merges and unmerges adapters with the base model, and migrates requests and adapters across worker replicas.

4.3 Speculative Decoding

Speculative decoding (Xia et al., 2024; Wang et al., 2024a) speeds up inference by generating draft tokens with smaller LLMs and verifying them in parallel with target LLM, reducing latency and costs without quality loss. SpecInfer (Miao et al., 2024a) uses tree-based speculative inference for faster distributed and single-GPU offloading inference. AdaServe (Li et al., 2025d) and SpecServe (Huang et al., 2025) are tailored for customized SLOs, employing different fine-grained speculative decoding methods for LLMs with varying SLO requirements.

5 Specific Scenarios

5.1 Long Context

As LLMs evolve, context lengths have expanded significantly, reaching hundreds of thousands or even millions of tokens (moonshot, 2023). This growth presents both opportunities and challenges for distributed deployment, computation, and storage, especially in parallel processing, attention computation, and KV cache management.

Parallel Processing. Loongserve (Wu et al., 2024a) enhances this with elastic sequence parallelism for efficient long-context LLM serving.

Attention Computation. The attention mechanism encounters significant challenges in parallel processing and resource management. RingAttention (Liu et al., 2023) uses blockwise self-attention and FFN computation to distribute long sequences across devices, overlapping KV communication with attention. StripedAttention (Brandon et al., 2023), an extension of RingAttention, addresses imbalances from causal attention’s triangular structure. DistAttention (Lin et al., 2024a) subdivides attention across GPUs, avoiding cache transfer during decoding and enabling partitioning for arbitrary sequence lengths with minimal data transfer. InstInfer (Pan et al., 2024b) offloads attention and data to Computational Storage Drives, reducing KV transfer overheads significantly.

KV Cache Management. Efficient storage for growing KV cache is crucial for generating new

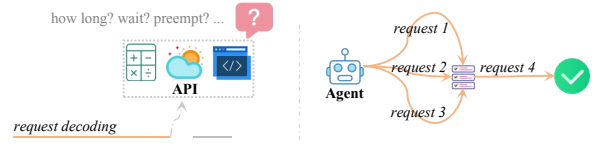


Figure 8: Overview of augmented LLMs, such as LLM with API, and LLM-based Agent.

tokens. Infinite-LLM (Lin et al., 2024a) manages dynamic LLM contexts by scheduling cache at the cluster level, balancing resources, and maximizing throughput. InfiniGen (Lee et al., 2024) optimizes cache management in CPU memory for offloading-based systems. Marconi (Pan et al., 2024a) introduces tailored admission and eviction policies for hybrid models, using experimental and theoretical analysis to show that personalized cache sizing per layer reduces memory usage significantly.

5.2 RAG

RAG enables LLMs to retrieve external knowledge for responses, but the diversity and complexity of processing pose challenges in optimizing latency and KV cache storage for large retrieval contexts.

Workflow Scheduling. Several recent innovations have focused on improving the efficiency, flexibility, and optimization of RAG workflows. PipeRAG (Jiang et al., 2024b) improves efficiency via pipeline parallelism, flexible retrieval intervals, and performance-driven quality adjustment. Teola (Tan et al., 2024) models LLM workflows as data flow nodes (e.g., Embedding, Indexing, Searching) for precise execution control. RaLMSpec (Zhang et al., 2024d) employs speculative retrieval with batched verification to reduce serving overhead. RAGServe (Ray et al., 2024) schedules queries and adjusts RAG configurations (e.g., text chunks, synthesis methods) to balance quality and latency.

Storage Optimization. Efficient storage management is critical for RAG systems, particularly in handling large-scale KV caches. Recent studies include RAGCache (Jin et al., 2024a), which employs knowledge trees and dynamic speculative pipelining to reduce redundancy. SparseRAG (Zhu et al., 2024) manages cache efficiently with pre-filling and selective decoding, focusing on relevant tokens. CacheBlend (Yao et al., 2024) reuses cache and selects tokens based on a fixed percentage to recompute KV values for partial updates, enhancing efficiency and reducing latency. In contrast to CacheBlend, EPIC (Hu et al., 2024c)

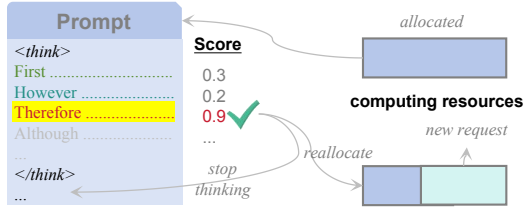


Figure 9: An example of efficient reasoning serving.

introduces position-independent context caching via static sparse computation, recomputing only a small number of tokens at the beginning of each block.

5.3 Augmented LLMs

LLMs increasingly integrate with external tools like APIs and Agents (as shown in Figure 8). APISERVE (Abhyankar et al., 2024) dynamically manages GPU resources for external APIs, while LAMPS (Shahout et al., 2024a) leverages predicting memory usage. Parrot (Lin et al., 2024b) optimizes scheduling by identifying request dependencies and commonalities, particularly in Agent scenarios, using Semantic Variables to tag each request. While Parrot is a pioneering approach to addressing agent-related challenges, it has significant limitations in practice. Tempo (Zhang et al., 2025) handles collective LLM requests by modeling execution dependencies as a DAG, with nodes as requests and edges as dependencies. Luo et al. (2025) considers the Agent as a program-level problem rather than a request-level one, which offers a valuable insight. Li et al. (2025c) uses a fluid queuing model to analyze Agent workload stability.

5.4 Test-Time Reasoning

Inference-time algorithms (OpenAI, 2024) enhance the reasoning ability (Ji et al., 2025; Qu et al., 2025) of LLMs but achieve this by generating a large number of tokens, which can strain computational resources (as shown in Figure 9). Dynasor (Fu et al., 2024b) introduces the Certainindex metric to dynamically track a model’s reasoning progress, adjust computational resources based on task difficulty, and proactively terminate unpromising requests. Damani et al. (2024) optimizes resource allocation (e.g., different LLMs, computing budgets) by using a built-in reward model to assess the marginal benefit of additional computation for each request.

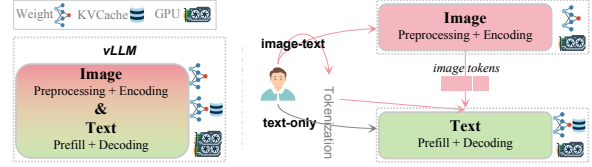


Figure 10: vLLM vs. decouple way in MLLM.

5.5 Multimodal LLM

As Multimodal LLM (MLLM) matures, the resulting surge in traffic presents significant deployment challenges. To address stage heterogeneity, image encoding latency, modality interference, and long-tail workload issues in multi-modal inference, ModServe (Qiu et al., 2025) introduces a decoupled architecture, as shown in Figure 10. It separates key stages such as image preprocessing, encoding, and LLM operations into independently scalable components, enabling fine-grained resource management and workload-aware optimization.

6 Future Works

Given the rapid evolution of LLM inference services, we present several recommendations for future research. *General Agent Scheduling*: Real-world demands for integrating audio, image, video, and text agents add complexity to inference serving. *Intelligent LLM Inference Service*: Utilizing the capabilities of smaller LLMs to optimize the deployment, scheduling, and storage management of larger LLMs. *Safety and Privacy*: As most services rely on cloud computing, it is essential to prevent cache leaks and ensure that any leaked data cannot be used to reconstruct user conversations. We hope that these suggestions will provide valuable insights for advancing future research.

7 Conclusion

This paper presents a comprehensive review of existing methods, covering general strategies, specialized architectures, and emerging scenarios. Specifically, the general approaches include model placement, scheduling, memory, PD disaggregation, and multiplexing optimization. We also discuss specialized architectures such as MoE, LoRA, and speculative decoding, as well as emerging scenarios like long-context problem, RAG, augmented LLM, test-time reasoning, and multimodal. We hope this work provides valuable insights for ongoing research in this critical area.

8 Limitations

This paper provides a summary and categorization of methods for LLM inference services, with the following limitations:

- While we have not conducted detailed experiments or comparative analyses of all approaches, we have presented critical insights on some key points.
- In the foundational sections, there may be some overlap with previous surveys, primarily involving a few classic papers. However, since the most recent survey (Zhou et al., 2024) only covers literature up to July 2024, many of the references we include were published after that date.
- The scope of this survey is time-sensitive, with coverage concluding at the end of April 2025.

References

- Marah Abidin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Xin Wang, Rachel Ward, Yue Wu, Dingli Yu, Cyril Zhang, and Yi Zhang. 2024. [Phi-4 technical report](#). *Preprint*, arXiv:2412.08905.
- Reyna Abhyankar, Zijian He, Vikranth Srivatsa, Hao Zhang, and Yiyang Zhang. 2024. [Infercept: Efficient intercept support for augmented large language model inference](#). *Preprint*, arXiv:2402.01869.
- Megha Agarwal, Asfandiyar Qureshi, Nikhil Sardana, Linden Li, Julian Quevedo, and Daya Khudia. 2023. [Llm inference performance engineering: Best practices](#).
- Saurabh Agarwal, Anyong Mao, Aditya Akella, and Shivaram Venkataraman. 2024. [Symphony: Improving memory management for llm inference workloads](#). *Preprint*, arXiv:2412.16434.
- Amey Agrawal, Anmol Agarwal, Nitin Kedia, Jayashree Mohan, Souvik Kundu, Nipun Kwatra, Ramachandran Ramjee, and Alexey Tumanov. 2024a. [Etalon: Holistic performance evaluation framework for llm inference systems](#). *Preprint*, arXiv:2407.07000.
- Amey Agrawal, Nitin Kedia, Jayashree Mohan, Ashish Panwar, Nipun Kwatra, Bhargav Gulavani, Ramachandran Ramjee, and Alexey Tumanov. 2024b. [Vidur: A large-scale simulation framework for llm inference](#). *Preprint*, arXiv:2405.05465.
- Amey Agrawal, Nitin Kedia, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav Gulavani, Alexey Tumanov, and Ramachandran Ramjee. 2024c. [Taming Throughput-Latency tradeoff in LLM inference with Sarathi-Serve](#). In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, pages 117–134, Santa Clara, CA. USENIX Association.
- Reza Yazdani Aminabadi, Samyam Rajbhandari, Minjia Zhang, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Jeff Rasley, Shaden Smith, Olatunji Ruwase, and Yuxiong He. 2022. [Deepspeed inference: Enabling efficient inference of transformer models at unprecedented scale](#). *Preprint*, arXiv:2207.00032.
- Abhimanyu Bambhaniya, Ritik Raj, Geonhwa Jeong, Souvik Kundu, Sudarshan Srinivasan, Midhilesh Elavazhagan, Madhu Kumar, and Tushar Krishna. 2024. [Demystifying platform requirements for diverse llm inference use cases](#). *Preprint*, arXiv:2406.01698.
- Fu Bang. 2023. [GPTCache: An open-source semantic cache for LLM applications enabling faster answers and cost savings](#). In *Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023)*, pages 212–218, Singapore. Association for Computational Linguistics.
- Alexander Borzunov, Max Ryabinin, Artem Chumachenko, Dmitry Baranchuk, Tim Dettmers, Younes Belkada, Pavel Samygin, and Colin Raffel. 2023. [Distributed inference and fine-tuning of large language models over the internet](#). *Preprint*, arXiv:2312.08361.
- William Brandon, Aniruddha Nrusimha, Kevin Qian, Zachary Ankner, Tian Jin, Zhiye Song, and Jonathan Ragan-Kelley. 2023. [Striped attention: Faster ring attention for causal transformers](#). *Preprint*, arXiv:2311.09431.
- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2024. [Longlora: Efficient fine-tuning of long-context large language models](#). *Preprint*, arXiv:2309.12307.
- Ke Cheng, Wen Hu, Zhi Wang, Peng Du, Jianguo Li, and Sheng Zhang. 2024a. [Enabling efficient batch serving for lmaas via generation length prediction](#). *arXiv preprint arXiv:2406.04785*.
- Ke Cheng, Wen Hu, Zhi Wang, Hongen Peng, Jianguo Li, and Sheng Zhang. 2024b. [Slice-level scheduling for high throughput and load balanced llm serving](#). *Preprint*, arXiv:2406.13511.
- Rongxin Cheng, Yifan Peng, Yuxin Lai, Xingda Wei, Rong Chen, and Haibo Chen. 2024c. [Kunserve: Elastic and efficient large language model serving with parameter-centric memory management](#). *Preprint*, arXiv:2412.18169.
- Weihao Cui, Zhenhua Han, Lingji Ouyang, Yichuan Wang, Ningxin Zheng, Lingxiao Ma, Yuqing Yang,

705	Fan Yang, Jilong Xue, Lili Qiu, Lidong Zhou, Quan	tivizing reasoning capability in llms via reinforce-	766
706	Chen, Haisheng Tan, and Minyi Guo. 2023. Opti-	ment learning. <i>Preprint</i> , arXiv:2501.12948.	767
707	mizing dynamic neural networks with brainstorm . In		
708	<i>17th USENIX Symposium on Operating Systems De-</i>	DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-	768
709	<i>sign and Implementation (OSDI 23)</i> , pages 797–815,	uan Wang, Bochao Wu, Chengda Lu, Chenggang	769
710	Boston, MA. USENIX Association.	Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,	770
		Damai Dai, Daya Guo, Dejian Yang, Deli Chen,	771
711	Mehul Damani, Idan Shenfeld, Andi Peng, Andreea	Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,	772
712	Bobu, and Jacob Andreas. 2024. Learning how hard	Fuli Luo, Guangbo Hao, Guanting Chen, Guowei	773
713	to think: Input-adaptive allocation of lm computation .	Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng	774
714	<i>Preprint</i> , arXiv:2410.04707.	Wang, Haowei Zhang, Honghui Ding, Huajian Xin,	775
		Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang,	776
715	DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,	Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang,	777
716	Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,	Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie	778
717	Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,	Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu,	779
718	Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong	Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean	780
719	Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue,	Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao,	781
720	Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu,	Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang,	782
721	Chenggang Zhao, Chengqi Deng, Chenyu Zhang,	Mingchuan Zhang, Minghua Zhang, Minghui Tang,	783
722	Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji,	Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang,	784
723	Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo,	Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu	785
724	Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang,	Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge,	786
725	Han Bao, Hanwei Xu, Haocheng Wang, Honghui	Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin	787
726	Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li,	Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao	788
727	Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang	Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu,	789
728	Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L.	Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu	790
729	Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai	Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou,	791
730	Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai	Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun,	792
731	Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong	W. L. Xiao, Wangding Zeng, Wanxia Zhao, Wei An,	793
732	Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan	Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu,	794
733	Zhang, Minghua Zhang, Minghui Tang, Meng Li,	Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang,	795
734	Miaojun Wang, Mingming Li, Ning Tian, Panpan	Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen,	796
735	Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen,	Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen,	797
736	Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan,	Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin	798
737	Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen,	Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu,	799
738	Shanghao Lu, Shangyan Zhou, Shanhuang Chen,	Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang,	800
739	Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng	Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li,	801
740	Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing	Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yan-	802
741	Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun,	hong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao	803
742	T. Wang, Wangding Zeng, Wanxia Zhao, Wen Liu,	Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu,	804
743	Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao	Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong,	805
744	Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan	Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yix-	806
745	Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin	uan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,	807
746	Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li,	Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue	808
747	Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin,	Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan	809
748	Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxi-	Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxi-	810
749	ang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang,	ang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z.	811
750	Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang	Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu,	812
751	Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng	Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan	813
752	Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,	Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhi-	814
753	Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang,	gang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu,	815
754	Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,	Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu,	816
755	Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yu-	Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi	817
756	jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,	Gao, and Zizheng Pan. 2024. Deepseek-v3 technical	818
757	Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu,	report . <i>Preprint</i> , arXiv:2412.19437.	819
758	Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu,	Deepspeed. 2023. Deepspeed-mii .	820
759	Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan,		
760	Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and	821
761	Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao,	Luke Zettlemoyer. 2023. Qlora: Efficient finetuning	822
762	Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi-	of quantized llms . <i>Preprint</i> , arXiv:2305.14314.	823
763	jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song,		
764	Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu	Jiangfei Duan, Runyu Lu, Haojie Duanmu, Xiuhong Li,	824
765	Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incen-	Xingcheng Zhang, Dahua Lin, Ion Stoica, and Hao	825
		Zhang. 2024. Muxserve: Flexible spatial-temporal	826

827 [multiplexing for multiple llm serving](#). *Preprint*,
828 [arXiv:2404.02015](#).

829 William Fedus, Barret Zoph, and Noam Shazeer. 2022.
830 [Switch transformers: Scaling to trillion parameter](#)
831 [models with simple and efficient sparsity](#). *Preprint*,
832 [arXiv:2101.03961](#).

833 Yao Fu, Leyang Xue, Yeqi Huang, Andrei-Octavian
834 Brabete, Dmitrii Ustiugov, Yuvraj Patel, and Luo
835 Mai. 2024a. [Serverlessllm: Low-latency server-](#)
836 [less inference for large language models](#). *Preprint*,
837 [arXiv:2401.14351](#).

838 Yichao Fu, Junda Chen, Siqi Zhu, Zheyu Fu, Zhong-
839 dongming Dai, Aurick Qiao, and Hao Zhang. 2024b.
840 [Efficiently serving llm reasoning programs with cer-](#)
841 [taindex](#). *Preprint*, [arXiv:2412.20993](#).

842 Yichao Fu, Siqi Zhu, Runlong Su, Aurick Qiao, Ion
843 Stoica, and Hao Zhang. 2024c. Efficient llm
844 scheduling by learning to rank. *arXiv preprint*
845 [arXiv:2408.15792](#).

846 Bin Gao, Zhuomin He, Puru Sharma, Qingxuan Kang,
847 Djordje Jevdjic, Junbo Deng, Xingkun Yang, Zhou
848 Yu, and Pengfei Zuo. 2024. {Cost-Efficient} large
849 language model serving for multi-turn conversations
850 with {CachedAttention}. In *2024 USENIX Annual*
851 *Technical Conference (USENIX ATC 24)*, pages 111–
852 126.

853 [ggml.org](#). 2022. [llamacpp](#).

854 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
855 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
856 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
857 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
858 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
859 tra, Archie Sravankumar, Artem Korenev, Arthur
860 Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-
861 driguez, Austen Gregerson, Ava Spataru, Baptiste
862 Roziere, Bethany Biron, Binh Tang, Bobbie Chern,
863 Charlotte Caucheteux, Chaya Nayak, Chloe Bi,
864 Chris Marra, Chris McConnell, Christian Keller,
865 Christophe Touret, Chunyang Wu, Corinne Wong,
866 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-
867 lonsius, Daniel Song, Danielle Pintz, Danny Livshits,
868 Danny Wyatt, David Esiobu, Dhruv Choudhary,
869 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
870 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,
871 Elina Lobanova, Emily Dinan, Eric Michael Smith,
872 Filip Radenovic, Francisco Guzmán, Frank Zhang,
873 Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-
874 derson, Govind Thattai, Graeme Nail, Gregoire Mi-
875 alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,
876 Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
877 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-
878 han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
879 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,
880 Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
881 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
882 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,
883 Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
884 Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-
885 teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,

Kartik Upasani, Kate Plawiak, Ke Li, Kenneth
Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,
Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal
Lakhotia, Lauren Rantala-Yeary, Laurens van der
Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,
Louis Martin, Lovish Madaan, Lubo Malo, Lukas
Blecher, Lukas Landzaat, Luke de Oliveira, Madeline
Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar
Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew
Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-
badur, Mike Lewis, Min Si, Mitesh Kumar Singh,
Mona Hassan, Naman Goyal, Narjes Torabi, Niko-
lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,
Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick
Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-
sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,
Praveen Krishnan, Punit Singh Koura, Puxin Xu,
Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj
Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,
Robert Stojnic, Roberta Raileanu, Rohan Maheswari,
Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-
nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan
Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-
hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-
hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-
ran Narang, Sharath Raparthy, Sheng Shen, Shengye
Wan, Shruti Bhosale, Shun Zhang, Simon Van-
denhende, Soumya Batra, Spencer Whitman, Sten
Sootla, Stephane Collet, Suchin Gururangan, Syd-
ney Borodinsky, Tamar Herman, Tara Fowler, Tarek
Sheasha, Thomas Georgiou, Thomas Scialom, Tobias
Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal
Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh
Ramanathan, Viktor Kerkez, Vincent Gouget, Vir-
ginie Do, Vish Voleti, Vitor Albiero, Vladan Petro-
vic, Weiwei Chu, Wenhan Xiong, Wenyan Fu, Whit-
ney Meers, Xavier Martinet, Xiaodong Wang, Xi-
aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-
feng Xie, Xuchao Jia, Xuwei Wang, Yaelle Gold-
schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,
Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,
Zacharie Delpierre Coudert, Zheng Yan, Zhengxing
Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-
vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,
Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,
Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
Baevski, Allie Feinstein, Amanda Kallet, Amit San-
gani, Amos Teo, Anam Yunus, Andrei Lupu, An-
dres Alvarado, Andrew Caples, Andrew Gu, Andrew
Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-
dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-
ajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,
Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-
dan, Beau James, Ben Maurer, Benjamin Leonhardi,
Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi
Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-
cock, Bram Wasti, Brandon Spence, Brani Stojkovic,
Brian Gamido, Britt Montalvo, Carl Parker, Carly
Burton, Catalina Mejia, Ce Liu, Changan Wang,
Changkyu Kim, Chao Zhou, Chester Hu, Ching-
Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-
ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,
Daniel Kreymer, Daniel Li, David Adkins, David
Xu, Davide Testuggine, Delia David, Devi Parikh,

950	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	Matthews, Timothy Chou, Tzook Shaked, Varun	1014
951	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	1015
952	Elaine Montgomery, Eleonora Presani, Emily Hahn,	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	1016
953	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	1017
954	ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	1018
955	Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	1019
956	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	1020
957	Seide, Gabriela Medina Florez, Gabriella Schwarz,	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	1021
958	Gada Badeer, Georgia Sweet, Gil Halpern, Grant	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	1022
959	Herman, Grigory Sizov, Guangyi, Zhang, Guna	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	1023
960	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	Yundi Qian, Yunlu Li, Yuze He, Zach Rait, Zachary	1024
961	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	1025
962	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd	1026
963	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	of models . <i>Preprint</i> , arXiv:2407.21783.	1027
964	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,		
965	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	Tyler Griggs, Xiaoxuan Liu, Jiaxiang Yu, Doyoung	1028
966	Geboski, James Kohli, Janice Lam, Japhet Asher,	Kim, Wei-Lin Chiang, Alvin Cheung, and Ion Stoica.	1029
967	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	2024. Mélange: Cost efficient large language model	1030
968	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	serving by exploiting gpu heterogeneity . <i>Preprint</i> ,	1031
969	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	arXiv:2404.14527.	1032
970	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-		
971	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	Vima Gupta, Kartik Sinha, Ada Gavrilovska, and	1033
972	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	Anand Padmanabha Iyer. 2024. Lynx: Enabling ef-	1034
973	delwal, Katayoun Zand, Kathy Matosich, Kaushik	ficient moe inference through dynamic batch-aware	1035
974	Veeraraghavan, Kelly Michelena, Keqian Li, Ki-	expert selection . <i>Preprint</i> , arXiv:2411.08982.	1036
975	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle		
976	Huang, Lailin Chen, Lakshya Garg, Lavender A,	Zixu Hao, Huiqiang Jiang, Shiqi Jiang, Ju Ren, and Ting	1037
977	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	Cao. 2024. Hybrid slm and llm for edge-cloud col-	1038
978	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	laborative inference . In <i>Proceedings of the Workshop</i>	1039
979	edt, Madian Khabisa, Manav Avalani, Manish Bhatt,	<i>on Edge and Mobile Foundation Models</i> , EdgeFM	1040
980	Martynas Mankus, Matan Hasson, Matthew Lennie,	'24, page 36–41, New York, NY, USA. Association	1041
981	Matthias Reso, Maxim Groshev, Maxim Naumov,	for Computing Machinery.	1042
982	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.		
983	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	Aaron Harlap, Deepak Narayanan, Amar Phanishayee,	1043
984	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	Vivek Seshadri, Nikhil Devanur, Greg Ganger, and	1044
985	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	Phil Gibbons. 2018. Pipedream: Fast and ef-	1045
986	Mo Metanat, Mohammad Rastegari, Munish Bansal,	ficient pipeline parallel dnn training . <i>Preprint</i> ,	1046
987	Nandhini Santhanam, Natascha Parks, Natasha	arXiv:1806.03377.	1047
988	White, Navyata Bawa, Nayan Singhal, Nick Egebo,		
989	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	Jiaao He, Jiezhong Qiu, Aohan Zeng, Zhilin Yang,	1048
990	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	Jidong Zhai, and Jie Tang. 2021. Fastmoe: A	1049
991	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	fast mixture-of-expert training system . <i>Preprint</i> ,	1050
992	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro	arXiv:2103.13262.	1051
993	Rittner, Philip Bontrager, Pierre Roux, Piotr		
994	Dollar, Polina Zvyagina, Prashant Ratanchandani,	Jiaao He and Jidong Zhai. 2024. Fastdecode: High-	1052
995	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	throughput gpu-efficient llm serving using heteroge-	1053
996	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	neous pipelines . <i>Preprint</i> , arXiv:2403.11421.	1054
997	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,		
998	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	Ying He, Jingcheng Fang, F. Richard Yu, and Victor C.	1055
999	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	Leung. 2024. Large language models (llms) infer-	1056
1000	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	ence offloading and resource allocation in cloud-edge	1057
1001	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	computing: An active inference approach . <i>IEEE</i>	1058
1002	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	<i>Transactions on Mobile Computing</i> , 23(12):11253–	1059
1003	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	11264.	1060
1004	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,		
1005	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	Soka Hisaharo, Yuki Nishimura, and Aoi Takahashi.	1061
1006	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	2024. Optimizing llm inference clusters for en-	1062
1007	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	hanced performance and energy efficiency . <i>Authorea</i>	1063
1008	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	<i>Preprints</i> .	1064
1009	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,		
1010	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	Connor Holmes, Masahiro Tanaka, Michael Wyatt, Am-	1065
1011	Subramanian, Sy Choudhury, Sydney Goldman, Tal	mar Ahmad Awan, Jeff Rasley, Samyam Rajbhan-	1066
1012	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	dari, Reza Yazdani Aminabadi, Heyang Qin, Arash	1067
1013	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	Bakhtiari, Lev Kurilenko, and Yuxiong He. 2024.	1068
		Deepspeed-fastgen: High-throughput text generation	1069
		for llms via mii and deepspeed-inference . <i>Preprint</i> ,	1070
		arXiv:2401.08671.	1071

1072	Ke Hong, Lufang Chen, Zhong Wang, Xiuhong Li, Qiuli	Kunal Jain, Anjaly Parayil, Ankur Mallick, Esha	1128
1073	Mao, Jianping Ma, Chao Xiong, Guanyu Wu, Buhe	Choukse, Xiaoting Qin, Jue Zhang, Íñigo Goiri, Rujia	1129
1074	Han, Guohao Dai, Yun Liang, and Yu Wang. 2025.	Wang, Chetan Bansal, Victor Rühle, Anoop Kulkarni,	1130
1075	semi-pd: Towards efficient llm serving via phase-	Steve Kofsky, and Saravan Rajmohan. 2024. Intel-	1131
1076	wise disaggregated computation and unified storage.	ligent router for llm workloads: Improving perform-	1132
1077	<i>Preprint</i> , arXiv:2504.19867.	ance through workload-aware scheduling. <i>Preprint</i> ,	1133
		arXiv:2408.13510.	1134
1078	Cunchen Hu, Heyang Huang, Junhao Hu, Jiang Xu,	Suhas Jayaram Subramanya, Daiyaan Arfeen, Shouxu	1135
1079	Xusheng Chen, Tao Xie, Chenxi Wang, Sa Wang,	Lin, Aurick Qiao, Zhihao Jia, and Gregory R Ganger.	1136
1080	Yungang Bao, Ninghui Sun, and Yizhou Shan.	2023. Sia: Heterogeneity-aware, goodput-optimized	1137
1081	2024a. Memserve: Context caching for disaggre-	ml-cluster scheduling. In <i>Proceedings of the 29th</i>	1138
1082	gated llm serving with elastic memory pool. <i>Preprint</i> ,	<i>Symposium on Operating Systems Principles</i> , pages	1139
1083	arXiv:2406.17565.	642–657.	1140
1084	Cunchen Hu, Heyang Huang, Liangliang Xu, Xusheng	Yixin Ji, Juntao Li, Hai Ye, Kaixin Wu, Kai Yao, Jia Xu,	1141
1085	Chen, Jiang Xu, Shuang Chen, Hao Feng, Chenxi	Linjian Mo, and Min Zhang. 2025. Test-time com-	1142
1086	Wang, Sa Wang, Yungang Bao, Ninghui Sun, and	pute: from system-1 thinking to system-2 thinking.	1143
1087	Yizhou Shan. 2024b. Inference without interference:	<i>Preprint</i> , arXiv:2501.02497.	1144
1088	Disaggregate llm inference for mixed downstream		
1089	workloads. <i>Preprint</i> , arXiv:2401.11181.		
1090	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Albert Q. Jiang, Alexandre Sablayrolles, Antoine	1145
1091	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	Roux, Arthur Mensch, Blanche Savary, Chris	1146
1092	Weizhu Chen. 2021. Lora: Low-rank adaptation of	Bamford, Devendra Singh Chaplot, Diego de las	1147
1093	large language models. <i>Preprint</i> , arXiv:2106.09685.	Casas, Emma Bou Hanna, Florian Bressand, Gi-	1148
1094	Junhao Hu, Wenrui Huang, Haoyi Wang, Weidong	anna Lengyel, Guillaume Bour, Guillaume Lam-	1149
1095	Wang, Tiancheng Hu, Qin Zhang, Hao Feng,	ple, Lélío Renard Lavaud, Lucile Saulnier, Marie-	1150
1096	Xusheng Chen, Yizhou Shan, and Tao Xie. 2024c.	Anne Lachaux, Pierre Stock, Sandeep Subramanian,	1151
1097	Epic: Efficient position-independent context caching	Sophia Yang, Szymon Antoniak, Teven Le Scao,	1152
1098	for serving large language models. <i>Preprint</i> ,	Théophile Gervet, Thibaut Lavril, Thomas Wang,	1153
1099	arXiv:2410.15332.	Timothée Lacroix, and William El Sayed. 2024a.	1154
1100	Haiyang Huang, Newsha Ardalani, Anna Sun, Liu	Mixtral of experts. <i>Preprint</i> , arXiv:2401.04088.	1155
1101	Ke, Hsien-Hsin S. Lee, Anjali Sridhar, Shruti Bhos-	Wenqi Jiang, Shuai Zhang, Boran Han, Jie Wang,	1156
1102	ale, Carole-Jean Wu, and Benjamin Lee. 2023. To-	Bernie Wang, and Tim Kraska. 2024b. Piperag: Fast	1157
1103	wards moe deployment: Mitigating inefficiencies	retrieval-augmented generation via algorithm-system	1158
1104	in mixture-of-expert (moe) inference. <i>Preprint</i> ,	co-design. <i>Preprint</i> , arXiv:2403.05676.	1159
1105	arXiv:2303.06182.		
1106	Kaiyu Huang, Hao Wu, Zubo Shi, Han Zou, Minchen	Youhe Jiang, Fangcheng Fu, Xiaozhe Yao, Guoliang	1160
1107	Yu, and Qingjiang Shi. 2025. Specserve: Ef-	He, Xupeng Miao, Ana Klimovic, Bin Cui, Binhang	1161
1108	ficient and slo-aware large language model serv-	Yuan, and Eiko Yoneki. 2025a. Demystifying cost-	1162
1109	ing with adaptive speculative decoding. <i>Preprint</i> ,	efficiency in llm serving over heterogeneous gpus.	1163
1110	arXiv:2503.05096.	<i>Preprint</i> , arXiv:2502.00722.	1164
1111	Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan	Youhe Jiang, Ran Yan, Xiaozhe Yao, Yang Zhou, Beidi	1165
1112	Firat, Mia Xu Chen, Dehao Chen, HyounJoong Lee,	Chen, and Binhang Yuan. 2024c. Hexgen: Genera-	1166
1113	Jiquan Ngiam, Quoc V. Le, Yonghui Wu, and Zhifeng	tive inference of large language model over heteroge-	1167
1114	Chen. 2019. Gpipe: Efficient training of giant neu-	neous environment. <i>Preprint</i> , arXiv:2311.11514.	1168
1115	ral networks using pipeline parallelism. <i>Preprint</i> ,		
1116	arXiv:1811.06965.	Youhe Jiang, Ran Yan, and Binhang Yuan. 2025b.	1169
1117	HuggingFace. 2023. Tgi (text generation inference).	Hexgen-2: Disaggregated generative inference of	1170
1118	Changho Hwang, Wei Cui, Yifan Xiong, Ziyue Yang,	LLMs in heterogeneous environment.	1171
1119	Ze Liu, Han Hu, Zilong Wang, Rafael Salas, Jithin	Chao Jin, Zili Zhang, Xuanlin Jiang, Fangyue Liu, Xin	1172
1120	Jose, Prabhat Ram, Joe Chau, Peng Cheng, Fan	Liu, Xuanzhe Liu, and Xin Jin. 2024a. Ragcache:	1173
1121	Yang, Mao Yang, and Yongqiang Xiong. 2023. Tu-	Efficient knowledge caching for retrieval-augmented	1174
1122	tel: Adaptive mixture-of-experts at scale. <i>Preprint</i> ,	generation. <i>Preprint</i> , arXiv:2404.12457.	1175
1123	arXiv:2206.03382.		
1124	Saki Imai, Rina Nakazawa, Marcelo Amaral, Sunyanan	Yibo Jin, Tao Wang, Huimin Lin, Mingyang Song,	1176
1125	Choochothaew, and Tatsuhiko Chiba. 2024. Pre-	Peiyang Li, Yipeng Ma, Yicheng Shan, Zhengfan	1177
1126	dicting llm inference latency: A roofline-driven ml	Yuan, Cailong Li, Yajing Sun, Tiandeng Wu, Xing	1178
1127	method.	Chu, Ruizhi Huan, Li Ma, Xiao You, Wenting Zhou,	1179
		Yunpeng Ye, Wen Liu, Xiangkun Xu, Yongsheng	1180
		Zhang, Tiantian Dong, Jiawei Zhu, Zhe Wang, Xi-	1181
		jian Ju, Jianxun Song, Haoliang Cheng, Xiaojing Li,	1182
		Jiandong Ding, Hefei Guo, and Zhengyong Zhang.	1183

1184	2024b. P/d-serve: Serving disaggregated large language model at scale . <i>Preprint</i> , arXiv:2408.08147.	Technical Conference (USENIX ATC 23), pages 945–959, Boston, MA. USENIX Association.	1239
1185			1240
1186	Yunho Jin, Chun-Feng Wu, David Brooks, and Gu-Yeon Wei. 2023. s³: Increasing gpu utilization during generative inference for higher throughput . <i>Advances in Neural Information Processing Systems</i> , 36:18015–18027.	Jiaxing Li, Chi Xu, Feng Wang, Isaac M von Riedemann, Cong Zhang, and Jiangchuan Liu. 2024c. Scalm: Towards semantic caching for automated chat services with large language models . <i>Preprint</i> , arXiv:2406.00025.	1241
1187			1242
1188			1243
1189			1244
1190			1245
1191	Kostis Kaffes, Timothy Chong, Jack Tigar Humphries, Adam Belay, David Mazières, and Christos Kozyrakis. 2019. Shinjuku: Preemptive scheduling for μsecond-scale tail latency . In <i>16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)</i> , pages 345–360, Boston, MA. USENIX Association.	Luchang Li, Sheng Qian, Jie Lu, Lunxi Yuan, Rui Wang, and Qin Xie. 2024d. Transformer-lite: High-efficiency deployment of large language models on mobile phone gpus . <i>Preprint</i> , arXiv:2403.20041.	1246
1192			1247
1193			1248
1194			1249
1195			
1196		Ning Li, Song Guo, Tuo Zhang, Muqing Li, Zicong Hong, Qihua Zhou, Xin Yuan, and Haijun Zhang. 2025a. The moe-empowered edge llms deployment: Architecture, challenges, and opportunities . <i>Preprint</i> , arXiv:2502.08381.	1250
1197			1251
1198	Kyoungmin Kim, Kijae Hong, Caglar Gulcehre, and Anastasia Ailamaki. 2024. The effect of scheduling and preemption on the efficiency of llm inference serving . <i>Preprint</i> , arXiv:2411.07447.		1252
1199			1253
1200			1254
1201			
1202	Vijay Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. 2022. Reducing activation recomputation in large transformer models . <i>Preprint</i> , arXiv:2205.05198.	Suyi Li, Hanfeng Lu, Tianyuan Wu, Minchen Yu, Qizhen Weng, Xusheng Chen, Yizhou Shan, Binhang Yuan, and Wei Wang. 2024e. Caraserve: Cpu-assisted and rank-aware lora serving for generative llm inference . <i>Preprint</i> , arXiv:2401.11240.	1255
1203			1256
1204			1257
1205			1258
1206			1259
1207	Ferdi Kossmann, Bruce Fontaine, Daya Khudia, Michael Cafarella, and Samuel Madden. 2024. Is the gpu half-empty or half-full? practical scheduling techniques for llms . <i>Preprint</i> , arXiv:2410.17840.	Wei Qing Li, Guochao Jiang, Xiangyong Ding, Zhangcheng Tao, Chuzhan Hao, Chenfeng Xu, Yuewei Zhang, and Hao Wang. 2025b. Flowkv: A disaggregated inference framework with low-latency kv cache transfer and load-aware scheduling . <i>Preprint</i> , arXiv:2504.03775.	1260
1208			1261
1209			1262
1210			1263
1211	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention . <i>Preprint</i> , arXiv:2309.06180.	Yueying Li, Jim Dai, and Tianyi Peng. 2025c. Throughput-optimal scheduling algorithms for llm inference and ai agents . <i>Preprint</i> , arXiv:2504.07347.	1264
1212			1265
1213			
1214			1266
1215			1267
1216			1268
1217	Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. 2024. Infinigen: Efficient generative inference of large language models with dynamic kv cache management . <i>Preprint</i> , arXiv:2406.19707.	Zhuohan Li, Lianmin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng, Xin Jin, Yanping Huang, Zhifeng Chen, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023b. Alpaserve: Statistical multiplexing with model parallelism for deep learning serving . <i>Preprint</i> , arXiv:2302.11665.	1269
1218			1270
1219			1271
1220			1272
1221	Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. Gshard: Scaling giant models with conditional computation and automatic sharding . <i>Preprint</i> , arXiv:2006.16668.	Zikun Li, Zhuofu Chen, Remi Delacourt, Gabriele Oliaro, Zeyu Wang, Qinghan Chen, Shuhuai Lin, April Yang, Zhihao Zhang, Zhuoming Chen, Sean Lai, Xupeng Miao, and Zhihao Jia. 2025d. Adaserve: Slo-customized llm serving with fine-grained speculative decoding . <i>Preprint</i> , arXiv:2501.12162.	1273
1222			1274
1223			1275
1224			1276
1225			1277
1226			1278
1227	Baolin Li, Yankai Jiang, Vijay Gadepally, and Devesh Tiwari. 2024a. Llm inference serving: Survey of recent advances and opportunities . <i>arXiv preprint arXiv:2407.12391</i> .	Bin Lin, Chen Zhang, Tao Peng, Hanyu Zhao, Wencong Xiao, Minmin Sun, Anmin Liu, Zhipeng Zhang, Lanbo Li, Xiafei Qiu, Shen Li, Zhigang Ji, Tao Xie, Yong Li, and Wei Lin. 2024a. Infinite-llm: Efficient llm service for long context with distattention and distributed kvcache . <i>Preprint</i> , arXiv:2401.02669.	1279
1228			1280
1229			1281
1230			1282
1231	Jialong Li, Shreyansh Tripathi, Lakshay Rastogi, Yiming Lei, Rui Pan, and Yiting Xia. 2024b. Optimizing mixture-of-experts inference time combining model deployment and communication scheduling . <i>Preprint</i> , arXiv:2410.17043.		1283
1232			1284
1233			1285
1234			1286
1235			
1236	Jiamin Li, Yimin Jiang, Yibo Zhu, Cong Wang, and Hong Xu. 2023a. Accelerating distributed MoE training and inference with lina . In <i>2023 USENIX Annual</i>	Chaofan Lin, Zhenhua Han, Chengruidong Zhang, Yuqing Yang, Fan Yang, Chen Chen, and Lili Qiu. 2024b. Parrot: Efficient serving of llm-based applications with semantic variable . <i>Preprint</i> , arXiv:2405.19888.	1287
1237			1288
1238			1289
			1290
			1291

1292	Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024c.	1347
1293	Awq: Activation-aware weight quantization for llm compression and acceleration. In <i>MLSys</i> .	1348
1294		1349
1295		1350
1296		1351
1297	Xue Lin, Zhibo Zhang, Peining Yue, Haoran Li, Jin Zhang, Baoyu Fan, Huayou Su, and Xiaoli Gong. 2024d.	1352
1298	Syncintellects: Orchestrating llm inference with progressive prediction and qos-friendly control .	1353
1299	In <i>2024 IEEE/ACM 32nd International Symposium on Quality of Service (IWQoS)</i> , pages 1–10.	
1300		
1301	Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. 2024e.	
1302	Qserve: W4a8kv4 quantization and system co-design for efficient llm serving . <i>Preprint</i> , arXiv:2405.04532.	
1303		
1304		
1305		
1306		
1307		
1308	Akide Liu, Jing Liu, Zizheng Pan, Yefei He, Gholamreza Haffari, and Bohan Zhuang. 2024a.	
1309	Minicache: Kv cache compression in depth dimension for large language models . <i>Preprint</i> , arXiv:2405.14366.	
1310		
1311		
1312	Hao Liu, Matei Zaharia, and Pieter Abbeel. 2023.	
1313	Ring attention with blockwise transformers for near-infinite context . <i>Preprint</i> , arXiv:2310.01889.	
1314		
1315	Jiacheng Liu, Peng Tang, Wenfeng Wang, Yuhang Ren, Xiaofeng Hou, Pheng-Ann Heng, Minyi Guo, and Chao Li. 2024b.	
1316	A survey on inference optimization techniques for mixture of experts models . <i>Preprint</i> , arXiv:2412.14219.	
1317		
1318		
1319		
1320	Yuhan Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, Michael Maire, Henry Hoffmann, Ari Holtzman, and Junchen Jiang. 2024c.	
1321	Cachegen: Kv cache compression and streaming for fast large language model serving . <i>Preprint</i> , arXiv:2310.07240.	
1322		
1323		
1324		
1325		
1326		
1327	Michael Luo, Xiaoxiang Shi, Colin Cai, Tianjun Zhang, Justin Wong, Yichuan Wang, Chi Wang, Yanping Huang, Zhifeng Chen, Joseph E. Gonzalez, and Ion Stoica. 2025.	
1328	Autellix: An efficient serving engine for llm agents as general programs . <i>Preprint</i> , arXiv:2502.13965.	
1329		
1330		
1331		
1332		
1333	Yixuan Mei, Yonghao Zhuang, Xupeng Miao, Juncheng Yang, Zhihao Jia, and Rashmi Vinayak. 2024.	
1334	Helix: Distributed serving of large language models via max-flow on heterogeneous gpus . <i>Preprint</i> , arXiv:2406.01566.	
1335		
1336		
1337		
1338	Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Hongyi Jin, Tianqi Chen, and Zhihao Jia. 2023.	
1339	Towards efficient generative large language model serving: A survey from algorithms to systems. <i>arXiv preprint arXiv:2312.15234</i> .	
1340		
1341		
1342		
1343	Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, Chun-shi, Zhuoming Chen, Daiyaan Arfeen, Reyna	
1344	Abhyankar, and Zhihao Jia. 2024a.	
1345	Specinfer: Accelerating large language model serving with tree-based speculative inference and verification . In <i>Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3</i> , ASPLOS '24, page 932–949. ACM.	1347
1346		1348
		1349
		1350
		1351
		1352
		1353
	Xupeng Miao, Chunan Shi, Jiangfei Duan, Xiaoli Xi, Dahua Lin, Bin Cui, and Zhihao Jia. 2024b.	1354
	Spotserve: Serving generative large language models on preemptible instances . In <i>Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2</i> , pages 1112–1127.	1355
		1356
		1357
		1358
		1359
		1360
	Microsoft. 2022.	1361
	DeepSpeed-mii .	
	Mistral. 2025.	1362
	Mistral-small-24b-instruct-2501 .	
	moonshot. 2023.	1363
	kimi .	
	Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Anand Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar Phanishayee, and Matei Zaharia. 2021.	1364
	Efficient large-scale language model training on gpu clusters using megatron-llm . <i>Preprint</i> , arXiv:2104.04473.	1365
		1366
		1367
		1368
		1369
		1370
	Sophia Nguyen, Beihao Zhou, and YD Liu. 2024.	1371
	S. towards sustainable large language model serving . In <i>ACM HotCarbon Workshop on Sustainable Computer Systems</i> .	1372
		1373
		1374
	NVIDIA. 2023.	1375
	Tensorrt-llm .	
	NVIDIA. 2024.	1376
	Context parallelism overview .	
	NVIDIA. 2025.	1377
	Multi-instance gpu .	
	OpenAI. 2024.	1378
	Introducing openai o1-preview .	
	Rui Pan, Zhuang Wang, Zhen Jia, Can Karakus, Luca Zancato, Tri Dao, Yida Wang, and Ravi Netravali. 2024a.	1379
	Marconi: Prefix caching for the era of hybrid llms . <i>Preprint</i> , arXiv:2411.19379.	1380
		1381
		1382
	Xiurui Pan, Endian Li, Qiao Li, Shengwen Liang, Yizhou Shan, Ke Zhou, Yingwei Luo, Xiaolin Wang, and Jie Zhang. 2024b.	1383
	Instinfer: In-storage attention offloading for cost-effective long-context llm inference . <i>Preprint</i> , arXiv:2409.04992.	1384
		1385
		1386
		1387
	Daon Park and Bernhard Egger. 2024.	1388
	Improving throughput-oriented llm inference with cpu computations . In <i>Proceedings of the 2024 International Conference on Parallel Architectures and Compilation Techniques</i> , PACT '24, page 233–245, New York, NY, USA. Association for Computing Machinery.	1389
		1390
		1391
		1392
		1393
	Pratyush Patel, Esha Choukse, Chaojie Zhang, Íñigo Goiri, Brijesh Warriar, Nithish Mahalingam, and Ricardo Bianchini. 2024a.	1394
	Characterizing power management opportunities for llms in the cloud . In	1395
		1396
		1397

1507	Mohammad Shoeybi, Mostofa Patwary, Raul Puri,	serving . In <i>18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)</i> ,	1561
1508	Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2020. Megatron-lm: Training multi-billion parameter language models using model parallelism .	pages 911–927, Santa Clara, CA. USENIX Association.	1562
1509	<i>Preprint</i> , arXiv:1909.08053.		1563
1510			1564
1511			
1512	Yixin Song, Zeyu Mi, Haotong Xie, and Haibo Chen.	Hui Wu, Yi Gan, Feng Yuan, Jing Ma, Wei Zhu, Yutao Xu, Hong Zhu, Yuhua Zhu, Xiaoli Liu, Jinghui Gu, and Peng Zhao. 2024d. Efficient llm inference solution on intel gpu . <i>Preprint</i> , arXiv:2401.05391.	1565
1513	2024. Powerinfer: Fast large language model serving with a consumer-grade gpu . <i>Preprint</i> , arXiv:2312.12456.		1566
1514			1567
1515			1568
1516	Vikranth Srivatsa, Zijian He, Reyna Abhyankar, Dongming Li, and Yiyang Zhang. 2024. Preble: Efficient distributed prompt scheduling for llm serving .	Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang, Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhi-fang Sui. 2024. Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding . <i>Preprint</i> , arXiv:2401.07851.	1569
1517	<i>Preprint</i> , arXiv:2407.00023.		1570
1518			1571
1519			1572
1520	Jovan Stojkovic, Chaojie Zhang, Íñigo Goiri, Josep Torrellas, and Esha Choukse. 2024. Dynamollm: Designing llm inference clusters for performance and energy efficiency . <i>arXiv preprint arXiv:2408.00741</i> .	Yi Xiong, Hao Wu, Changxu Shao, Ziqing Wang, Rui Zhang, Yuhong Guo, Junping Zhao, Ke Zhang, and Zhenxuan Pan. 2024. Layerkv: Optimizing large language model serving with layer-wise kv cache management . <i>Preprint</i> , arXiv:2410.00428.	1574
1521			1575
1522			1576
1523			1577
1524	Foteini Strati, Sara Mcallister, Amar Phanishayee, Jakub Tarnawski, and Ana Klimovic. 2024. Déjàvu: Kv-cache streaming for fast, fault-tolerant generative llm serving . <i>Preprint</i> , arXiv:2403.01876.	Daliang Xu, Hao Zhang, Liming Yang, Ruiqi Liu, Gang Huang, Mengwei Xu, and Xuanzhe Liu. 2024. Fast on-device llm inference with npus . <i>Preprint</i> , arXiv:2407.05858.	1579
1525			1580
1526			1581
1527			1582
1528	Yi Su, Yuechi Zhou, Quantong Qiu, Juntao Li, Qingrong Xia, Ping Li, Xinyu Duan, Zhefeng Wang, and Min Zhang. 2025. Accurate kv cache quantization with outlier tokens tracing . <i>Preprint</i> , arXiv:2505.10938.	Huan Yang, Deyu Zhang, Yudong Zhao, Yuanchun Li, and Yunxin Liu. 2024a. A first look at efficient and secure on-device llm inference against kv leakage . <i>Preprint</i> , arXiv:2409.04040.	1583
1529			1584
1530			1585
1531			1586
1532	Biao Sun, Ziming Huang, Hanyu Zhao, Wencong Xiao, Xinyi Zhang, Yong Li, and Wei Lin. 2024. Llumnix: Dynamic scheduling for large language model serving . <i>Preprint</i> , arXiv:2406.03243.	Zheming Yang, Yuanhao Yang, Chang Zhao, Qi Guo, Wenkai He, and Wen Ji. 2024b. Perllm: Personalized inference scheduling with edge-cloud collaboration for diverse llm services . <i>Preprint</i> , arXiv:2405.14636.	1587
1533			1588
1534			1589
1535			1590
1536	Xin Tan, Yimin Jiang, Yitao Yang, and Hong Xu. 2024. Teola: Towards end-to-end optimization of llm-based applications . <i>Preprint</i> , arXiv:2407.00326.	Jiayi Yao, Hanchen Li, Yuhan Liu, Siddhant Ray, Yihua Cheng, Qizheng Zhang, Kuntai Du, Shan Lu, and Junchen Jiang. 2024. Cacheblend: Fast large language model serving for rag with cached knowledge fusion . <i>Preprint</i> , arXiv:2405.16444.	1591
1537			1592
1538			1593
1539	Jikai Wang, Yi Su, Juntao Li, Qingrong Xia, Zi Ye, Xinyu Duan, Zhefeng Wang, and Min Zhang. 2024a. Opt-tree: Speculative decoding with adaptive draft tree structure . <i>Preprint</i> , arXiv:2406.17276.	Wangsong Yin, Mengwei Xu, Yuanchun Li, and Xuanzhe Liu. 2024. Llm as a system service on mobile devices . <i>Preprint</i> , arXiv:2403.11805.	1596
1540			1597
1541			1598
1542			
1543	Wenxiao Wang, Wei Chen, Yicong Luo, Yongliu Long, Zhengkai Lin, Liye Zhang, Binbin Lin, Deng Cai, and Xiaofei He. 2024b. Model compression and efficient inference for large language models: A survey . <i>Preprint</i> , arXiv:2402.09748.	Chengye Yu, Tianyu Wang, Zili Shao, Linjie Zhu, Xu Zhou, and Song Jiang. 2024. Twinpilots: A new computing paradigm for gpu-cpu parallel llm inference . In <i>Proceedings of the 17th ACM International Systems and Storage Conference, SYSTOR '24</i> , page 91–103, New York, NY, USA. Association for Computing Machinery.	1599
1544			1600
1545			1601
1546			1602
1547			1603
1548	Bingyang Wu, Shengyu Liu, Yinmin Zhong, Peng Sun, Xuanzhe Liu, and Xin Jin. 2024a. Loongserve: Efficiently serving long-context large language models with elastic sequence parallelism . <i>Preprint</i> , arXiv:2404.09526.	Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. Orca: A distributed serving system for Transformer-Based generative models . In <i>16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)</i> , pages 521–538, Carlsbad, CA. USENIX Association.	1606
1549			1607
1550			1608
1551			1609
1552			1610
1553	Bingyang Wu, Yinmin Zhong, Zili Zhang, Shengyu Liu, Fangyue Liu, Yuanhang Sun, Gang Huang, Xuanzhe Liu, and Xin Jin. 2024b. Fast distributed inference serving for large language models . <i>Preprint</i> , arXiv:2305.05920.	Hanfei Yu, Xingqi Cui, Hong Zhang, Hao Wang, and Hao Wang. 2025a. fmoe: Fine-grained expert offloading for large mixture-of-experts serving . <i>Preprint</i> , arXiv:2502.05370.	1613
1554			1614
1555			1615
1556			1616
1557			
1558	Bingyang Wu, Ruidong Zhu, Zili Zhang, Peng Sun, Xuanzhe Liu, and Xin Jin. 2024c. dLoRA: Dynamically orchestrating requests and adapters for LoRA LLM		
1559			
1560			

1617	Shan Yu, Jiarong Xing, Yifan Qiao, Mingyuan Ma,	Zhen Zheng, Xin Ji, Taosong Fang, Fanghao Zhou,	1670
1618	Yangmin Li, Yang Wang, Shuo Yang, Zhiqiang Xie,	Chuanjie Liu, and Gang Peng. 2024c. Batchllm:	1671
1619	Shiyi Cao, Ke Bao, Ion Stoica, Harry Xu, and	Optimizing large batched llm inference with global	1672
1620	Ying Sheng. 2025b. Prism: Unleashing gpu shar-	prefix sharing and throughput-oriented token batch-	1673
1621	ing for cost-efficient multi-llm serving. <i>Preprint,</i>	ing. <i>Preprint,</i> arXiv:2412.03594.	1674
1622	arXiv:2505.04021.		
1623	Zhihang Yuan, Yuzhang Shang, Yang Zhou, Zhen Dong,	Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu,	1675
1624	Zhe Zhou, Chenhao Xue, Bingzhe Wu, Zhikai Li,	Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang.	1676
1625	Qingyi Gu, Yong Jae Lee, et al. 2024. Llm inference	2024. Distserve: Disaggregating prefill and decoding	1677
1626	unveiled: Survey and roofline model insights. <i>arXiv</i>	for goodput-optimized large language model serving.	1678
1627	<i>preprint arXiv:2402.16363.</i>	<i>arXiv preprint arXiv:2401.09670.</i>	1679
1628	Hailin Zhang, Xiaodong Ji, Yilin Chen, Fangcheng Fu,	Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yan-	1680
1629	Xupeng Miao, Xiaonan Nie, Weipeng Chen, and Bin	ping Huang, Vincent Zhao, Andrew Dai, Zhifeng	1681
1630	Cui. 2024a. Pqcache: Product quantization-based	Chen, Quoc Le, and James Laudon. 2022. Mixture-	1682
1631	kvcache for long context llm inference. <i>Preprint,</i>	of-experts with expert choice routing. <i>Preprint,</i>	1683
1632	arXiv:2407.12820.	arXiv:2202.09368.	1684
1633	Mingjin Zhang, Jiannong Cao, Xiaoming Shen, and	Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Ji-	1685
1634	Zeyang Cui. 2024b. Edgeshard: Efficient llm in-	aming Xu, Shiyao Li, Yuming Lou, Luning Wang,	1686
1635	ference via collaborative edge computing. <i>Preprint,</i>	Zhihang Yuan, Xiuhong Li, et al. 2024. A survey on	1687
1636	arXiv:2405.14371.	efficient inference for large language models. <i>arXiv</i>	1688
1637	Wei Zhang, Zhiyu Wu, Yi Mu, Banruo Liu, Myungjin	<i>preprint arXiv:2404.14294.</i>	1689
1638	Lee, and Fan Lai. 2025. Tempo: Application-aware	Yun Zhu, Jia-Chen Gu, Caitlin Sikora, Ho Ko, Yinx-	1690
1639	llm serving with mixed slo requirements. <i>Preprint,</i>	iao Liu, Chu-Cheng Lin, Lei Shu, Liangchen Luo,	1691
1640	arXiv:2504.20068.	Lei Meng, Bang Liu, and Jindong Chen. 2024.	1692
1641	Xiaojin Zhang, Yulin Fei, Yan Kang, Wei Chen, Lixin	Accelerating inference of retrieval-augmented gen-	1693
1642	Fan, Hai Jin, and Qiang Yang. 2024c. No free	eration via sparse context selection. <i>Preprint,</i>	1694
1643	lunch theorem for privacy-preserving llm inference.	arXiv:2405.16178.	1695
1644	<i>Preprint,</i> arXiv:2405.20681.	Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong,	1696
1645	Zhihao Zhang, Alan Zhu, Lijie Yang, Yihua Xu, Lant-	Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and	1697
1646	ing Li, Phitchaya Mangpo Phothilimthana, and Zhi-	Xia Hu. 2024. Kivi : Plug-and-play 2bit kv cache	1698
1647	hao Jia. 2024d. Accelerating retrieval-augmented	quantization with streaming asymmetric quantiza-	1699
1648	language model serving with speculation. <i>Preprint,</i>	tion.	1700
1649	arXiv:2401.14021.	Longwei Zou, Tingfeng Liu, Kai Chen, Jiangang Kong,	1701
1650	Juntao Zhao, Borui Wan, Yanghua Peng, Haibin Lin,	and Yangdong Deng. 2024. Instcache: A predictive	1702
1651	and Chuan Wu. 2024a. Llm-pq: Serving llm on het-	cache for llm serving. <i>Preprint,</i> arXiv:2411.13820.	1703
1652	erogeneous clusters with phase-aware partition and	Małgorzata Łazuka, Andreea Anghel, and Thomas Par-	1704
1653	adaptive quantization. <i>Preprint,</i> arXiv:2403.01136.	nell. 2024. Llm-pilot: Characterize and optimize	1705
1654	Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn	performance of your llm inference services. <i>Preprint,</i>	1706
1655	Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy,	arXiv:2410.02425.	1707
1656	Tianqi Chen, and Baris Kasikci. 2024b. Atom: Low-		
1657	bit quantization for efficient and accurate llm serving.		
1658	<i>Preprint,</i> arXiv:2310.19102.		
1659	Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue		
1660	Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos		
1661	Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark		
1662	Barrett, and Ying Sheng. 2024a. Sglang: Efficient		
1663	execution of structured language model programs.		
1664	<i>Preprint,</i> arXiv:2312.07104.		
1665	Zangwei Zheng, Xiaozhe Ren, Fuzhao Xue, Yang		
1666	Luo, Xin Jiang, and Yang You. 2024b. Response		
1667	length perception and sequence scheduling: An llm-		
1668	empowered llm inference pipeline. <i>Advances in Neu-</i>		
1669	<i>ral Information Processing Systems</i> , 36.		

A LLM Inference Serving in Cluster

This section focuses on cluster-level deployment and scheduling, as well as cloud-based cluster serving, as detailed in Figure 11.

A.1 Cluster Optimization

Internal optimizations for homogeneous devices require more machines as parameter scale increases, while heterogeneous machines are preferred for their flexibility, efficiency, and cost-effectiveness (Mei et al., 2024). External optimizations, like service-oriented cluster scheduling, further enhance internal optimizations.

Architecture and Optimization for Heterogeneous Resources. Jayaram Subramanya et al. (2023) propose a joint optimization framework for adaptive task allocation across GPU types and batch sizes, demonstrating significant throughput improvements over static configurations. Helix (Mei et al., 2024) models the execution of LLM services on heterogeneous GPUs and networks as a maximum flow problem in a directed weighted graph, where nodes represent GPU instances and edges encode GPU and network heterogeneity through capacity constraints. LLM-PQ (Zhao et al., 2024a) advocates an adaptive quantization and phase-aware partition scheme tailored for heterogeneous GPU clusters. HexGen (Jiang et al., 2024c) supports asymmetric parallel execution on GPUs with different computing capabilities. Splitwise (Patel et al., 2024b), DistServe (Zhong et al., 2024) and HEXGEN-2 (Jiang et al., 2025b) optimize computation on heterogeneous disaggregated architectures, with the latter focusing on LLM serving via constraint-based scheduling and graph-based resource optimization. Hisaharo et al. (2024) integrate advanced interconnect technology, high-bandwidth memory, and energy-efficient power management.

Service-Aware Scheduling. DynamoLLM (Stojkovic et al., 2024) optimizes service clusters by adjusting instances, parallelization, and GPU frequencies based on input/output lengths. Splitwise (Patel et al., 2024b) proposes cluster-level scheduling across prefill and decoding on separate devices.

A.2 Load Balancing

Cluster-level load balancing optimizes request distribution to prevent node overload or underutilization, improving throughput and service qual-

ity. While most frameworks (Yu et al., 2022; Kwon et al., 2023) rely on traditional methods like Round Robin and Random (Deepspeed, 2023), recent advances in heuristic, dynamic, and predictive scheduling provide more sophisticated solutions.

Heuristic Algorithm. SCLS (Cheng et al., 2024b) employs a max-min algorithm (Radunovic and Le Boudec, 2007) to balance the workloads. It assigns the batch with the longest estimated serving time to the instance with the lowest score, where the score represents the total serve time of all batches in the instance’s queue. SAL (Kossmann et al., 2024) quantifies the load on two key factors: (1) the number of queued prefill tokens and (2) the available memory. This ensures that requests are dispatched to the server with the lowest load, addressing scenarios where delays occur due to either a full token batch or insufficient memory.

Dynamic Scheduling. Llumnix (Sun et al., 2024) dynamically reschedules requests across model instances during runtime to handle request heterogeneity and unpredictability. It uses real-time migration to transfer requests and memory states, enabling mid-operation migration to the least loaded instance based on real-time load growth.

Intelligent Predictive Scheduling. Jain et al. (2024) propose a reinforcement learning-based router that models request routing as a Markov Decision Process, aiming to derive an optimal policy for maximizing discounted rewards. It integrates response length prediction, workload impact estimation, and reinforcement learning.

A.3 Cloud-Based LLM Serving

If local LLM deployment lacks resources, cloud services offer a more economical alternative, with recent research focusing on optimizing cloud deployment and edge collaboration for efficiency.

Deployment and Computing Effective. To reduce LLM deployment costs, spot instances are used despite preemption risks. SpotServe (Miao et al., 2024b) mitigates this with dynamic parallelization, parameter reuse, and stateful inference recovery. ServerlessLLM (Fu et al., 2024a) tackles serverless cold start latency via optimized checkpoints, live migration, and locality-aware scheduling. Mélange (Griggs et al., 2024) optimizes GPU allocation based on request patterns, lowering costs. POLCA (Patel et al., 2024a) boosts efficiency through power management, while Imai

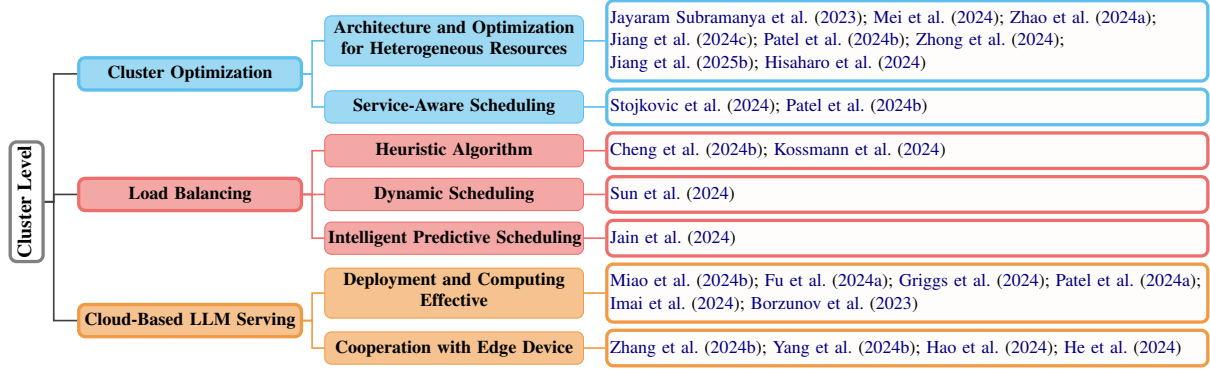


Figure 11: Taxonomy of Cluster-Level strategies for LLM inference serving.

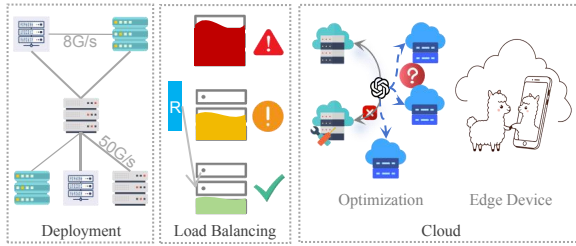


Figure 12: Overview of the deployment, load balancing, and cloud-based in cluster-level optimization. R represents a request.

et al. (2024) predict inference latency to enhance cluster management. Borzunov et al. (2023) propose a way to integrate idle resources through geodistributed devices connected via the internet. Jiang et al. (2025a) found that optimizing GPU combinations, deployment setups, and workload distribution boosts LLM cost efficiency on heterogeneous cloud platforms.

Cooperation with Edge Device. To meet SLOs amid cloud latency and bandwidth limits, edge computing offers solutions. EdgeShard (Zhang et al., 2024b) leverages collaboration between distributed edge devices and cloud servers. PreLLM (Yang et al., 2024b) uses a multi-armed bandit framework for personalized scheduling. Hao et al. (2024) integrate small edge models with cloud LLM to address memory constraints, while He et al. (2024) apply deep reinforcement learning for efficient, latency-aware inference offloading.

B Miscellaneous Areas

Other niche but important directions (Figure 13) are also advancing LLM inference toward a more comprehensive and far-reaching future.

B.1 Hardware

Recent advancements in optimizing LLM inference have focused on improving efficiency, speed, and resource utilization in various hardware techniques.

Peng et al. (2024) propose a mixed-precision, multi-level caching system (HBM, DRAM, SSDs) and a model modularization algorithm to enable LLM inference on resource-constrained, outdated hardware. Wu et al. (2024d) explore inference service solutions on Intel GPUs. LLM-Pilot (Łazuka et al., 2024) benchmarks LLM inference across GPUs and recommends the most cost-effective GPU for unseen LLMs. GenZ (Bambhaniya et al., 2024) is an analytical tool for studying the relationship between LLM inference performance and various hardware platform design parameters. Li et al. (2024d) present Transformer-Lite, an innovative inference engine optimized for mobile GPUs, designed to enhance the efficiency and inference speed of LLM deployment on mobile devices. LLMS (Yin et al., 2024) is an innovative system on mobile devices that, under stringent memory constraints, implements fine-grained, chunk-based KV cache compression and a globally optimized swapping mechanism to decouple applications from LLM memory management, thereby minimizing the overhead of context switching. Xu et al. (2024) utilize on-device Neural Processing Unit (NPU) offloading to enhance NPU offloading efficiency and reduce prefill latency.

B.2 Privacy

Protecting user conversation content in LLMs from potential leakage is an important issue. Yang et al. (2024a) adopt weight permutation to shuffle KV pairs, preventing attackers from reconstructing the entire context. Zhang et al. (2024c) quantify the trade-off between privacy protection and utility



Figure 13: Taxonomy of Miscellaneous Areas for LLM inference serving.

loss, pointing out that privacy protection mechanisms (such as randomization) can reduce privacy leakage but will introduce utility loss. MARILL (Rathee et al., 2024) achieves substantial reductions in the costly operations required for secure inference within multi-party computation by optimizing the architecture of LLMs during the fine-tuning phase.

B.3 Simulator

Considering the diversity of computing devices and their associated high costs, a comprehensive simulator is indispensable for conducting trials in virtual environments. Agrawal et al. (2024b) introduce Vidur, a scalable, high-fidelity simulation framework for evaluating LLM performance under various deployment configurations, alongside Vidur-Search, a tool for optimizing deployments to meet performance constraints and reduce costs. The Helix system (Mei et al., 2024), featuring an event-based simulator, enables accurate simulation of LLM inference in heterogeneous GPU clusters by adjusting factors like network conditions, machine heterogeneity, and cluster scale, providing rapid and cost-effective deployment evaluations.

B.4 Fairness

In LLM inference services, request frequency limits are typically imposed on each client (e.g., user or application) to ensure fair resource allocation. These limits prevent excessive requests from monopolizing resources and degrading service quality for others. However, they may also result in underutilized resources. Sheng et al. (2024) propose a novel fairness definition, based on a cost function considering input and output tokens. Additionally, a new scheduling algorithm, the Virtual Token Counter (VTC), introduces fair scheduling through a continuous batching mechanism.

B.5 Energy

Given the substantial power demands of LLM computations, optimizing energy usage is a critical challenge that must be addressed. Nguyen et al. (2024) investigate the carbon emissions of LLMs from operational and embodied perspectives, aiming to promote sustainable LLM services. Researchers analyzed the performance and energy consumption of the LLaMA model across varying parameter scales and batch sizes, incorporating the carbon intensity of different power grid regions. This study provides insights into the environmental impact of LLMs and explores opportunities to optimize sustainable LLM systems.

C Common Inference Framework

We present a comprehensive comparison of vLLM (Kwon et al., 2023), TensorRT-LLM (NVIDIA, 2023), SGLang (Zheng et al., 2024a), TGI (HuggingFace, 2023), DeepSpeed-MII (Microsoft, 2022), and llama.cpp (ggml org, 2022) in terms of hardware support, core optimization focus, and typical use cases (as shown in Table 3).

Metric	Definition	Key Focus
TTFT	Latency from input to first token.	Critical for real-time apps (e.g., chatbots).
TBT	Time interval between consecutive tokens.	Reflects step-by-step responsiveness.
TPOT	Average time per token during decoding.	Measures token generation efficiency.
Throughput	Tokens generated per second across all requests.	Evaluates system capacity under high load.
Capacity	Maximum throughput while meeting SLOs.	Represents system's upper performance limit.
Normalized Latency	Total execution time divided by token count.	Holistic view of system efficiency.
Percentile Metrics	Latency distribution (e.g., P50, P90, P99).	Evaluates stability and performance bounds.

Table 2: LLM Inference Service Evaluation Metrics.

Tool	Hardware Support	Core Optimization Focus	Typical Use Case
vLLM	GPU	PagedAttention, optimized KV cache, high concurrency	Multi-turn chat, OpenChat, FastChat, scalable SaaS
TensorRT-LLM	NVIDIA GPU	Low latency, high throughput	Enterprise-grade GPU server deployment
SGLang	GPU / CPU	Structured prompt execution optimization	Agent-based reasoning, complex logic inference
TGI	GPU	Continuous batching, cloud integration	HuggingFace model-as-a-service (SaaS)
DeepSpeed-MII	GPU	Memory compression, distributed inference	High-throughput API serving
llama.cpp	CPU / Low-end GPU	Quantization, lightweight deployment	Local development / edge computing

Table 3: Comparison of LLM Inference Frameworks.

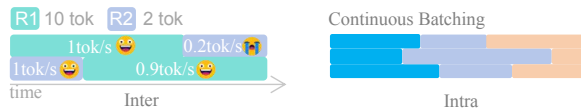


Figure 14: An example of inter- and intra-scheduling. **R** represents a request. In Inter-request scheduling, two requests, **R1** (10 toks) and **R2** (2 toks), arrive simultaneously. Ignoring the prefill process, if **R1** is processed first, its generation rate is 1 tok/s, and **R2**'s rate is 0.2 tok/s. Reversing the order gives **R2** a rate of 1 tok/s and **R1** 0.9 tok/s. The default decoding speed is 1 token/s.