# Serving Heterogeneous LoRA Adapters in Distributed LLM Inference Systems

Shashwat Jaiswal[1], Shrikara Arun[2], Anjaly Parayil[2], Ankur Mallick[2], Spyros Mastorakis[2], Alind Khare[2],
Chloi Alverti[3], Renee St Amant[2], Chetan Bansal[2], Victor Rühle[2], Josep Torrellas[1]
{sj74, torrella}@illinois.edu, {t-sarun, aparayil, ankurmallick, smastorakis,
alindkhare, reneestamant, chetanb, virueh}@microsoft.com, xalverti@cslab.ece.ntua.gr
[1]University of Illinois Urbana-Champaign, [2]Microsoft, [3]National Technical University of Athens

*Abstract*—**Low-Rank Adaptation (LoRA) has become the de facto method for parameter-efficient fine-tuning of large language models (LLMs), enabling rapid adaptation to diverse domains. In production, LoRA-based models are served at scale, creating multi-tenant environments with hundreds of adapters sharing a base model. However, state-of-the-art serving systems co-batch heterogeneous adapters without accounting for rank (size) variability, leading to severe performance skew, which ultimately requires adding more GPUs to satisfy service-level objectives (SLOs). Existing optimizations, focused on loading, caching, and kernel execution, ignore this heterogeneity, leaving GPU resources underutilized. We present LORASERVE, a workload-aware dynamic adapter placement and routing framework designed to tame rank diversity in LoRA serving. By dynamically rebalancing adapters across GPUs and leveraging GPU Direct RDMA for remote access, LORASERVE maximizes throughput and minimizes tail latency under real-world workload drift. Evaluations on production traces from Company X show that LORASERVE elicits up to 2× higher throughput, up to 9× lower TTFT, while using up to 50% fewer GPUs under SLO constraints compared to state-of-the-art systems.**

## I. INTRODUCTION

Large language models (LLMs) are reshaping modern applications. Models like GPT [38] and Llama [23], pre-trained on vast corpora, exhibit strong general-purpose capabilities. These pre-trained (base) LLMs are increasingly fine-tuned for specific domains to optimize particular use cases, for instance, adapting Llama-2 for improved code generation [40]. Emerging LLM platforms [1]–[4] offer fine-tuning APIs and services that let developers customize models and build domain-specific applications. For example, OpenAI provides APIs to fine-tune GPT-4 and a Completions API to serve those fine-tuned models [1].

Low-Rank Adaptation (LoRA) [21], [25] has emerged as a widely used parameter-efficient fine-tuning method for LLMs [47]. Instead of updating all model weights, LoRA learns low-dimensional updates captured by pairs of small matrices (called LoRA adapters) while the base model remains frozen. Practically, this method of LLM fine-tuning involves training only a LoRA adapter for the target domain, thereby significantly reducing the parameter count and compute requirement. For example, compared to fully fine-tuning GPT-3 (175B), LoRA reduces trainable parameters by over two orders of magnitude and lowers GPU usage by approximately 3×, while maintaining comparable model quality [25].

While LoRA is cost-effective in training, its adapters introduce significant challenges at inference, specifically when multiple LoRA adapters need to be served off the same base model instance. Typical LoRA serving systems maintain thousands of LoRA adapters on inference servers, either in CPU or GPU memory, offering APIs for end users to access specialized LLM inference as per application requirements [41]. The inference server multiplexes different adapters to enable such multi-tenant LoRA serving using specialized CUDA kernels [18], [41] to co-batch requests to different LoRA adapters and execute them together. Consequently, adapters served on the same LLM instance experience performance interference; as shown in Fig 1, requests to lower-rank adapters incur higher latency when co-served with larger-rank requests.

Using these custom CUDA kernels [18], [41] along with hundreds of adapters on a cluster of LLM inference servers is the state-of-the-art and is used at Company X and other cloud providers serving multiple LoRA adapters as shown in Fig 2. This entails two challenges, namely, *managing the inherent scale and heterogeneity (adapter size, traffic etc.) of such workloads* and *management of hundreds of LoRA adapters at a cluster level*.
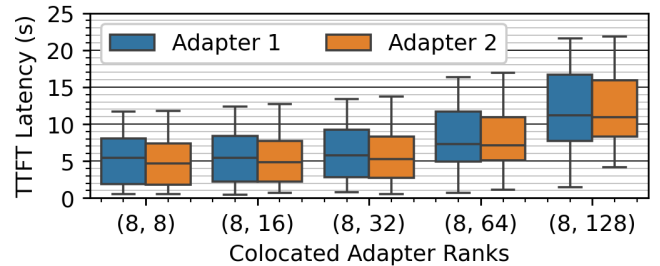


Fig. 1. ***Performance per adapter when two adapters are co-served from the same Llama 7B instance.*** *Higher rank heterogeneity leads to greater variability.*

*Adapter and Workload Heterogeneity.* Each adapter has a different size (rank) and popularity (demand traffic). Prior work has shown that co-batching adapters of different ranks causes the smaller adapters to suffer in performance [33]. Moreover, different adapters have different traffic (ingest rate, bursts etc.) depending on the application they serve, further compounding the heterogeneity. For example, Fig 1 shows that co-serving two adapters of rank-8 and rank-128 from the same Llama 7B LLM inference server causes the P95

1

TTFT (time-to-first-token) of the rank-8 requests to increase by 84% when compared to a server serving the same amount of purely rank-8 requests. This latency increase forces systems to throttle request rates to meet SLOs (Service Level Objectives), reducing overall throughput and GPU efficiency. However, state of the art serving systems like vLLM [5] and S-LoRA [41] remain oblivious to rank-induced heterogeneity, leading to performance degradation under real-workloads.
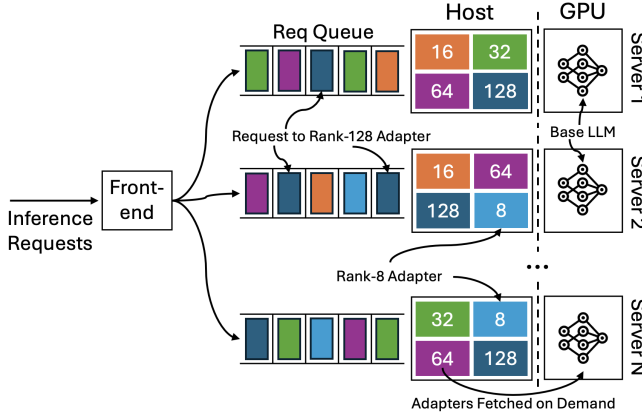


Fig. 2. *A conventional LLM cluster serving LoRA workloads.*

*Managing Hundreds of LoRA Adapters.* Company X uses hundreds of LoRA adapters to support a broad suite of services across hundreds of products used by millions of customers. Large cloud providers also store and serve their customers' LoRA Adapters as a service [6]. These are deployed using a cluster of LLM inference servers each running one base LLM instance and storing a subset of adapters in their CPU memory to be fetched on demand as shown in Fig 2. Each of these adapters see highly fluctuating request demand and statically placing a subset of adapters on each LLM server leads to load imbalance and SLO violations. Another way is to replicate all adapters in the CPU memory of all servers. This is used in state-of-the-art [33] to enable flexible request routing, but this significantly sacrifices CPU memory capacity on every server and may even be impossible at scale due to limited CPU memory. Considering an average general-purpose LLM has approximately 200B parameters [7] and is quantized in 8-bits [20], it should take up approximately 200GB of space. LoRA adapters are usually around 1% of model size [18], [25] which comes to 2GB per adapter and 1TB for five hundred adapters in this case. Storing such high amount of data in every server's CPU memory is wasteful especially since CPU memory is a precious resource for LLM inference scenarios involving long contexts [37] and retrieval-augmented-generation [22].

LORASERVE addresses these performance degradation and memory pressure challenges by managing which adapters are served on which LLM inference server in the cluster. To do this, it takes into account adapter ranks and their request demands to ensure minimal heterogeneity within the adapters being served on a server. Moreover, it places locally only the adapters *actually needed* by the server to alleviate memory pressure. The adapters are distributed across the cluster and

they are fetched on demand when needed.

Designing a system to address these issues efficiently has two technical challenges. The first is the determination of the optimal adapter placement, i.e., which adapter should be served on which (and how many) LLM servers. The second is keeping track of all adapters in the cluster and using this information when a necessary adapter is not present locally in a server. We propose a novel cluster orchestrator, LORASERVE, tailored for LoRA serving to address these two problems. To address the first problem, we propose a dynamic adapter placement technique that determines the ideal set of adapters to be placed on a particular LLM server taking into account the adapter rank and its request arrival history. This ensures load balancing, minimizes interference, and improves resource efficiency. To address the second problem, we propose simple indexing schemes to maintain a logical view of all the adapters in the cluster on every server and leverage GPUDirect RDMA over InfiniBand [8] interconnects to access them when necessary.

In summary, we make the following contributions:

- We analyze interference between adapters of diverse ranks, analyzing impact across multiple dimensions like deployment characteristics (e.g., parallelism), model size, and input size, and reveal implications of adapter rank heterogeneity on system throughput and SLO compliance.
- Further, we characterize real-world LoRA serving workloads using production traces from Company X, uncovering critical insights such as heavy-tailed request distributions, skewed adapter popularity across regions, and distinct arrival patterns.
- We design LORASERVE, a rank-aware, workload-adaptive framework for dynamic adapter placement and routing. LORASERVE dynamically rebalances adapters across LLM inference servers to mitigate rank-induced interference and optimize resource utilization under workload drift.
- We evaluate LORASERVE on production workloads from Company X and open source production traces [9], demonstrating up to $2\times$ higher throughput, up to $9\times$ lower TTFT, while using up to 50% fewer GPUs to meet SLOs compared to state-of-the-art systems. We will release our code and traces publicly after acceptance.

## II. BACKGROUND

### A. Large Language Models

*1) LLM Inference:* LLM inference generates output tokens from a given prompt in two stages: prefill and decode. In the prefill stage, the model reads the entire prompt, produces the first output token, and builds a key–value (KV) cache for each token. The decode stage then iteratively uses this KV cache to produce subsequent tokens, updating the cache as it generates each output token. Decoding stops once a termination condition is reached, typically when an end-of-sequence ($< eos >$) token is generated.

*2) Parameter Efficient Fine-tuning:* Finetuning a model is a standard process which involves training a pre-trained model on specialized datasets to make them adapt to domain specific tasks. Parameter-Efficient Fine-Tuning (PEFT) is a method which achieves comparable results by freezing most of the original parameters and training only a small subset of new or existing parameters [47]. This approach significantly reduces computational and memory requirements compared to full fine-tuning, making it more accessible and cost-effective. PEFT techniques like LoRA [25] and QLoRA [21] allow models to achieve performance comparable to full fine-tuning for specific downstream tasks.

*3) Low-Rank Adaptation:* Low-Rank Adaptation fine-tunes a base model by inserting trainable low-rank matrices into its layers. The key driver of performance is the adapter size determined by the *rank* of those matrices. Using a higher rank can capture richer updates and often improves accuracy. As tasks differ in complexity and their accuracy needs, application developers typically choose adapters with different ranks for different tasks while sharing the base model [41], [44].

Given that adapters are typically $< 1\%$ of their base model in size [18], [25], this technique also reduces the memory requirements of online systems catering to diverse tasks. Moreover, recent works [18], [41] also enable executing multiple adapters in the same batch further improving throughput.

### B. LoRA Serving Systems

*1) Kernel Optimizations:* Research on kernel-level optimizations for LoRA adapter serving in LLM inference has largely focused on overcoming the batching and memory inefficiencies that arise when many requests share the same base model but request different adapters. Early LoRA serving approaches bound each request to a separate kernel invocation, which destroyed batch efficiency and led to GPU underutilization. Recent systems such as Punica [18] and S-LoRA [41] address this by redesigning the compute path. Punica introduces a segmented-gather GEMV (SGMV) kernel that fuses heterogeneous LoRA deltas into a single matmul, allowing the GPU to process many different adapters in one batched operation. This makes multi-tenant LoRA serving compute-efficient even when adapters and adapter ranks differ.

*2) System Optimizations:* S-LoRA [41] pairs custom heterogeneous LoRA kernels with unified paging, a GPU-aware memory virtualization layer that stores a large number of adapters in CPU memory and pages only small active slices to the GPU alongside KV-cache blocks. This unified allocator reduces fragmentation and enables thousands of adapters to coexist without increasing GPU memory footprint. Toppings [33] simultaneously uses CPUs to compute the low-rank adaptation for prefilling as the requested LoRA adapter is being loaded onto GPUs, to improve the TTFT. It also proposes a rank-aware scheduling policy for load balancing requests but assumes that all adapters are replicated on each server in the cluster. dLoRA [45] improves inference efficiency by dynamically merging adapters with the base model and migrating requests and adapters between different worker replicas.

Chameleon [26] introduces adapter caching in GPU memory to prevent unnecessary loading of adapters from system memory. None of these systems take adapter rank heterogeneity into consideration while making scheduling decisions, leading to performance degradation due to interference between requests of different ranks.

### III. MOTIVATION

#### A. Performance

In this section, we quantify the interference between adapters of diverse ranks, analyzing its impact across multiple dimensions, including parallelism, model size, and input size, and reveal implications of adapter rank heterogeneity on system throughput and SLO compliance.

*1) Input Sizes:* Serving adapters of larger ranks involves loading larger chunks from memory and computation on larger matrices, which naturally makes it slower than serving smaller ranks. Fig 3 shows these trends for a single request running in isolation. We observe that a rank-128 adapter takes $2.7\times$ the prefill time of a rank-8 adapter to process a prompt of size 2000. For decodes, the effect is subtle as the operation is largely memory bound. However, we still see a steady drop in performance with increase in input size. As input size grows, the performance impact of adapter rank becomes more pronounced because LoRA is applied to the Q, K, V, and O projection layers, whose compute cost grows with sequence length [21].
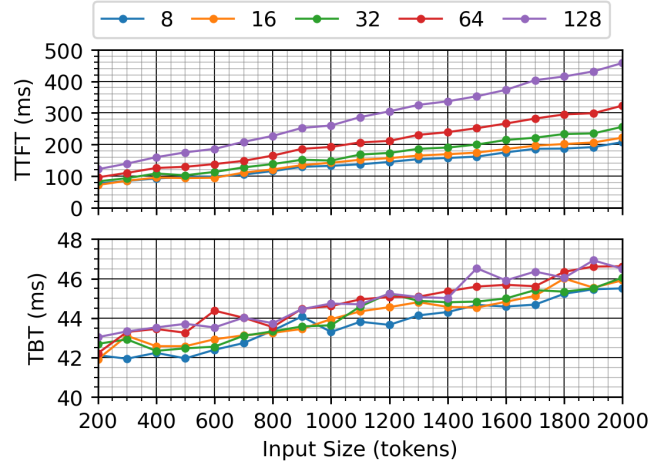


Fig. 3. *Time-To-First-Token (TTFT) (top) and Time-Between-Tokens (TBT) (bottom) of different ranks against Input Size for Llama 7B.*

*2) Model Sizes:* As model size increases, the impact of rank heterogeneity grows significantly, as illustrated in Fig 4. Larger models introduce higher computational complexity, with wider hidden dimensions, deeper transformer stacks, and greater memory bandwidth demands. Consequently, LoRA's low-rank matrices become larger and are injected into more projection layers, amplifying both compute and data-movement overhead. This increased pressure on the system magnifies the performance penalty from rank heterogeneity, reaching up to a 45% degradation on Llama 70B.

*3) Tensor Parallelism:* Production LLM deployments commonly rely on tensor parallelism (TP) [10]. With increasing TP degree, the LoRA adapters are sharded across the GPUs [41] and hence the added computational cost of adaptation is divided across GPUs. Fig 5 shows the performance impact of adapter ranks for different degrees of TP, where relative TTFT is computed by dividing the TTFT of each rank by that of rank-8. Though increasing TP diminishes the impact of rank heterogeneity, we still see a substantial 20% increase in TTFT when using a rank-128 adapter with TP=8, on Llama 7B.

*4) Impact on Service Level Objectives (SLOs):* Cloud providers operate under strict SLOs for tail latencies. For example, a common SLO for LLM Inference is that the P95 TTFT must be under 10s [28]. A direct consequence of varying performance with adapter ranks is that the same hardware can process fewer tokens per unit time when operating on higher ranks. Thus, to operate under SLO, more resources must be dedicated for larger ranks making their cost per token greater. Fig 6 shows the performance of a simple 4 RPS (requests per second) Poisson arrival workload on different ranks using the same underlying hardware. Assuming a P95 TTFT SLO of 20s, Fig 6 shows that high ranks like 64 and 128 violate SLO for the same workload whereas others do not.

*5) Effect of Co-serving Adapters of Different Ranks:* The kernels used for multi-tenant LoRA serving [18], [41] enable co-batching requests from different LoRA adapters. However, both kernels size their compute tiles and MMA pipelines to the maximum LoRA rank in the batch. In Punica [18], the BGMV forces all requests, regardless of their actual rank, to execute GEMMs padded to the largest rank, inflating register usage, shared-memory requirements, and tile dimensions. On the other hand, S-LoRA's MBGMV [41] mitigates some padding overhead, but still inherits the same dependency, i.e., the fused kernel's scheduling and memory layout are dictated by the highest rank present. As a result, the presence of high-rank requests in the batch causes low-rank requests to slow down, despite needing far fewer FLOPs. Kernel-level profiling confirms that throughput and latency track the maximum rank [33]. Consequently, co-batching heterogeneous ranks in these kernels bottlenecks performance, causing smaller-rank adapters to pay the computational cost of the largest one.

Fig 1 depicts the prefill performance when two different adapters are co-served on a single host. Evidently, co-serving rank 8 with rank 128 makes the smaller rank's tail performance slower by 84%. Moreover, greater rank heterogeneity leads to higher variability in performances, as indicated by the larger
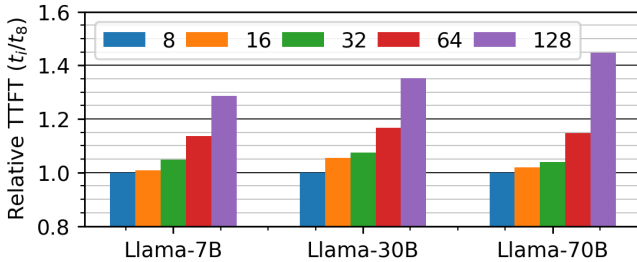


Fig. 4. *Relative TTFT of different ranks on increasing model size. Input size = 2000 and TP = 8.*
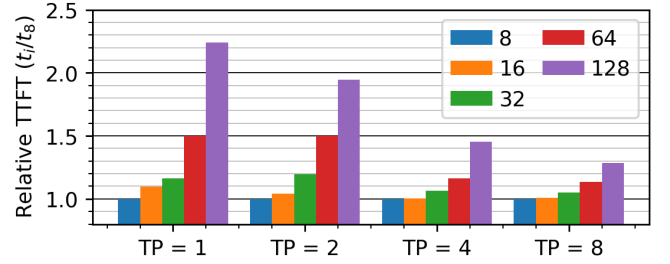


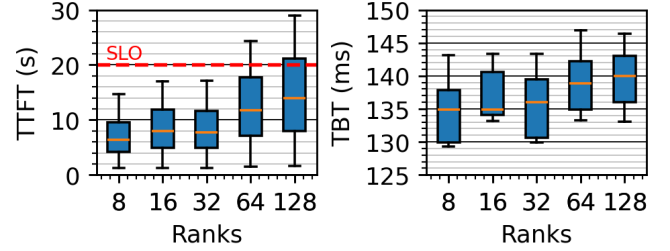Fig. 5. *Relative TTFT of different ranks on Llama 7B on increasing TP. Input size = 2000.*



Fig. 6. *4 RPS Poisson workload on Llama 7B with varying adapter ranks and fixed input (512) and output (128) lengths.*

ranges in the box plot. This calls for improvements at the system level to alleviate these issues.

### B. Workloads

We analyze production traces from LoRA deployments across regions for the week of October 22–29, 2025 at Company X . The characterization provides a comprehensive view of adapter distribution, memory footprint, traffic patterns, and arrival trends, offering critical insights into workload heterogeneity and its implications for system design and resource allocation.
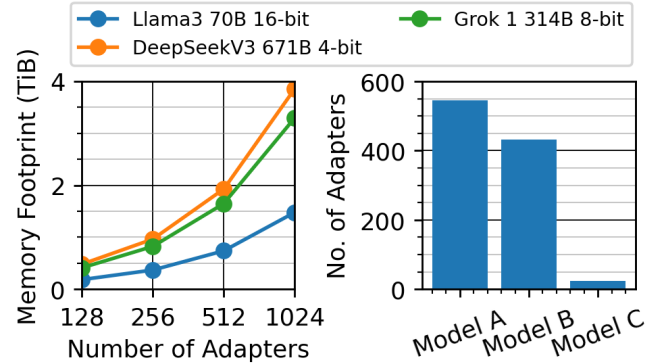


Fig. 7. *Estimated memory footprint of adapters (left) and no. of LoRA adapters for three base models at Company X (right).* Popular base models host hundreds of LoRA adapters with substantial memory footprints, making full colocation infeasible and motivating selective placement strategies.

*1) LoRA in Production:* Fig 7 shows the number of adapters for the three most popular base LLMs deployed at Company X. Certain base models host significantly more adapters than others, making it impractical to place all adapters on every LLM inference server [33], especially due to high memory footprint as illustrated in Fig 7. Moreover, out of the over thousand adapters currently deployed in production, merely the top 5 account for more than 70% of traffic as illustrated in Fig 8. Hence, they can be merged into the base model for serving on dedicated LLM Inference servers due to high
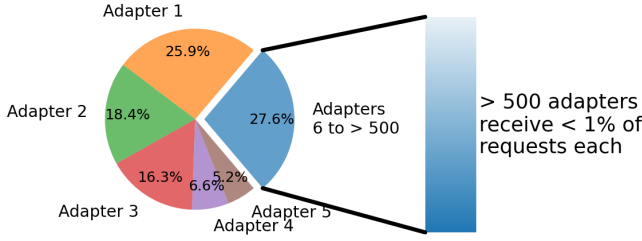
Fig. 8. ***Each adapter's request share for Model A at Company X.*** *Top 5 adapters account for more than 70% of requests, highlighting the need for demand-aware capacity allocation and co-location of low-demand adapters.*

demand [28]. The remaining 27.6% traffic is distributed across the rest of the adapter, each of which receive way less than 1% of the total traffic. Dedicating exclusive LLM inference servers for these may lead to resource under-utilization due to highly fluctuating request demands. Thus, in production, each server hosts a subset of adapters.

As illustrated in Fig 9, a few models dominate resource consumption, and this capacity is further concentrated in specific regions. Unlike vanilla LLM deployments, this imbalance arises because models are fine-tuned on customer data and are subject to data boundary constraints. Consequently, since capacity is regionally constrained rather than globally distributed, optimizing adapter placement at the server level offers greater potential for improving GPU utilization than global routing or load balancing strategies typically employed for vanilla LLMs.

*Long-tail model popularity*: Fig. 8 shows the request distribution for each adapter over one week from Model A in Region A. The top five adapters account for more than 70% of all requests. This implies that tailoring capacity allocation to adapter-specific demand patterns enables additional GPU savings. Co-locating multiple low-demand adapters on the same GPU, both spatially and temporally, can significantly improve resource utilization.
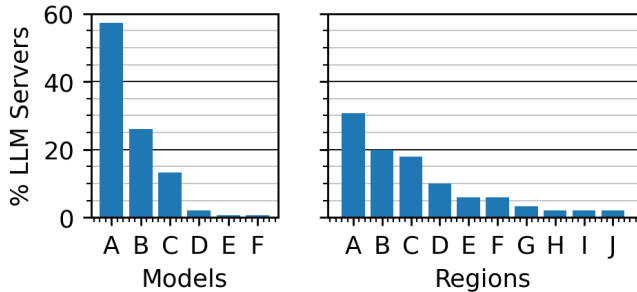


Fig. 9. ***Percentage of LLM servers used for each model (left) and at each region (right) at Company X.*** *Resource consumption is concentrated in a few models and regions due to data boundary constraints, emphasizing that server-level optimization is more effective than global routing.*

*Arrival Patterns of Top 5 Adapters*: Next, we examine the arrival patterns of the top five adapters for Model A in Region A. As shown in Fig 10, each adapter exhibits a distinct trend in request arrivals over the observed week, captured through the moving average of requests per minute Most adapters display smooth activity with gradual drifts, though anomalies are present. Notably, Adapter 1 and 3 show varying load trends while adapter 5 follows diurnal patterns. Adapter 4 exhibits

relatively stable demand whereas adapter 4 remains stable initially but suddenly sees a load surge towards the end of the trace. These fluctuating patterns highlight the potential benefits of dynamic rebalancing of resources allocated to adapters. This can improve the efficiency of capacity allocation, thereby improving GPU utilization and reducing resource overhead compared to static allocation strategies.
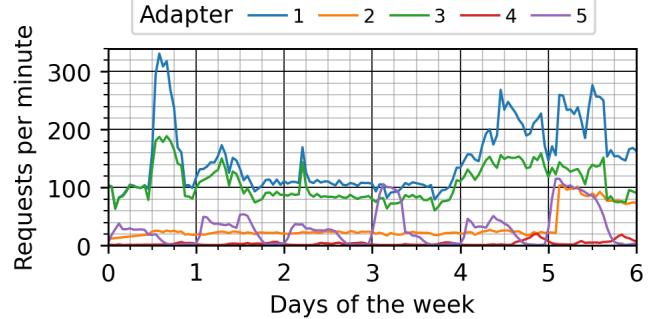


Fig. 10. ***Requests per minute for five adapters over one week in October 2025.*** *Arrival patterns of top adapters vary significantly, with drifts over time, highlighting the importance of workload-aware strategies for dynamic resource allocation.*

*2) Implications for System Design:* The presence of a few adapters that make up the majority of traffic conveys that GPU resource allocation should be based on adapter demand and not a simple uniform distribution across all adapters. The large number of adapters with little traffic must be efficiently multiplexed for effective resource utilization. This multiplexing cannot be static, since we see varying load patterns across adapters and across time, necessitating dynamic adapter placement based on adapter demand.

## IV. LORASERVE DESIGN

Based on insights from our characterization (Section III), we design LORASERVE, an LLM Inference cluster orchestrator designed for serving LoRA adapters at scale on a pool of LLM Instances. LORASERVE is designed to address the unique challenges faced by enterprises like Company X while scaling heterogeneous LoRA adapter serving. Broadly, LO-RASERVE addresses two challenges, i.e., (1) *the overheads of co-serving heterogeneous adapters* and (2) *the memory capacity pressure of storing all adapters in each server*.

To address the first issue, LORASERVE proposes a novel adapter placement and request routing algorithm that minimizes heterogeneity of adapters being served on a specific server while simultaneously balancing the request load evenly across the cluster by taking into account the popularity and workload demands of each adapter, at every time step. To accommodate thousands of adapters in a quickly accessible manner, LORASERVE only stores the adapters needed by a server in the current time step in its local host memory, ensuring that every adapter is present in at least one of the servers in the cluster. When the workload drifts and LORASERVE rebalances the adapter placement in the cluster, the new adapters are fetched over InfiniBand links using GPU-Direct RDMA on first access and then stored locally. The placement and routing module and the adapter migration module work synergistically,
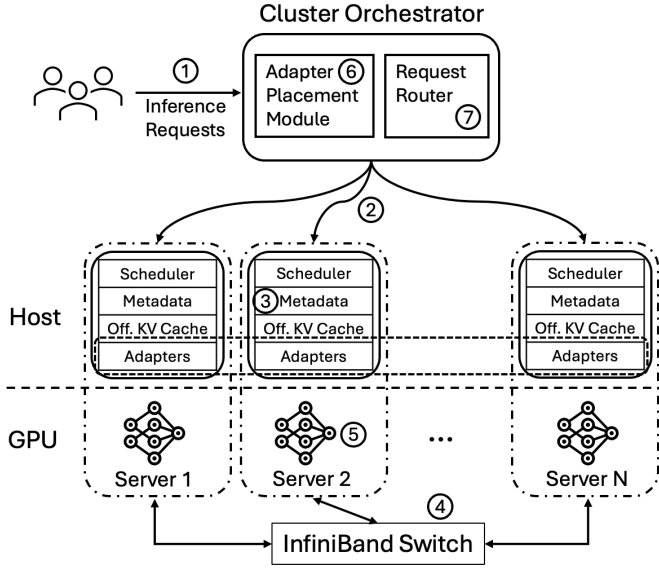
Fig. 11. **LORASERVE** *architecture overview. Each server stores and serves a subset of adapters managed by* LORASERVE.

maximizing the system's throughput and balancing the request load in the system.

**Architecture Overview.** Fig 11 shows an overview of the LORASERVE architecture. For sake of explanation, we assume that one LLM inference server runs one instance only. LO-RASERVE maintains a routing table in the cluster orchestrator that contains tuples of $(adapter\_id, server\_id, \phi)$ for every adapter. The adapter corresponding to $adapter\_id$ is stored on all servers corresponding to $server\_id$ across such tuples. An inference request ① specifying a particular adapter is received at the cluster orchestrator and routed ② to $server\_id$ with probability $\phi$. It is ensured that for every $adapter\_id$, $\Sigma\phi = 1$, implying the request is guaranteed to be executed on one of the servers where its adapter is allocated. In the host, LORASERVE looks up the adapter metadata to identify where the adapter is present ③. If the adapter is not in the local memory of the host, LORASERVE initiates a transfer over InfiniBand using GPUDirect RDMA to fetch it from the remote server where it is present ④. Finally, the LLM instance runs the inference ⑤ and returns the result. At the cluster orchestrator level, the LORASERVE service keeps track of how many requests each adapter receives. At every time step (configurable by the cluster admin), LORASERVE estimates the TPS (tokens per second) demand of every adapter and uses it, along with information regarding the rank and current location, to generate a balanced placement of adapters ⑥. Subsequently, the new mappings are updated in the request router ⑦ and used for subsequent requests.

### A. LORASERVE *Adapter Placement*

The adapter placement algorithm takes into consideration the load expected for every adapter in the cluster and makes allocation decisions to determine which server to use to serve each adapter. Allocation of an adapter to a server in

LORASERVE is a function of the projected workload demand of of that adapter, its rank and the target utilization of the cluster. By taking these factors into account, LORASERVE is able to generate adapter placements that minimize heterogeneity and balance request load across the cluster. For every *adapter* with projected load $\ell$, the expression below shows the output of the algorithm where $s_1, s_2 \ldots s_k$ are the $k$ servers in the cluster and $\phi_1, \phi_2 \ldots \phi_k$ denote the fractional load (probability of serving) allocated to each server where $\Sigma\phi = 1$.

$$f(adapter, \ell) = \big[(s_1, \phi_1), (s_2, \phi_2), \ldots, (s_k, \phi_k)\big]$$

---

**Algorithm 1** LORASERVE Algorithm

1: **function** ASSIGNLORASERVE (servers, adapters, requestHistory, operatingPoints, prevAssignment)
2:     demandTPS ← {}           ▷ Step 1
3:     **for each** adapter in adapters **do**
4:         adapterTPS ← GETPREVTIMESTEPTPS(adapter, requestHistory)
5:         demandTPS[adapter] ← EXTRAPOLATE(TPSHistory[adapter], adapterTPS)
6:     **end for**
7:     targetUtil ← 0
8:     rankUtil ← {}
9:     **for each** rank in UNIQUERANKS(adapters) **do**
10:        rankUtil ←
11: $\Sigma_{rank}$demandTPS/operatingPoints[rank]
12:        targetUtil ← targetUtil + rankUtil[rank]
13:     **end for**
14:     targetUtil ← targetUtil/LEN(servers)
15:     rankServerBudget ← {}        ▷ Step 2
16:     **for each** rank in UNIQUERANKS(adapters) **do**
17:        rankServerBudget[rank] ← ROUND(rankUtil[rank]/targetUtil)
18:     **end for**
19:     assignment ← {}          ▷ Step 3
20:     **for each** rank in rankServerBudget **do**
21:    FRACTIONALBINPACKING(adapters, rankServerBudget[rank], targetUtil, assignment)
22:     **end for**
23:     leftovers ← GETUNASSIGNEDADAPTERS(adapters, assignment)        ▷ Step 4
       SORTDESCENDINGRANK(leftovers)
24:     **for each** adapter in leftovers **do**
       ALLOCATEHIGHESTMAXRANK(adapter, assignment)    ▷ Allocates the adapter to the server with highest max rank and least utilization
25:     **end for**
       PERMUTEASSIGNMENT(assignment, prevAssignment)        ▷ Step 5
       UPDATEROUTINGTABLE(assignment)    ▷ Step 6
26:     **for each** server in servers **do**
       UPDATEADAPTERMAPPING(assignment, server)
27:     **end for**
28: **end function**

---

We profile the servers a priori, to estimate the operating point of each rank under SLO constraints, i.e., the maximum number of tokens per second the LLM inference server can process using an adapter of a specific rank under SLO constraints. Algorithm 1 illustrates LoRASERVE's adapter placement logic, using necessary abstractions for brevity. The algorithm follows the following basic steps at every time step:

1) Estimate the TPS demand per adapter and calculate average target utilization per server using operating point information.
2) Calculate the server budget per rank, i.e., the number of servers that can be dedicated for that rank.
3) Pack the adapters of ranks which have non-zero server budget assigned using fractional bin packing.
4) Preferentially allocate each of the remaining adapters on a server with higher maximum rank, if possible.
5) Permute the placement across servers to ensure minimal deviation with previous placement.
6) Update routing table and local adapter metadata.



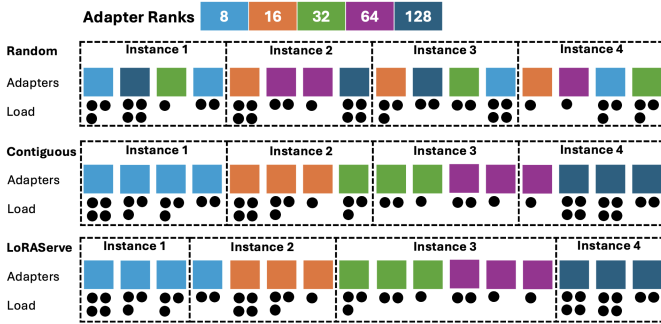Fig. 12. *A high-level illustration of adapter placement by* LoRASERVE *and baselines.*

Fig 12 shows a high level visualization of LoRASERVE's adapter placement along with baselines. *Random* placement does a good job of balancing the request load but does poorly in alleviating rank heterogeneity whereas naive rank-based *Contiguous* placement does a good job of mitigating rank heterogeneity but is not able to balance load. This results in the former underutilizing resources and the latter having poor tail latency. LoRASERVE on the other hand is able to optimize for both dimensions achieving better performance.

### B. LoRASERVE *Distributed Adapter Pool*

The LoRASERVE Adapter Placement (IV-A) and routing mechanism ensures that an LLM instance is only responsible for processing requests corresponding to a relatively small subset of adapters for every time step duration. Moreover, it also ensures that the union of adapters assigned to every instance is the universal set of adapters that the cluster may need to use. The adapters assigned to a specific instance are stored in the main memory of the host. Combined, all locally stored adapters give the abstraction of the *distributed adapter pool* in LoRASERVE. An adapter may be assigned to one or more LLM servers depending on its popularity and demand.

The LoRASERVE service running on the cluster orchestrator maintains an in-memory map of all the adapters in the cluster to help identify where they are stored. At every timestep, the adapter placement routine updates the placement and this map based on the latest adapter popularity patterns.
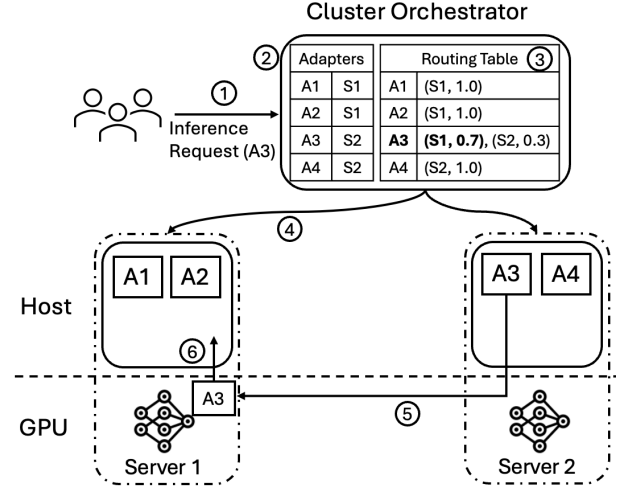


Fig. 13. LoRASERVE *Distributed Adapter Pool.* The cluster orchestrator maintains a table containing the location of each adapter in the cluster.

Consider the example illustrated in Fig 13. Assume that the adapter placement module recently changed the allocation of $A3$ from being served entirely on Server 2 ($A3 : (S2, 1.0)$), to serving only 30% of its traffic on Server 2 and the rest on Server 1 ($A3 : (S1, 0.7), (S2, 0.3)$), but the actual migration hasn't occurred yet. When a request for adapter $A3$ arrives at the cluster orchestrator ①, LoRASERVE looks up the adapters table ② to figure out where the adapter is present. Subsequently, it looks up the routing table ③ and decides which server to send the request to. Assume that $S1$ is chosen for this example. The router then sends the adapter location information along with the inference request to $S1$ ④. Since the adapter is absent in $S1$, it is fetched from $S2$ over InfiniBand interconnects using GPUDirect-RDMA. To do so, the adapter must first be copied from host memory to GPU in $S2$ and then transferred to $S1$ ⑤. Once the request is processed, the adapter is saved in the host memory of $S1$ for future requests. Given that these operations are deterministic, the cluster orchestrator updates the adapter table to $A3 : (S1, S2)$ upon successful completion of the request. In case if the adapter was no longer needed at $S2$ (i.e., the routing table had ($A3 : (S1, 1.0)$)), $A3$ would be deleted from $S2$ after being copied to $S1$.

Transferring adapters over InfiniBand [8] between GPUs of different servers incurs similar latency as transferring it over the same host's memory channel from local memory to GPU. Moreover, we explored the possibility of replicating all adapters in the local SSD of each server but the latency and bandwidth was found to be prohibitive. Fig 14 benchmarks this trend. InfiniBand enabled us to perform the transfer operation without increasing the overheads.
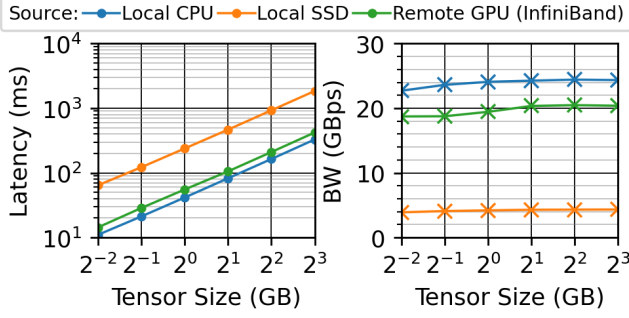
Fig. 14. *Latency of fetching a tensor from different sources.*

## V. EVALUATION

We implement LORASERVE on top of S-LoRA [41], a recent system built to serve thousands of LoRA adapters at scale. It is compatible with any LLM Inference framework that supports LoRA adapters and can also be integrated into common LLM cluster orchestrators [11], [12] used in production due to its modular design.

### A. Hardware

We run our experiments on Azure Machine Learning Compute Clusters [13], one containing Standard_ND96asr_v4 VM [14] nodes and the other containing Standard_NC24ads_A100_v4 VM [15] nodes. Each node in the former is equipped with $8\times$ Nvidia A100 SXM 40GB GPUs connected through NVLINK 3.0, an AMD EPYC 7V12 CPU with 96 cores along with 900 GiB of host memory and 6TiB disk storage. Internode communication is through InfiniBand [8] and Azure Accelerated Networking [16]. In the latter, each node contains $1\times$ Nvidia A100 80GB GPU, an AMD EPYC 7V13 CPU with 24 cores, 220GiB of host memory and 960GiB of disk storage.

### B. Experiment Setup

Due to limited resources, we run most of our experiments on a cluster of 4 LLM inference servers and validate scalability of LORASERVE on up to 12 servers. We use Llama 7B with TP = 4 on a uniform poisson trace of specified RPS, unless stated otherwise.

### C. Models

We evaluate LORASERVE on the Llama family of models [7]. For majority of our experiments, we use Llama-7B. However, we also show LORASERVE's sensitivity to model sizes using Llama-7B along with Llama-30B and Llama-70B. While we present the results with Llama series of models, we also experimented with other models like OPT [51]. We observed similar performance trends.

### D. Baselines

We compare LORASERVE against the following state-of-the-art baselines, keeping the underlying kernel the same across LORASERVE and the baselines for a fair comparison:

1) *S-LoRA Random* allocates adapters randomly to servers in the cluster, each running S-LoRA [41]. This adapter placement resembles the one used at Company X .
2) *S-LoRA Contiguous* orders the adapters being served on the cluster by their ranks and places an equal number of adapters on each server contiguously. This ensures that ranks close to each other are co-located.
3) *Toppings* [33] is a state of the art LoRA serving system that proposes a new load aware scheduler. The basic idea is to route the incoming request to the globally optimal server after considering the requests currently being served and queued at each server.
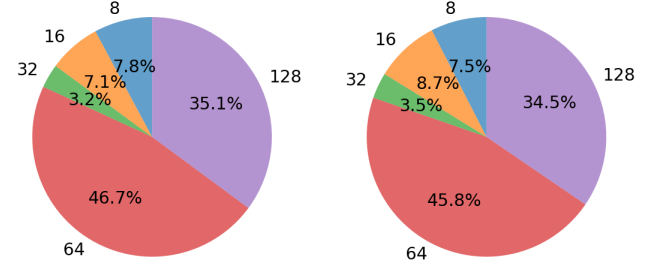
### E. Traces



Fig. 15. *Rank-wise request (left) and token (right) distribution of production trace.*

We evaluate LORASERVE on production traces of real users interacting with Company X's services that use LoRA adapters at scale. The trace contains 250,138 requests over 8 hours going to 5 adapters of different ranks used in production services as illustrated in Fig 15. We divide each adapter's requests in the trace by annotating the traces with different adapter names of the same rank following a power law distribution for adapter counts within a rank, with $\alpha = 1$ totaling to 50, 100 and 200 to model three scenarios involving different total number of adapters. For evaluating on a desired RPS, we scale the timestamps proportionally to retain the original arrival pattern. The traces have the following attributes for every request: request_id, model, adapter, prompt_length, output_length and timestamp.

We also evaluate LORASERVE on open-source Azure Public Dataset [9] traces from 2024. As these traces lack timestamps and adapter names, we annotate them with these attributes ourselves using some common assumptions used in prior work. For timestamps, we assume the arrival pattern of the requests to be either uniformly distributed or following a Poisson process [41] as used in many prior works on LLM inference. We use a total of 25 adapters of ranks 8, 16, 32, 64 and 128 in our experiments as used by prior work [26], [33]. For adapters' rank popularity, we evaluate on the following:

1) Uniform: This assumes that all adapter rank popularities are uniformly distributed throughout the trace.
2) Shifting Skew: This is illustrated in Fig 16. In the beginning, rank 128 gets half the traffic and the other half is distributed uniformly between the other four

ranks. This skew shifts linearly until, at the end of the trace, rank 8 sees half the traffic and the rest is divided uniformly among the other ranks.
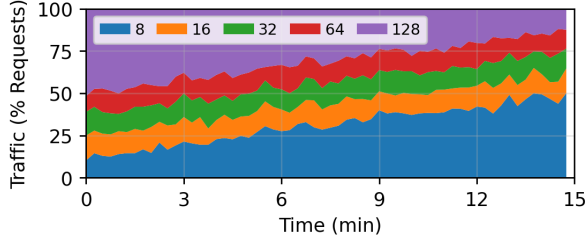


Fig. 16. *Shifting skew in adapter rank popularities.*

3) Exponential: This assumes that the popularities of ranks are exponentially distributed [26] with the smaller ranks being more popular as done by prior work.

Combining both dimensions, i.e., request arrival pattern and adapter popularity, we get six unique traces that we use for subsequent evaluation.
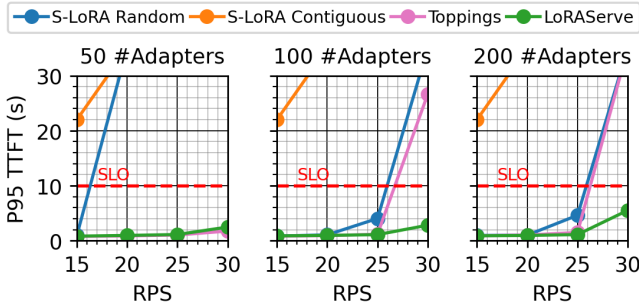


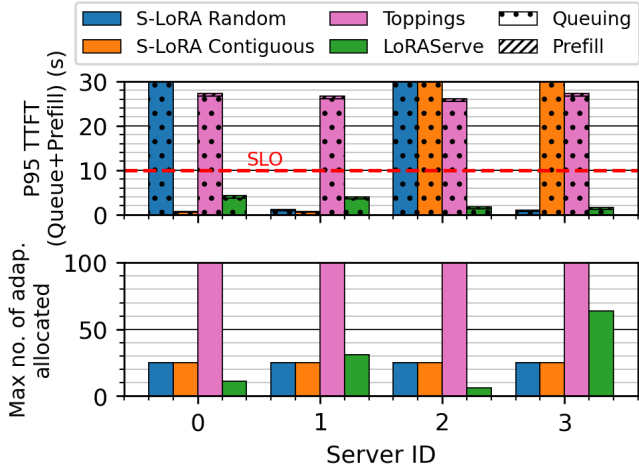Fig. 17. *Evaluation on production traces from Company X .*



Fig. 18. *Server-wise queuing and prefill time (top) and maximum no. of adapters used (bottom)* for 30 RPS 100 #Adapters experiment on production traces.

### F. Performance

Fig 17 shows the performance of LORASERVE and baselines on production traces from Company X. As described earlier, we model three scenarios where we have a total of 50, 100 and 200 adapters each. We observe that LORASERVE is

able to serve upto 20% higher request throughput with respect to Toppings [33] and upto $2\times$ higher than S-LoRA [41], under SLA. This translates into GPU savings of 17% compared to Toppings and 50% compared to S-LoRA when serving the same workload. This can be attributed to LORASERVE's rank aware adapter allocation for load balancing. Fig 18 (top) shows the tail TTFT latency observed on each LLM inference server in the cluster. It is evident that S-LoRA Random and S-LoRA Contiguous, because of static allocation of adapters to servers, fail to balance the load evenly causing some servers to have excessively high load and time out. Toppings, on the other hand, performs load balancing at the request level, which is more granular, and does a much better job at having uniform performance on all servers. However, due to its rank agnostic behavior, it ends up sending high rank requests to all servers, causing interference leading to high amount of queuing and high TTFT. As LORASERVE is rank-aware, it is able to balance adapter placement dynamically and reduce queuing leading to upto $9\times$ lower TTFT. Furthermore, Fig 18 (bottom) shows the maximum amount of adapters needed on every server by the four systems. We observe that LORASERVE reduces the amount of space needed for adapter storage by upto $16\times$ with respect to state of the art Toppings. This is because LORASERVE allocates adapters to servers after taking into account both, the demand and rank of every adapter. Thus, a small number of same-ranked highly popular adapters get a dedicated server (server 0) and many sparsely used adapters get put together (server 3).

Fig 19 and Fig 20 show the performance of LO-RASERVE on the six traces derived from open-source production data (section V-E). While LORASERVE beats the baselines in all cases, the general trend in Fig 19 shows that S-LoRA Contiguous has reasonable performance when the popularity of adapters are uniformly distributed as it places approximately equal number of adapters having ranks close to one another on every server. This ensures good load balancing as every adapter receives similar traffic. However, it times out very quickly on other distributions. On the other hand, we observe S-LoRA Random to have passable performance in all cases except when the adapter demand changes (Shifting Skew) in the request trace. Moreover, S-LoRA Random also times out at high request rates for most workloads and isn't able to scale at high traffic.

As LORASERVE accounts for both adapter rank and its demand, it is able to elicit much better performance and serve significantly higher request throughput on the same hardware. Quantitatively, we observe a TTFT improvement of upto $9\times$. Furthermore, these improvements are achieved without any drop in request TBT (Fig 20). Instead, LORASERVE delivers up to a 15% TBT improvement in some cases over the baselines due to better load balancing across servers in the cluster.

### G. Scalability

Fig 21 shows the performance of LORASERVE when scaled to large number of nodes, i.e., clusters of size 8 and 12 LLM
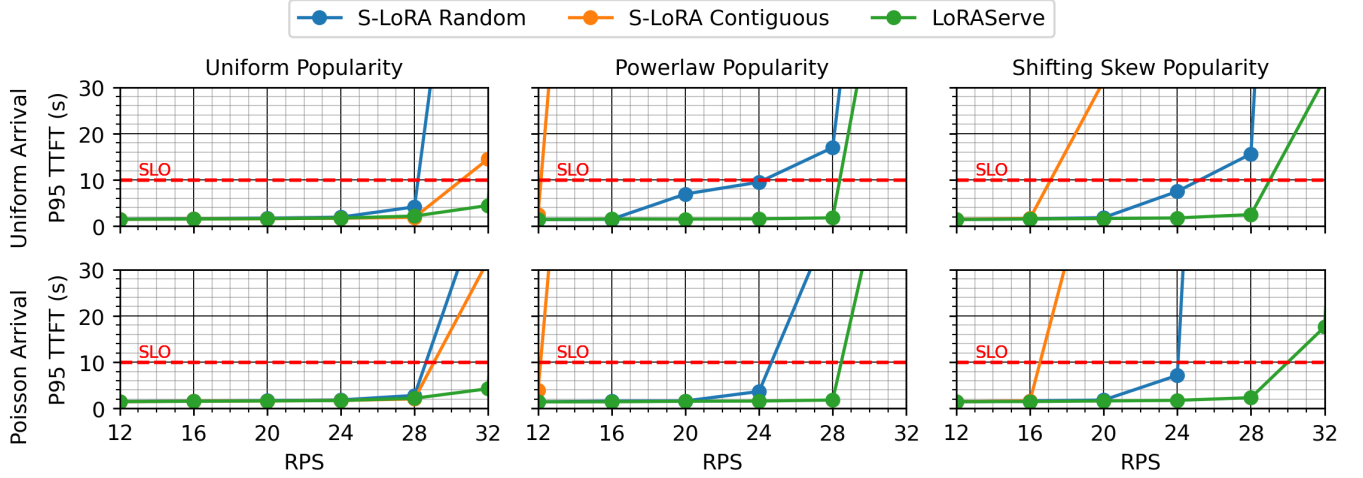
Fig. 19. *TTFT performance on different workloads.* LORASERVE *acheives up to 9× improvement in TTFT over baselines.*
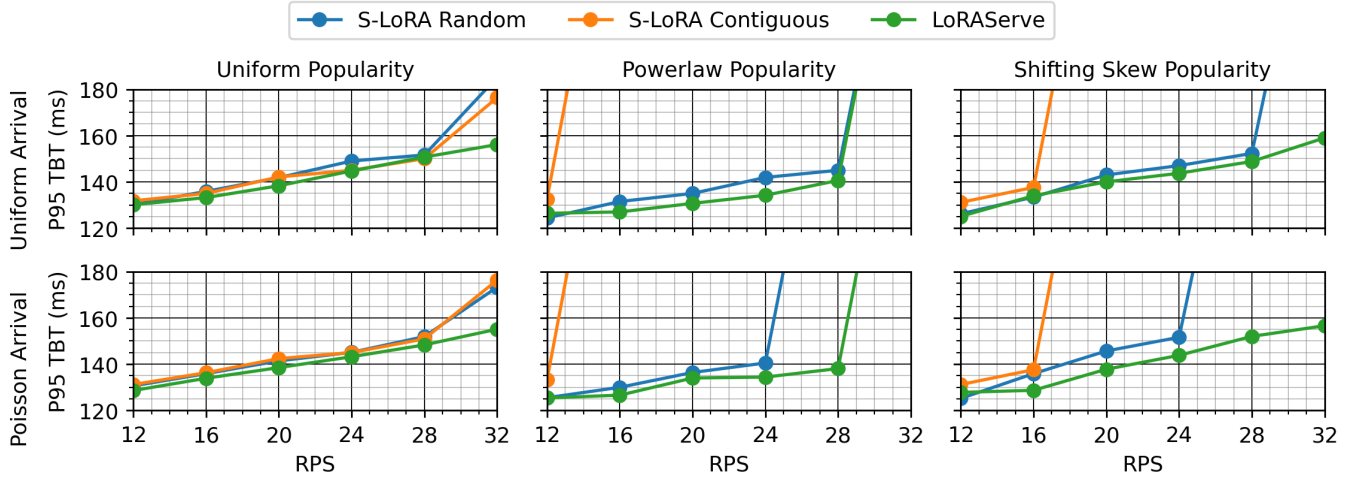


Fig. 20. *TBT performance on different workloads. TBT either remains similar in most cases or improves by up to 15%.*
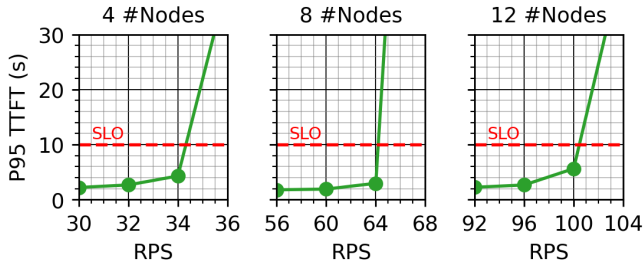


Fig. 21. *Scalability of LORASERVE on up to 12 servers.*

servers along with 4. We also scale the number of adapters and request traffic being served on the cluster proportionally with the number of nodes to evaluate weak scaling [24], [29]. It can be observed that LORASERVE scales well with the number of nodes and is able to sustain a proportional amount of traffic with every added node. For example, in Fig 21, LORASERVE is able to serve approximately 32 RPS workload within a 10s P95 TTFT SLO on a cluster of 4 LLM servers, which translates into 64 and 96 RPS for clusters of size 8 and 12 servers respectively as confirmed by the experiment.

### H. Sensitivity to Rank Skews

Fig 22 shows performance trends on traces with varying rank skews where the adapter popularity follows a power-law distribution ($y = x^{-\alpha}$) with smaller ranks receiving more requests [26]. We observe that although LORASERVE always shows tail TTFT performance well within SLO, upon increasing $\alpha$ from 1/3 to 3, the tail behavior slightly improves. This is because at lower $\alpha$, larger ranks receive a sizable share ($\geq 16\%$) of total requests whereas as $\alpha$ increases, the share presence of the largest rank in the trace diminishes and the trace is mostly dominated by small ranks as illustrated in Fig 22. Thus, the performance gains can be attributed to smaller rank requests being naturally lighter. However, varying rank skews do not have any adverse effect on LORASERVE's performance unlike the baselines.

The baselines, on the other hand, suffer with increasing skew in adapter popularity. Contiguous adapter placement shows poor tail performance and times out at higher skews due to load imbalance across nodes in the cluster as it only considers adapter ranks. Random adapter placement depicts performance similar to LORASERVE at smaller skews but
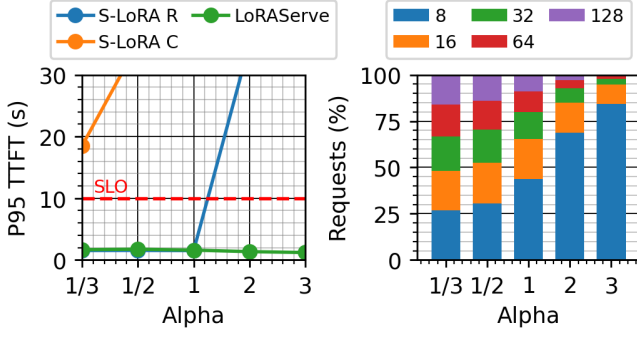
Fig. 22. *Varying α in power law distribution for adapter popularity on a 36 RPS Poisson arrival trace with 100 adapters (20 of each rank).*

times out as the skew increases. This is because, as a small number of adapters become highly popular, it is easy for a bad random allocation to happen, causing severe load imbalance.
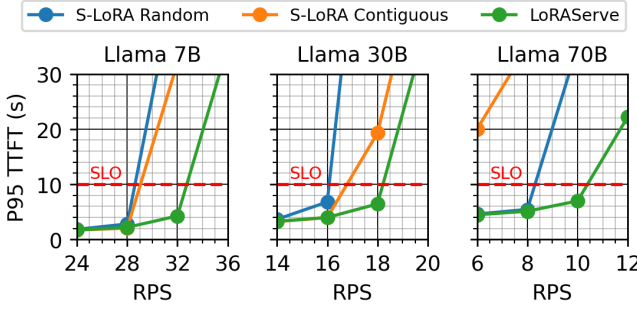
### I. Sensitivity to Model Sizes



Fig. 23. *Performance of LORASERVE using models of different sizes.*

General purpose LLMs used in production at Company X are generally larger in size, often upto 200B parameters. Fig 4 shows the performance of LORASERVE and baselines on larger models from the Llama family. We observe that the performance trends are quite similar to what is discussed in Section V-E w.r.t. Llama 7B, except that the larger models serve lower throughput in general which is orthogonal to the design of LORASERVE. The baseline performance trends also translate similarly.
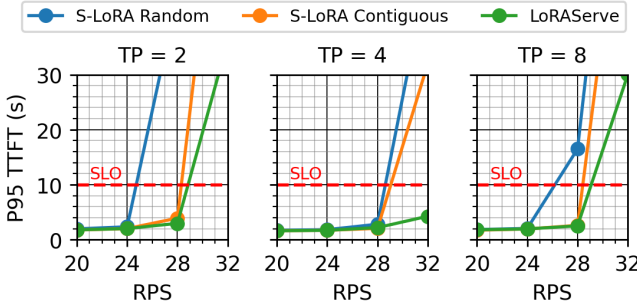
### J. Sensitivity to TP



Fig. 24. *Performance of LORASERVE using Llama 7B with different TP configurations.*

Production LLMs are usually deployed under different configurations. We show that LORASERVE's performance gains are orthogonal to the model's deployment configuration as it

operates on a cluster scale. Fig 24 shows performance trends when models are deployed using various tensor parallelism configurations. In every case, LORASERVE complements the baselines, enabling higher throughput and lower TTFT.

## VI. RELATED WORK

### A. LLM Inference

Making LLM inference more efficient is the target of many recent studies. A vast corpora of recent papers have proposed hardware [32], [50], software [5], kernel-level [19], [48] and system-level [28] optimizations for improving the performance [39], [49], cost [28] and energy efficiency [42] of LLM Inference. These include optimizations in batching [17], [52], scheduling [39], [49], load-balancing [27], [39] and KV-caching [35], [46] among others. While most of these works optimize for a monolithic LLM Instance deployed on a single GPU server, there are only a handful which consider cluster level optimizations. However, LORASERVE optimizes serving large scale serving of LoRA adapter workloads for a cluster of LLMs, and can be used in conjunction with them to improve the cluster's utilization and performance.

### B. Parameter Efficient Fine-Tuning

LoRA [25] and its derivatives such as QLoRA [21] are the most popular of several parameter-efficient fine tuning [47] techniques [30], [31], [34], [36], [43]. LoRA allows for zero overhead during inference by merging the adapter weights into the base model, as suggested in the original work [25].

While this approach is suitable when serving a single adapter, it does not scale to multiple adapters, especially of varying ranks and traffic patterns where multiplexing of adapters becomes necessary. Several new inference techniques address this issue. Punica [18] and S-LoRA [41] develop special CUDA kernels to batch requests from different adapters and decouple base model and adapter computation. S-LoRA also unifies paging across KV cache and LoRA adapters stored on the GPU, extending the paging seen in serving systems like vLLM [5]. dLoRA [45] dynamically merges and unmerges adapters with the base model and migrates requests and adapters between replicas. Chameleon [26] presents adapter caching on the GPU and adapter-aware scheduling. Toppings [33] uses the CPU for LoRA computation while adapters are being loaded to the GPU to improve TTFT.

## VII. CONCLUSION

Rank heterogeneity and memory pressure in multi-tenant LoRA serving lead to severe performance degradation and inefficient resource utilization. Guided by insights from production like skewed adapter popularity and rank-induced interference, LORASERVE dynamically allocates and migrates adapters based on their rank and demand, optimizing performance and alleviating memory capacity pressure. Across diverse workloads and configurations, LORASERVE delivers up to 2× higher throughput, up to 9× lower TTFT, and uses up to 50% fewer GPUs under SLO constraints, while reducing adapter storage memory footprint by up to 16×, consistently outperforming state-of-the-art systems.

<center>REFERENCES</center>

[1] OpenAI Documentation: Model Optimization. [Online]. Available: https://platform.openai.com/docs/guides/model-optimization

[2] Anyscale: Fine-Tune LLM Models at Scale. [Online]. Available: https://www.anyscale.com/use-case/llm-fine-tuning

[3] AWS Documentation: Amazon SageMaker AI - Fine-Tune a Model. [Online]. Available: https://docs.aws.amazon.com/sagemaker/latest/dg/jumpstart-fine-tune.html

[4] Microsoft Azure: Fine-tune models with Azure AI Foundry. [Online]. Available: https://learn.microsoft.com/en-us/azure/ai-foundry/concepts/fine-tuning-overview

[5] vLLM. [Online]. Available: https://github.com/vllm-project/vllm

[6] AWS Blogs: Manage hundreds of LoRA adapters with SageMaker. [Online]. Available: https://aws.amazon.com/blogs/machine-learning/easily-deploy-and-manage-hundreds-of-lora-adapters-with-sagemaker-efficient-multi-adapter-inference/

[7] Wikipedia Llama (language model). [Online]. Available: https://en.wikipedia.org/wiki/Llama_(language_model)

[8] Nvidia InfiniBand. [Online]. Available: https://network.nvidia.com/pdf/whitepapers/IB_Intro_WP_190.pdf

[9] Azure Public Dataset. [Online]. Available: https://github.com/Azure/AzurePublicDataset

[10] Meta Blogs: Scaling LLM Inference. [Online]. Available: https://engineering.fb.com/2025/10/17/ai-research/scaling-llm-inference-innovations-tensor-parallelism-context-parallelism-expert-parallelism/

[11] LLM-d. [Online]. Available: https://llm-d.ai/

[12] vLLM Production Stack. [Online]. Available: https://github.com/vllm-project/production-stack

[13] Azure Machine Learning Compute Cluster. [Online]. Available: https://learn.microsoft.com/en-us/azure/machine-learning/how-to-create-attach-compute-cluster?view=azureml-api-2&tabs=python

[14] Azure NDasrA100_v4 Series Specifications. [Online]. Available: https://learn.microsoft.com/en-us/azure/virtual-machines/sizes/gpu-accelerated/ndasra100v4-series?tabs=sizebasic

[15] Azure NC24ads_A100_v4 Series Specifications. [Online]. Available: https://learn.microsoft.com/en-us/azure/virtual-machines/sizes/gpu-accelerated/nca100v4-series?tabs=sizebasic

[16] Azure Accelerated Networking. [Online]. Available: https://learn.microsoft.com/en-us/azure/virtual-network/accelerated-networking-overview

[17] A. Agrawal, N. Kedia, A. Panwar, J. Mohan, N. Kwatra, B. Gulavani, A. Tumanov, and R. Ramjee, "Taming {Throughput-Latency} tradeoff in {LLM} inference with {Sarathi-Serve}," in *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, 2024, pp. 117–134.

[18] L. Chen, Z. Ye, Y. Wu, D. Zhuo, L. Ceze, and A. Krishnamurthy, "Punica: Multi-tenant lora serving," 2023. [Online]. Available: https://arxiv.org/abs/2310.18547

[19] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré, "Flashattention: Fast and memory-efficient exact attention with io-awareness," *Advances in neural information processing systems*, vol. 35, pp. 16344–16359, 2022.

[20] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, "Llm.int8(): 8-bit matrix multiplication for transformers at scale," 2022. [Online]. Available: https://arxiv.org/abs/2208.07339

[21] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "Qlora: Efficient finetuning of quantized llms," 2023. [Online]. Available: https://arxiv.org/abs/2305.14314

[22] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, and H. Wang, "Retrieval-augmented generation for large language models: A survey," 2024. [Online]. Available: https://arxiv.org/abs/2312.10997

[23] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Upasani, K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla, K. Lakhotia, L. Rantala-Yeary, L. van der Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo, L. Blecher, L. Landzaat, L. de Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri, M. Kardas, M. Tsimpoukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si, M. K. Singh, M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, N. Zhang, O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Krishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Stojnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly, R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim, S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang, S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler, T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami, V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti, V. Albiero, V. Petrovic, W. Chu, W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenberg, A. Baevski, A. Feinstein, A. Kallet, A. Sangani, A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani, A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury, A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang, B. Loyd, B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty, D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood, E.-T. Le, E. Brinkman, E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos, F. Ozgenel, F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri, H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan, I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski, J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan, K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A, L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani, M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi, M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev, N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy, R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked,

V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait, Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma, "The llama 3 herd of models," 2024. [Online]. Available: https://arxiv.org/abs/2407.21783

[24] J. L. Gustafson, "Reevaluating amdahl's law," *Commun. ACM*, vol. 31, no. 5, p. 532–533, May 1988. [Online]. Available: https://doi.org/10.1145/42411.42415

[25] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," 2021. [Online]. Available: https://arxiv.org/abs/2106.09685

[26] N. Iliakopoulou, J. Stojkovic, C. Alverti, T. Xu, H. Franke, and J. Torrellas, "Chameleon: Adaptive caching and scheduling for many-adapter llm inference environments," in *Proceedings of the 58th IEEE/ACM International Symposium on Microarchitecture*, ser. MICRO '25. New York, NY, USA: Association for Computing Machinery, 2025, p. 217–231. [Online]. Available: https://doi.org/10.1145/372584 3.3756083

[27] K. Jain, A. Parayil, A. Mallick, E. Choukse, X. Qin, J. Zhang, I. Goiri, R. Wang, C. Bansal, V. Ruehle, A. Kulkarni, S. Kofsky, and S. Rajmohan, "Performance aware llm load balancer for mixed workloads," in *EuroMLSys 2025*, April 2025. [Online]. Available: https://www.microsoft.com/en-us/research/publication/performance-awa re-llm-load-balancer-for-mixed-workloads/

[28] S. Jaiswal, K. Jain, Y. Simmhan, A. Parayil, A. Mallick, R. Wang, R. S. Amant, C. Bansal, V. Rühle, A. Kulkarni, S. Kofsky, and S. Rajmohan, "Sageserve: Optimizing llm serving on cloud data centers with forecast aware auto-scaling," 2025. [Online]. Available: https://arxiv.org/abs/2502.14617

[29] S. Jaiswal, S. Raj, S. Sidhanta, and Y. Simmhan, "Aerialdb: A federated peer-to-peer spatio-temporal edge datastore for drone fleets," *Pervasive and Mobile Computing*, vol. 114, p. 102109, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1574119225000987

[30] R. Karimi Mahabadi, J. Henderson, and S. Ruder, "Compacter: Efficient low-rank hypercomplex adapter layers," *Advances in neural information processing systems*, vol. 34, pp. 1022–1035, 2021.

[31] B. Lester, R. Al-Rfou, and N. Constant, "The power of scale for parameter-efficient prompt tuning," 2021. [Online]. Available: https://arxiv.org/abs/2104.08691

[32] J. Li, J. Xu, S. Huang, Y. Chen, W. Li, J. Liu, Y. Lian, J. Pan, L. Ding, H. Zhou *et al.*, "Large language model inference acceleration: A comprehensive hardware perspective," *arXiv preprint arXiv:2410.04466*, 2024.

[33] S. Li, H. Lu, T. Wu, M. Yu, Q. Weng, X. Chen, Y. Shan, B. Yuan, and W. Wang, "Toppings: Cpu-assisted, rank-aware adapter serving for llm inference," in *Proceedings of the 2025 USENIX Conference on Usenix Annual Technical Conference*, ser. USENIX ATC '25. USA: USENIX Association, 2025.

[34] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," 2021. [Online]. Available: https://arxiv.org/abs/2101.00190

[35] A. Liu, J. Liu, Z. Pan, Y. He, G. Haffari, and B. Zhuang, "Minicache: Kv cache compression in depth dimension for large language models," *Advances in Neural Information Processing Systems*, vol. 37, pp. 139 997–140 031, 2024.

[36] H. Liu, D. Tam, M. Muqeeth, J. Mohta, T. Huang, M. Bansal, and C. A. Raffel, "Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 1950–1965, 2022.

[37] J. Liu, D. Zhu, Z. Bai, Y. He, H. Liao, H. Que, Z. Wang, C. Zhang, G. Zhang, J. Zhang, Y. Zhang, Z. Chen, H. Guo, S. Li, Z. Liu, Y. Shan, Y. Song, J. Tian, W. Wu, Z. Zhou, R. Zhu, J. Feng, Y. Gao, S. He, Z. Li, T. Liu, F. Meng, W. Su, Y. Tan, Z. Wang, J. Yang, W. Ye, B. Zheng, W. Zhou, W. Huang, S. Li, and Z. Zhang, "A comprehensive survey on long context language modeling," 2025. [Online]. Available: https://arxiv.org/abs/2503.17407

[38] OpenAI, J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, R. Avila, I. Babuschkin, S. Balaji, V. Balcom, P. Baltescu, H. Bao, M. Bavarian, J. Belgum, I. Bello, J. Berdine, G. Bernadett-Shapiro, C. Berner, L. Bogdonoff, O. Boiko, M. Boyd, A.-L. Brakman, G. Brockman, T. Brooks, M. Brundage, K. Button, T. Cai, R. Campbell,

A. Cann, B. Carey, C. Carlson, R. Carmichael, B. Chan, C. Chang, F. Chantzis, D. Chen, S. Chen, R. Chen, J. Chen, M. Chen, B. Chess, C. Cho, C. Chu, H. W. Chung, D. Cummings, J. Currier, Y. Dai, C. Decareaux, T. Degry, N. Deutsch, D. Deville, A. Dhar, D. Dohan, S. Dowling, S. Dunning, A. Ecoffet, A. Eleti, T. Eloundou, D. Farhi, L. Fedus, N. Felix, S. P. Fishman, J. Forte, I. Fulford, L. Gao, E. Georges, C. Gibson, V. Goel, T. Gogineni, G. Goh, R. Gontijo-Lopes, J. Gordon, M. Grafstein, S. Gray, R. Greene, J. Gross, S. S. Gu, Y. Guo, C. Hallacy, J. Han, J. Harris, Y. He, M. Heaton, J. Heidecke, C. Hesse, A. Hickey, W. Hickey, P. Hoeschele, B. Houghton, K. Hsu, S. Hu, X. Hu, J. Huizinga, S. Jain, S. Jain, J. Jang, A. Jiang, R. Jiang, H. Jin, D. Jin, S. Jomoto, B. Jonn, H. Jun, T. Kaftan, Łukasz Kaiser, A. Kamali, I. Kanitscheider, N. S. Keskar, T. Khan, L. Kilpatrick, J. W. Kim, C. Kim, Y. Kim, J. H. Kirchner, J. Kiros, M. Knight, D. Kokotajlo, Łukasz Kondraciuk, A. Kondrich, A. Konstantinidis, K. Kosic, G. Krueger, V. Kuo, M. Lampe, I. Lan, T. Lee, J. Leike, J. Leung, D. Levy, C. M. Li, R. Lim, M. Lin, S. Lin, M. Litwin, T. Lopez, R. Lowe, P. Lue, A. Makanju, K. Malfacini, S. Manning, T. Markov, Y. Markovski, B. Martin, K. Mayer, A. Mayne, B. McGrew, S. M. McKinney, C. McLeavey, P. McMillan, J. McNeil, D. Medina, A. Mehta, J. Menick, L. Metz, A. Mishchenko, P. Mishkin, V. Monaco, E. Morikawa, D. Mossing, T. Mu, M. Murati, O. Murk, D. Mély, A. Nair, R. Nakano, R. Nayak, A. Neelakantan, R. Ngo, H. Noh, L. Ouyang, C. O'Keefe, J. Pachocki, A. Paino, J. Palermo, A. Pantuliano, G. Parascandolo, J. Parish, E. Parparita, A. Passos, M. Pavlov, A. Peng, A. Perelman, F. de Avila Belbute Peres, M. Petrov, H. P. de Oliveira Pinto, Michael, Pokorny, M. Pokrass, V. H. Pong, T. Powell, A. Power, B. Power, E. Proehl, R. Puri, A. Radford, J. Rae, A. Ramesh, C. Raymond, F. Real, K. Rimbach, C. Ross, B. Rotsted, H. Roussez, N. Ryder, M. Saltarelli, T. Sanders, S. Santurkar, G. Sastry, H. Schmidt, D. Schnurr, J. Schulman, D. Selsam, K. Sheppard, T. Sherbakov, J. Shieh, S. Shoker, P. Shyam, S. Sidor, E. Sigler, M. Simens, J. Sitkin, K. Slama, I. Sohl, B. Sokolowsky, Y. Song, N. Staudacher, F. P. Such, N. Summers, I. Sutskever, J. Tang, N. Tezak, M. B. Thompson, P. Tillet, A. Tootoonchian, E. Tseng, P. Tuggle, N. Turley, J. Tworek, J. F. C. Uribe, A. Vallone, A. Vijayvergiya, C. Voss, C. Wainwright, J. J. Wang, A. Wang, B. Wang, J. Ward, J. Wei, C. Weinmann, A. Welihinda, P. Welinder, J. Weng, L. Weng, M. Wiethoff, D. Willner, C. Winter, S. Wolrich, H. Wong, L. Workman, S. Wu, J. Wu, M. Wu, K. Xiao, T. Xu, S. Yoo, K. Yu, Q. Yuan, W. Zaremba, R. Zellers, C. Zhang, M. Zhang, S. Zhao, T. Zheng, J. Zhuang, W. Zhuk, and B. Zoph, "Gpt-4 technical report," 2024. [Online]. Available: https://arxiv.org/abs/2303.08774

[39] P. Patel, E. Choukse, C. Zhang, A. Shah, I. Goiri, S. Maleki, and R. Bianchini, "Splitwise: Efficient generative llm inference using phase splitting," in *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2024, pp. 118–132.

[40] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. C. Ferrer, A. Grattafiori, W. Xiong, A. Défossez, J. Copet, F. Azhar, J. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve, "Code llama: Open foundation models for code," 2024. [Online]. Available: https://arxiv.org/abs/2308.12950

[41] Y. Sheng, S. Cao, D. Li, C. Hooper, N. Lee, S. Yang, C. Chou, B. Zhu, L. Zheng, K. Keutzer, J. E. Gonzalez, and I. Stoica, "S-lora: Serving thousands of concurrent lora adapters," 2024. [Online]. Available: https://arxiv.org/abs/2311.03285

[42] J. Stojkovic, C. Zhang, I. Goiri, J. Torrellas, and E. Chouske, "Dynamollm: Designing llm inference clusters for performance and energy efficiency," in *2025 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, Mar. 2025, p. 1348–1362. [Online]. Available: http://dx.doi.org/10.1109/HPCA61900 .2025.00102

[43] Y. Wang, S. Agarwal, S. Mukherjee, X. Liu, J. Gao, A. Hassan, and J. Gao, "Adamix: Mixture-of-adaptations for parameter-efficient model tuning," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 5744–5760.

[44] Z. Wang, J. Liang, R. He, Z. Wang, and T. Tan, "Lora-pro: Are low-rank adapters properly optimized?" 2025. [Online]. Available: https://arxiv.org/abs/2407.18242

[45] B. Wu, R. Zhu, Z. Zhang, P. Sun, X. Liu, and X. Jin, "dLoRA: Dynamically orchestrating requests and adapters for LoRA LLM serving," in *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. Santa Clara, CA: USENIX

Association, Jul. 2024, pp. 911–927. [Online]. Available: https://www.usenix.org/conference/osdi24/presentation/wu-bingyang

[46] H. Wu and K. Tu, "Layer-condensed KV cache for efficient inference of large language models," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 11175–11188. [Online]. Available: https://aclanthology.org/2024.acl-long.602/

[47] L. Xu, H. Xie, S.-Z. J. Qin, X. Tao, and F. L. Wang, "Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment," 2023. [Online]. Available: https://arxiv.org/abs/2312.12148

[48] Z. Ye, L. Chen, R. Lai, W. Lin, Y. Zhang, S. Wang, T. Chen, B. Kasikci, V. Grover, A. Krishnamurthy *et al.*, "Flashinfer: Efficient and customizable attention engine for llm inference serving," *arXiv preprint arXiv:2501.01005*, 2025.

[49] G.-I. Yu, J. S. Jeong, G.-W. Kim, S. Kim, and B.-G. Chun, "Orca: A distributed serving system for Transformer-Based generative models," in *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. Carlsbad, CA: USENIX Association, Jul. 2022, pp. 521–538. [Online]. Available: https://www.usenix.org/conference/osdi22/presentation/yu

[50] H. Zhang, A. Ning, R. B. Prabhakar, and D. Wentzlaff, "Llmcompass: Enabling efficient hardware design for large language model inference," in *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2024, pp. 1080–1096.

[51] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, T. Mihaylov, M. Ott, S. Shleifer, K. Shuster, D. Simig, P. S. Koura, A. Sridhar, T. Wang, and L. Zettlemoyer, "Opt: Open pre-trained transformer language models," 2022. [Online]. Available: https://arxiv.org/abs/2205.01068

[52] Z. Zheng, X. Ji, T. Fang, F. Zhou, C. Liu, and G. Peng, "Batchllm: Optimizing large batched llm inference with global prefix sharing and throughput-oriented token batching," 2025. [Online]. Available: https://arxiv.org/abs/2412.03594